

Techniques for Dimensionality Reduction

PCA and Other Matrix
Factorization Methods

Outline

- Principle Components Analysis (PCA)
 - Example (Bishop, ch 12)
 - PCA as a mixture model variant
 - With a continuous latent variable
 - Breaking down PCA
 - Optimization problem
 - Solution
 - Intuitions
- General matrix factorization
 - Application to collaborative filtering
 - Algorithms
 - Wrap-up

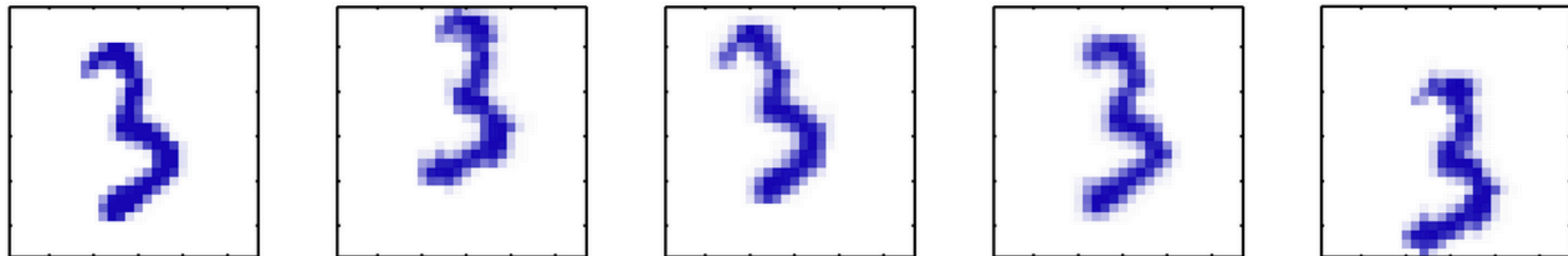
A Motivating Example

- The MNist digits problem was simplified because the digits were
 - Centered
 - In a canonical position
 - Scaled to the same size
- What if they weren't?



A Motivating Example

- Take a *single* 64*64 digit and create a dataset by repeatedly
 - Move it to a 100*100 image
 - Shift by x,y and rotate by θ
- Dataset has 10,000 features but really only needs 3

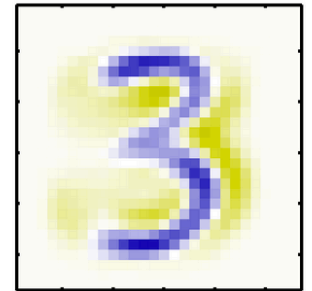
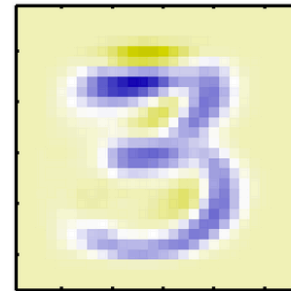
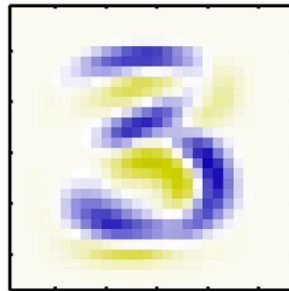
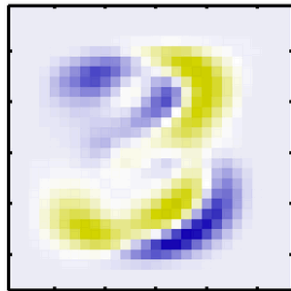
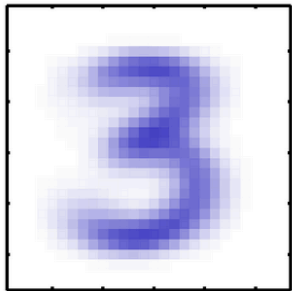


A Motivating Example

“prototype” = a vector of the same dimension as the instances

- PCA: reduces each instance to a linear combination of a few “prototypes” (blue+, green-). These are the first 5:

*A specific choice of prototypes are the **principle components***



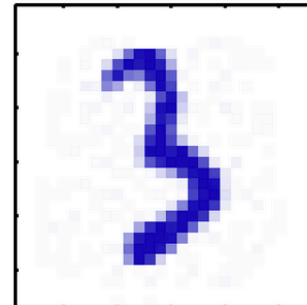
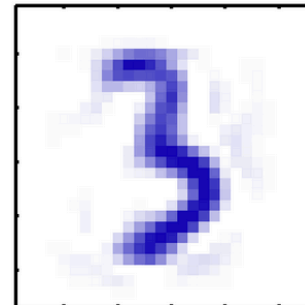
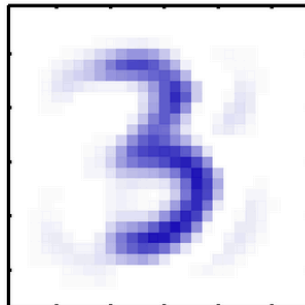
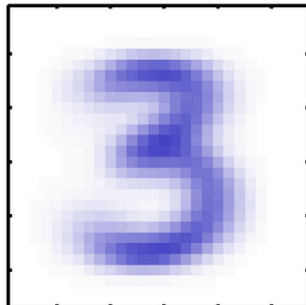
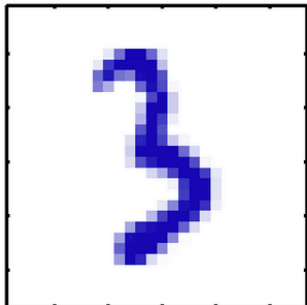
Original

$M = 1$

$M = 10$

$M = 50$

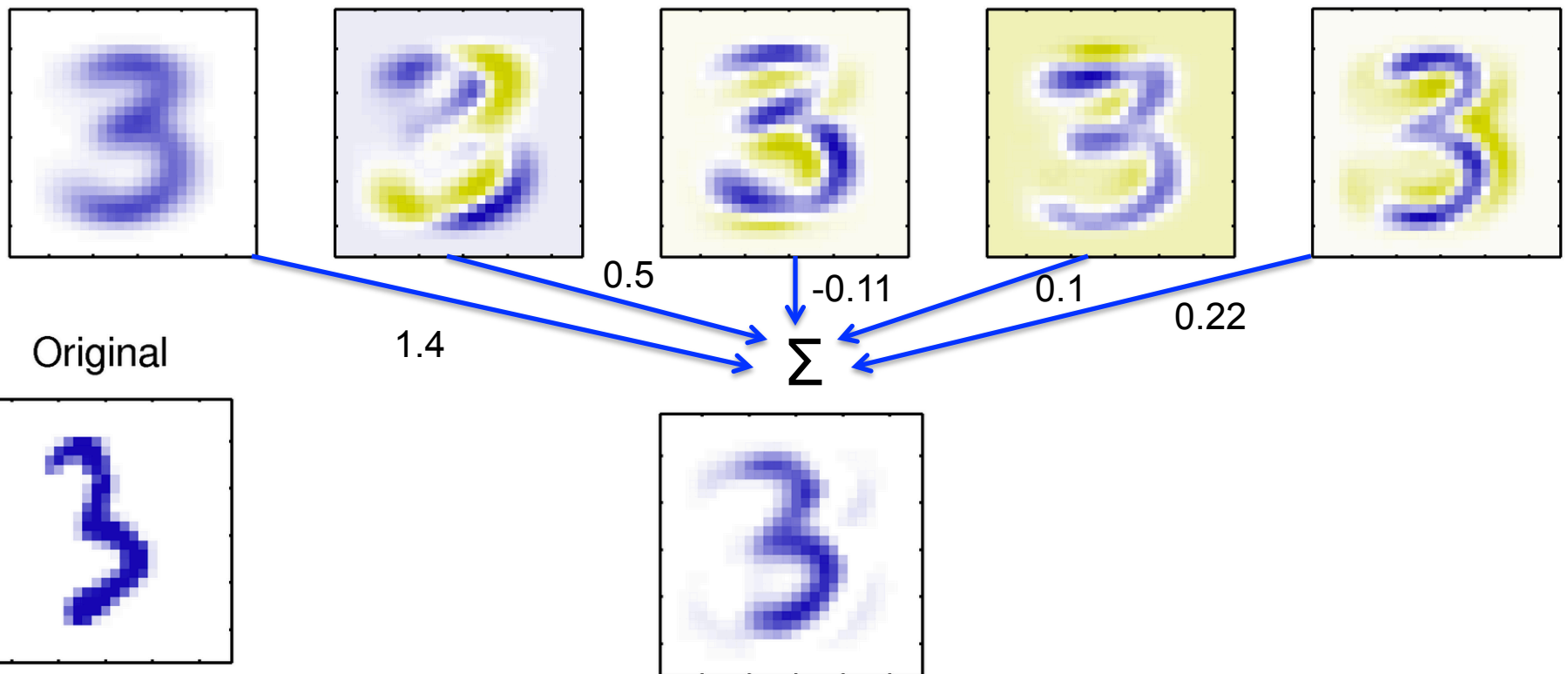
$M = 250$



A Motivating Example

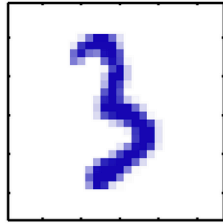
“prototype” = a vector of the same dimension as the instances

- PCA: reduces each instance to a linear combination of a few “prototypes” (blue+, green-). These are the first 5:

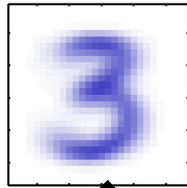


PCA as matrices

Original



PC1



2 prototypes

10,000 pixels

1000 * 10,000,00

1000 images

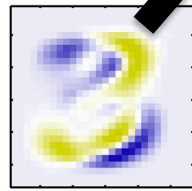
x1 y1
x2 y2
.. ..
.. ..
xn yn

\times

a1 a2 am
b1 b2 bm

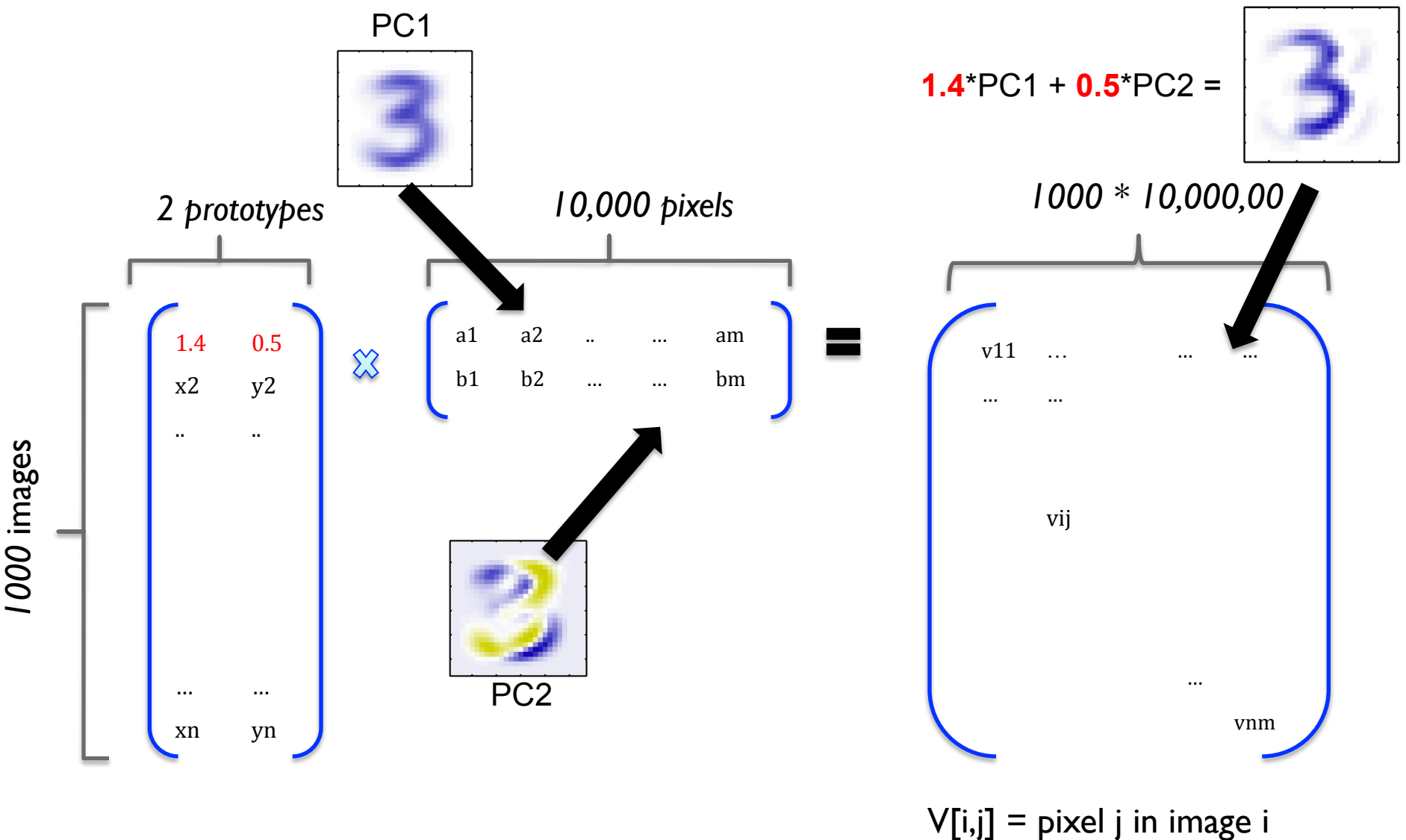
\approx

v11
... ..
vij
...
vnm

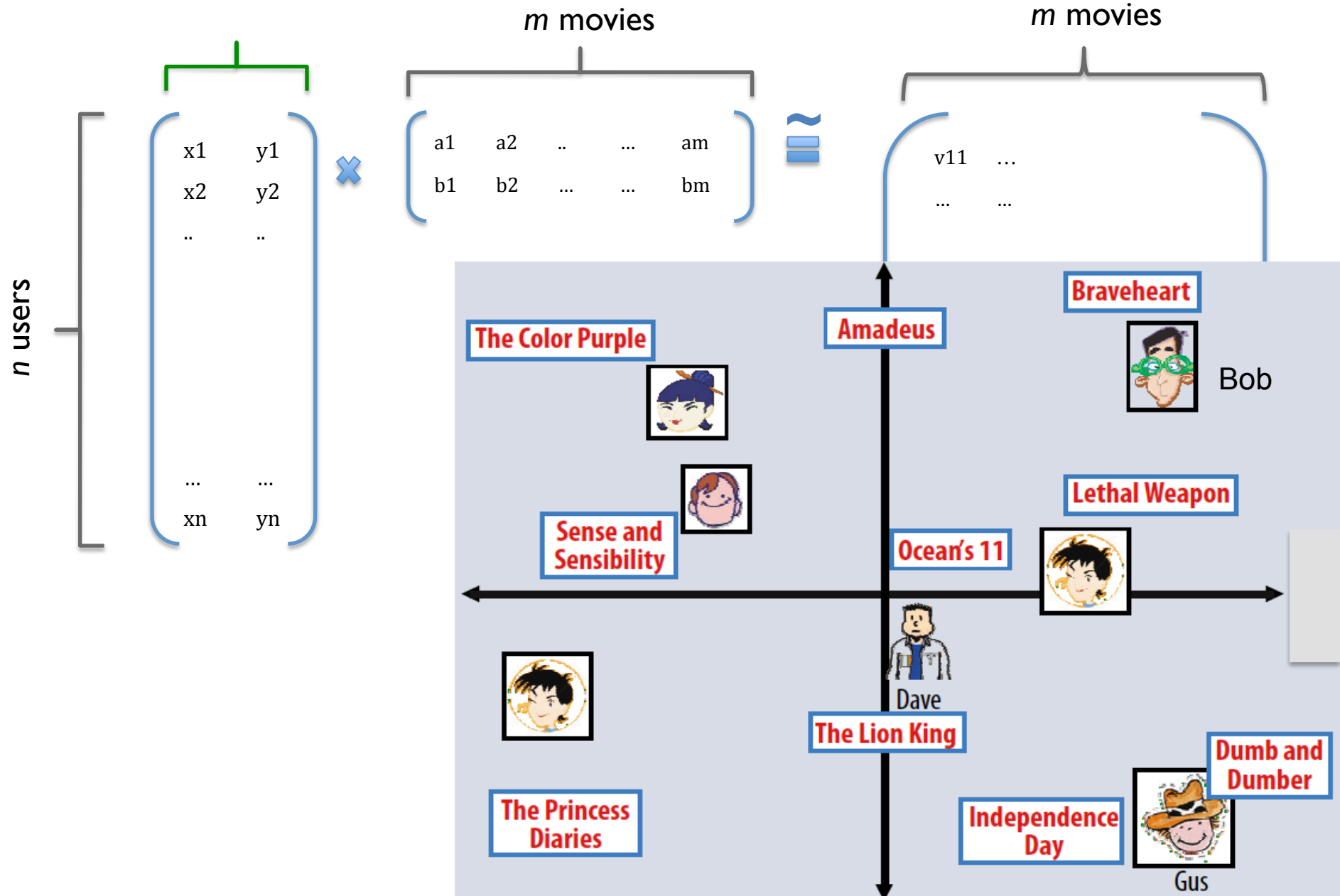


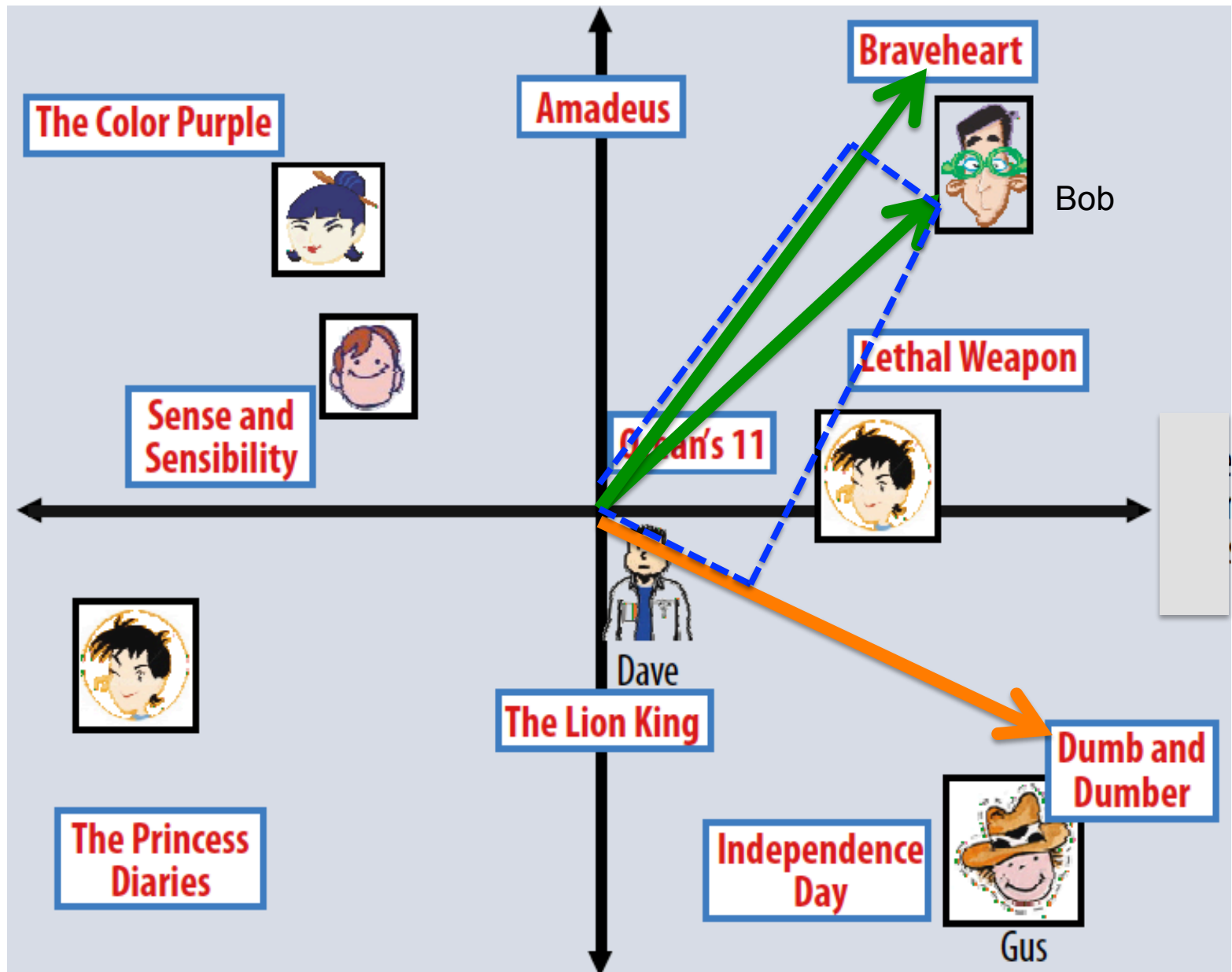
PC2

$V[i,j] = \text{pixel } j \text{ in image } i$

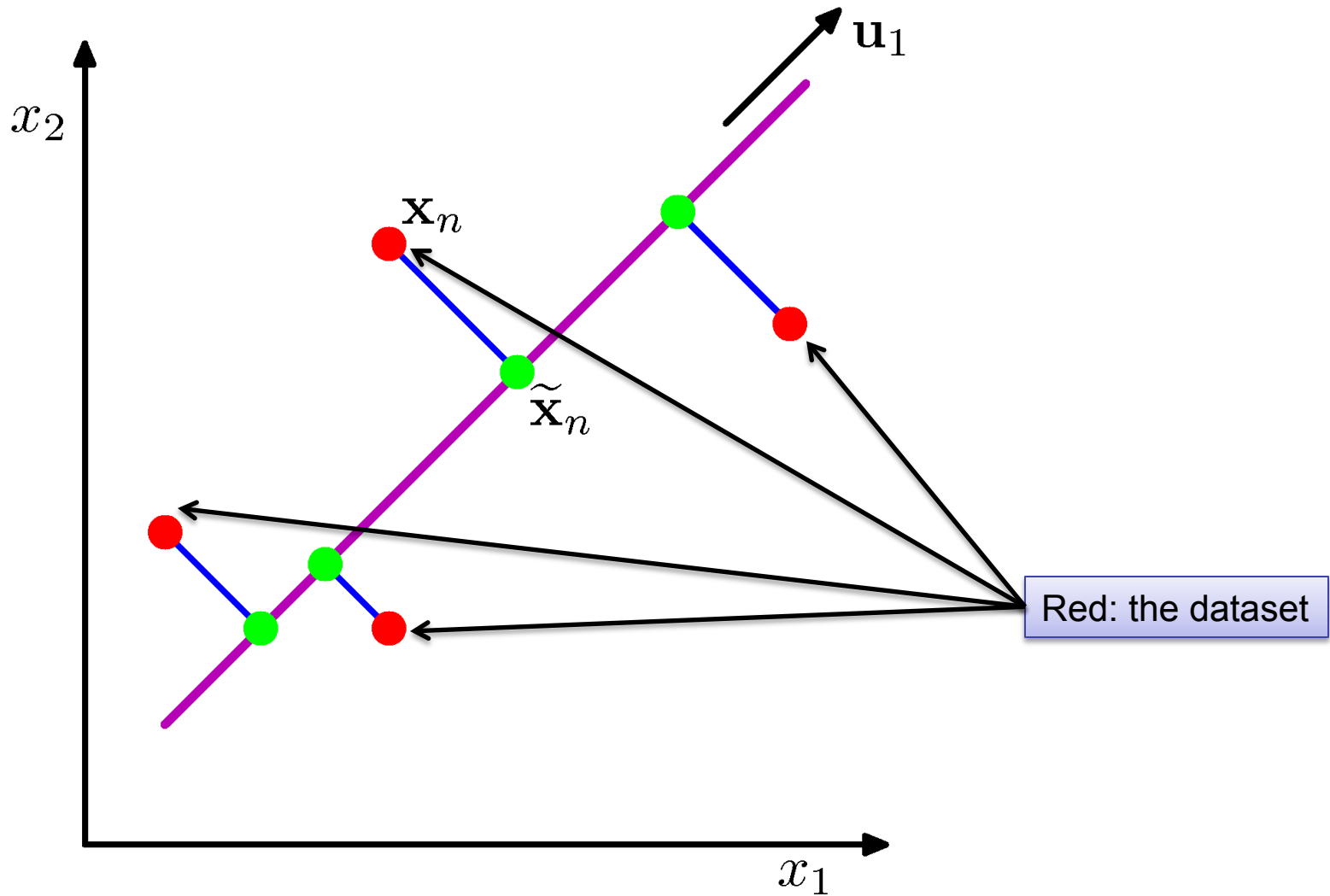


PCA for movie recommendation...

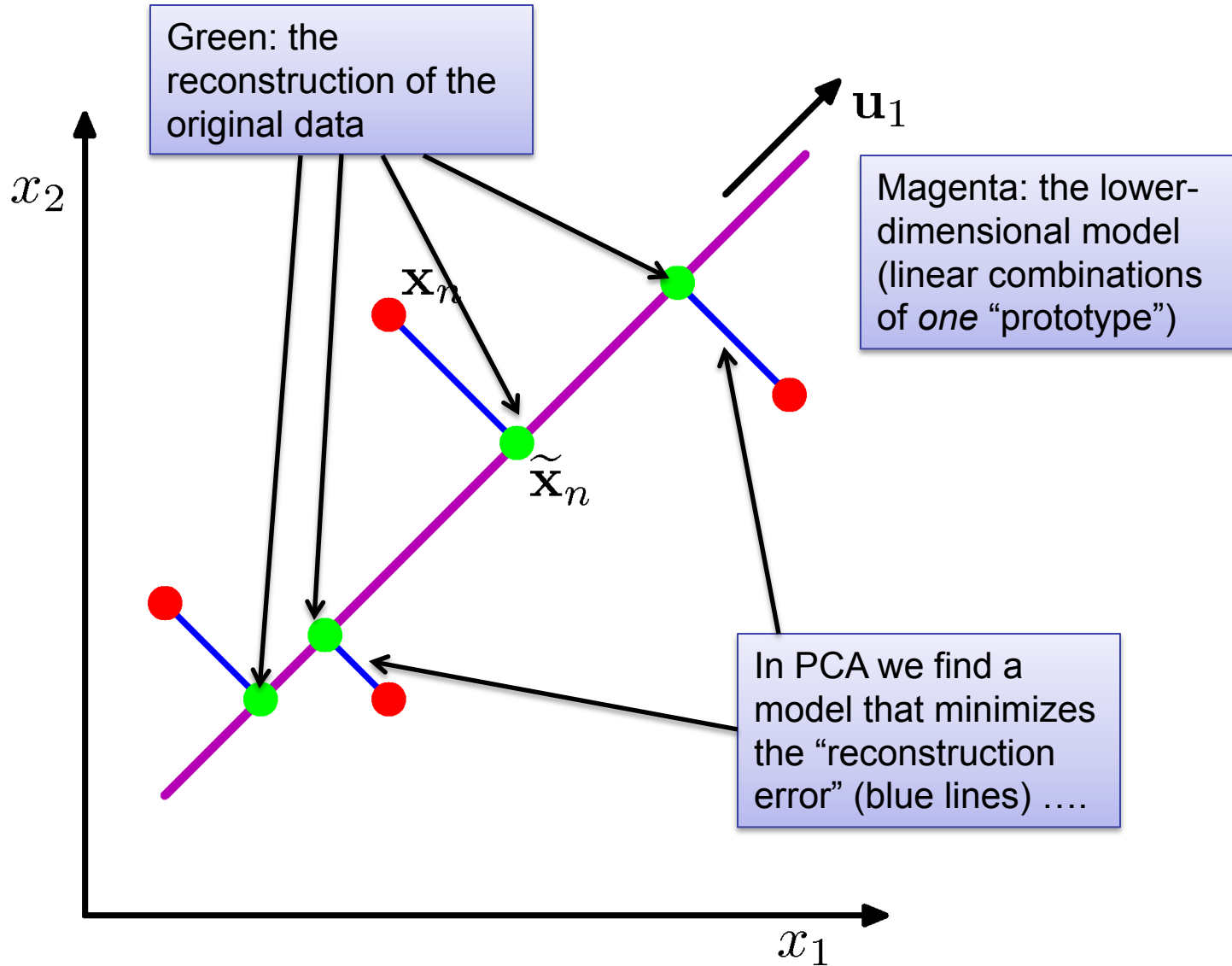




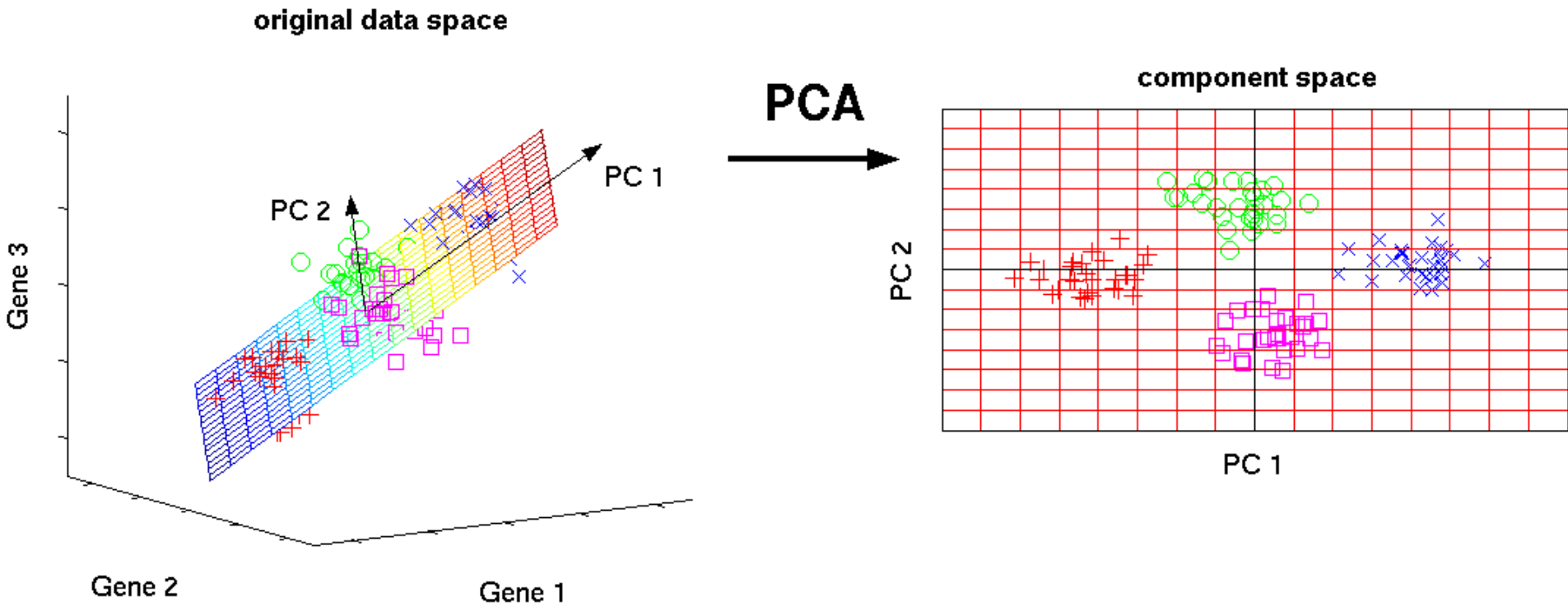
A Cartoon of PCA



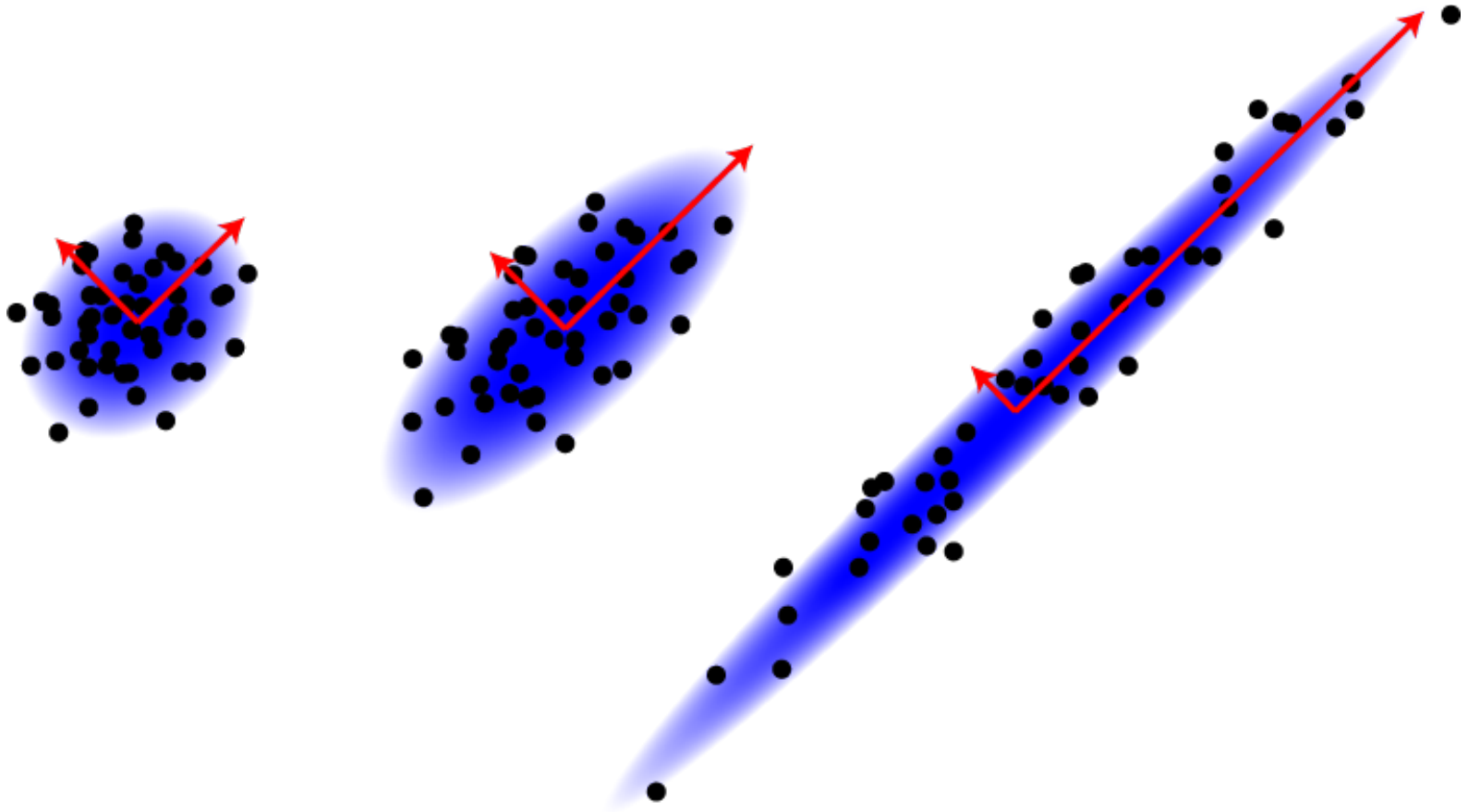
A Cartoon of PCA



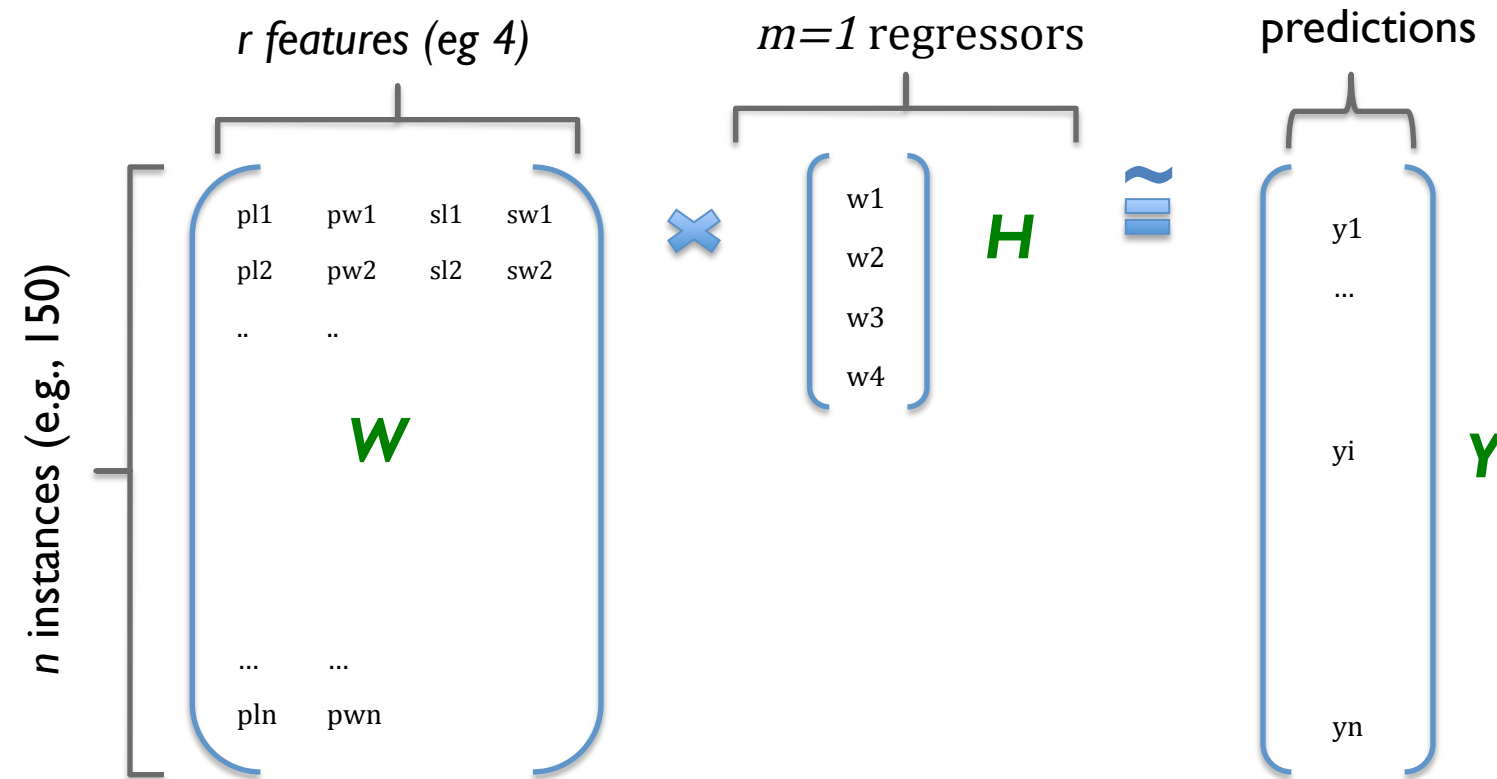
A 3D Cartoon of PCA



Some more cartoons

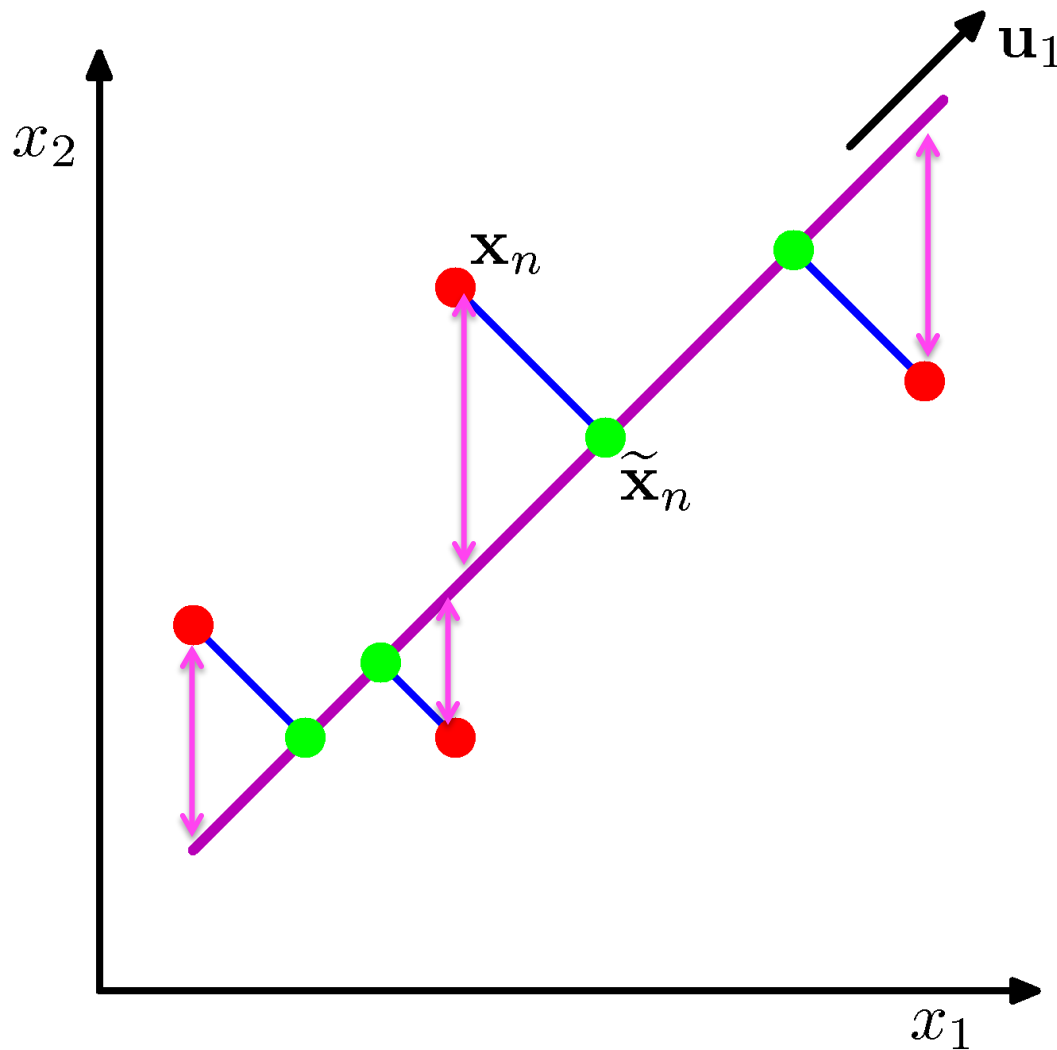


PCA vs Linear Regression



$Y[i,1]$ = instance i 's prediction

PCA vs Linear Regression



In contrast: in regression we'd minimize square error on *one* dimension (x_2) using a linear combination the *other dimensions*

PCA vs mixture of Gaussians

Mixture of Gaussians

For each point:

- Pick the index of the (latent) Gaussian $Z=k$
- Pick the the point \mathbf{x} from that the k-th Gaussian, $\mathbf{x} \sim N(\mu_k, \Sigma_k)$

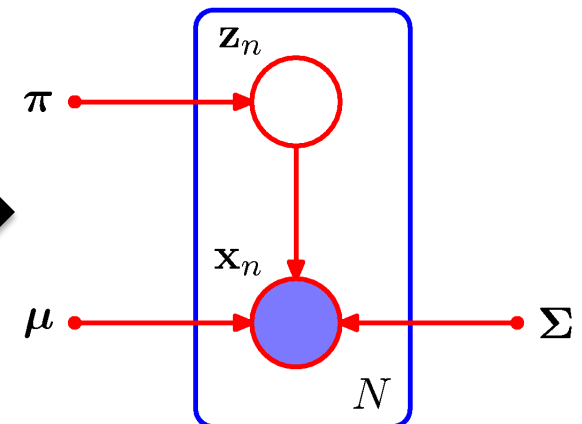
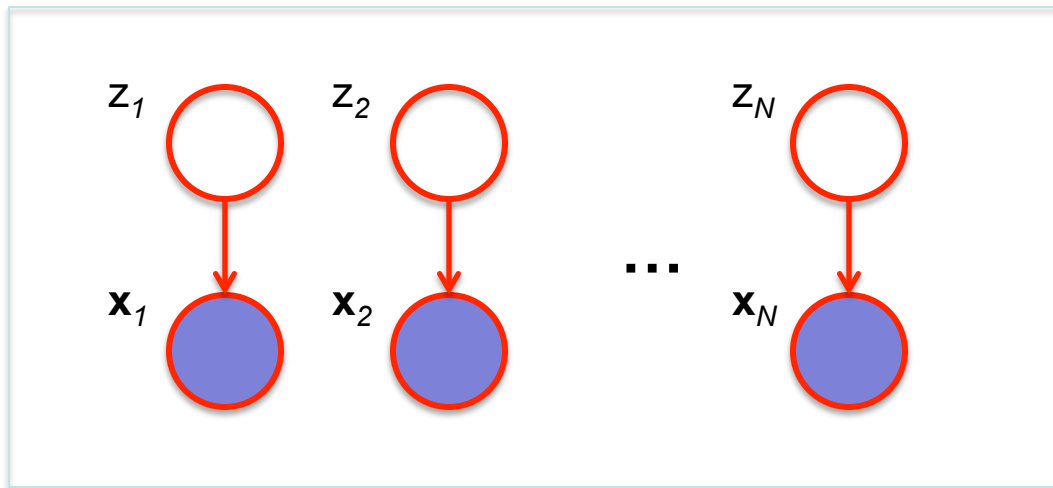
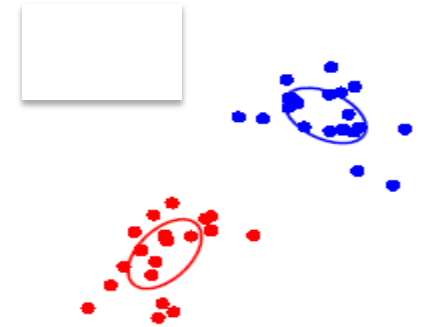
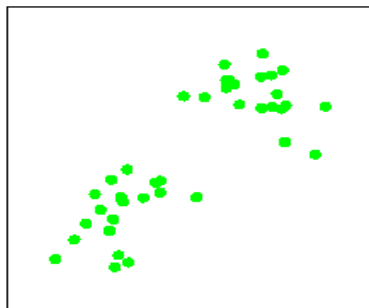
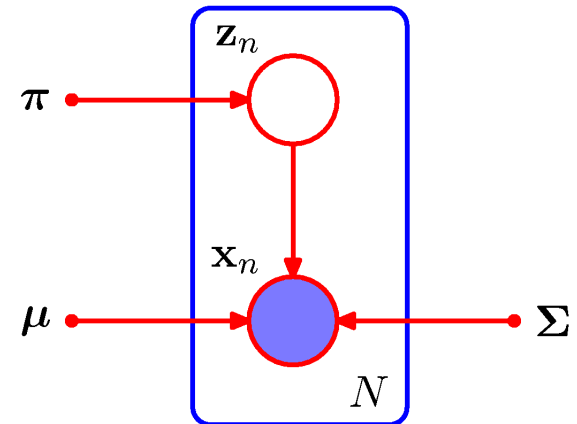


Plate notation

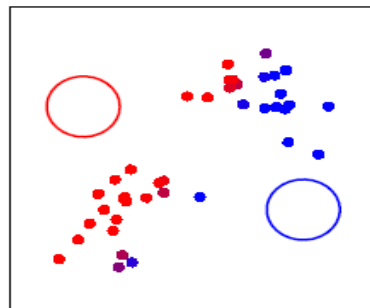
PCA vs mixture of Gaussians

Mixture of Gaussians

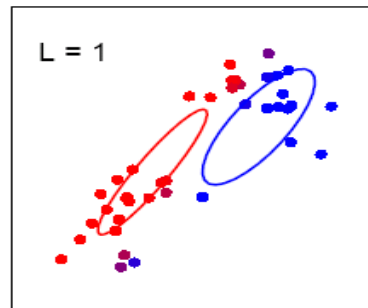
- Pick the index of the (latent) Gaussian $Z=k$
- Pick the point \mathbf{x} from that the k -th Gaussian, $\mathcal{N}(\mu_k, \Sigma_k)$



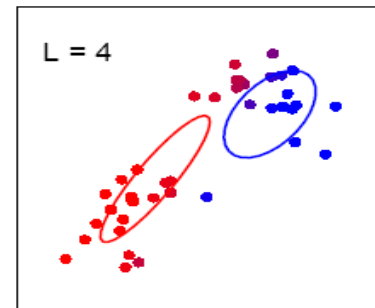
(a)



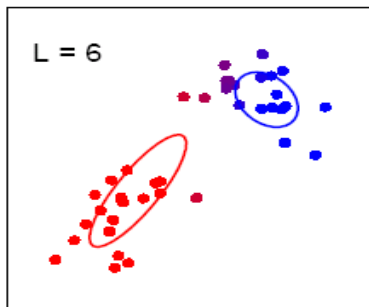
(c)



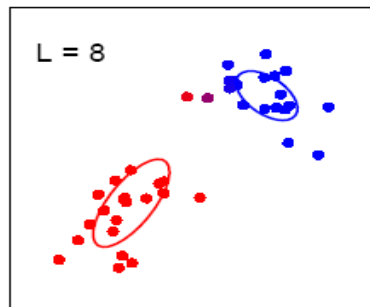
(d)



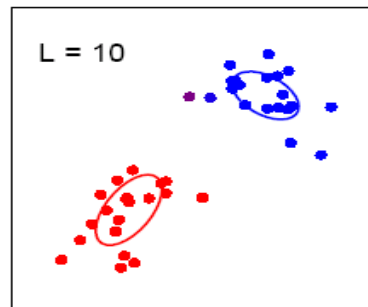
(e)



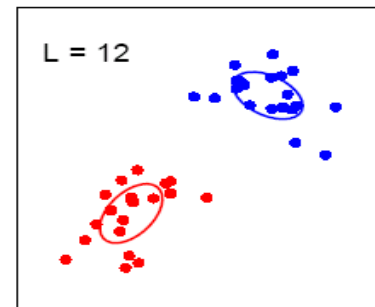
(f)



(g)



(h)

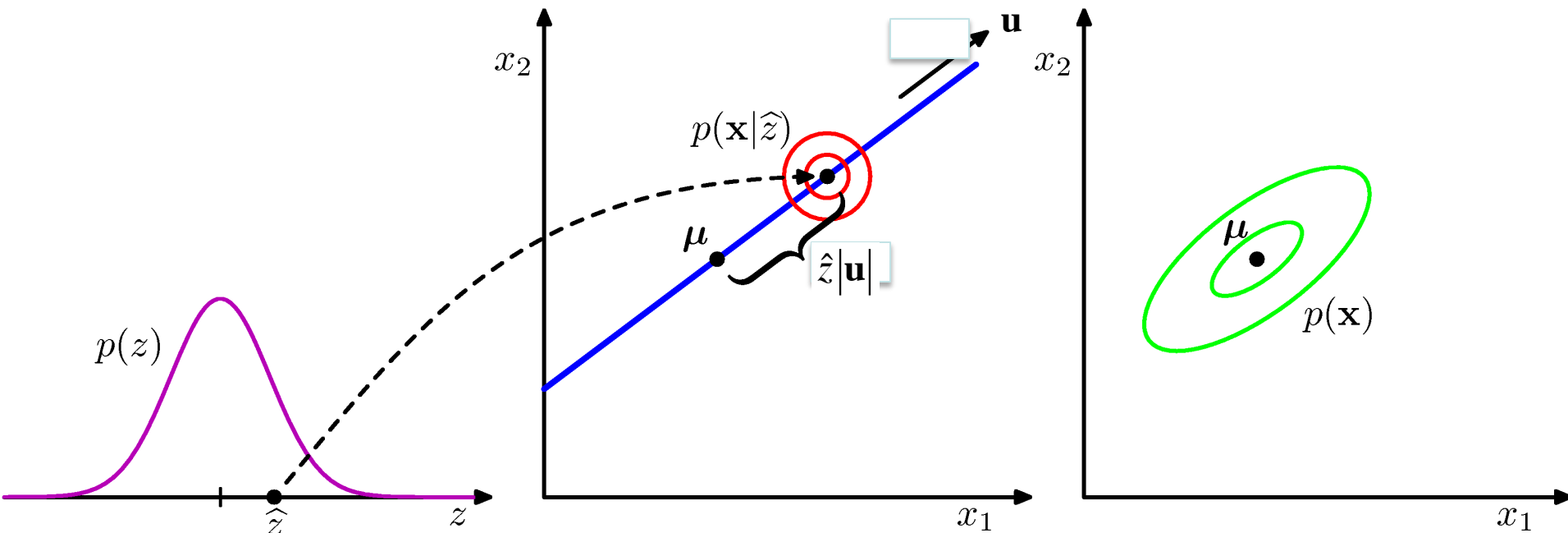
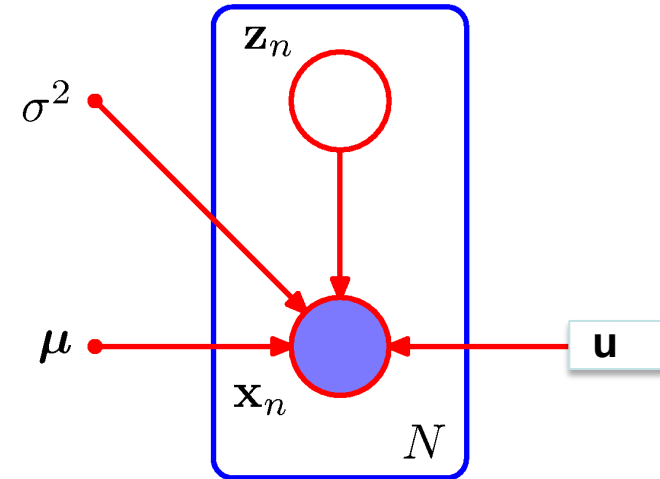


(i)

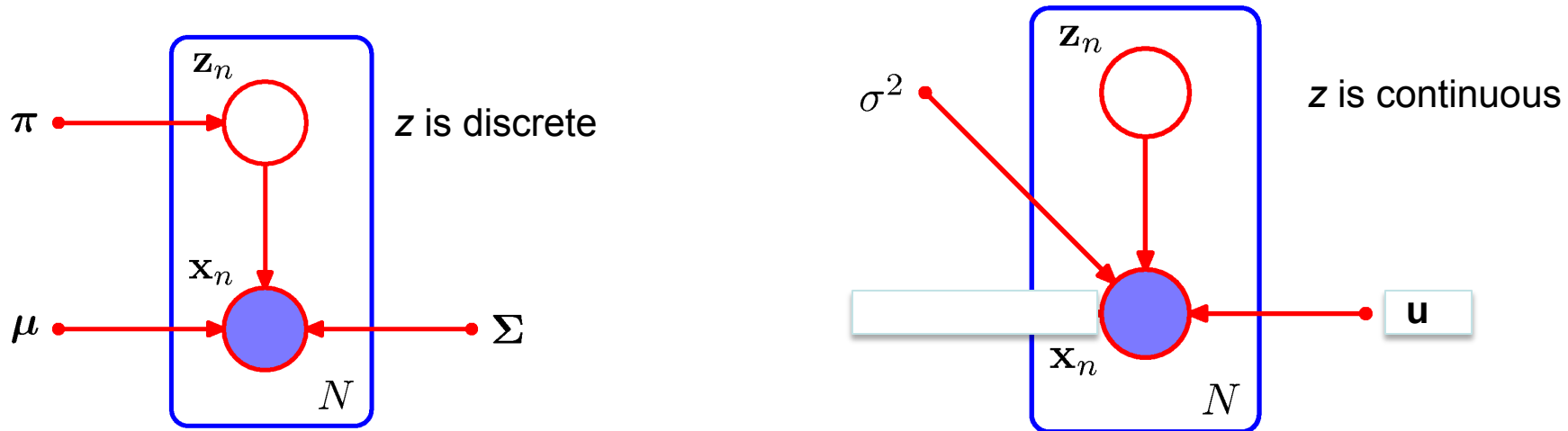
PCA vs mixture of Gaussians

PCA

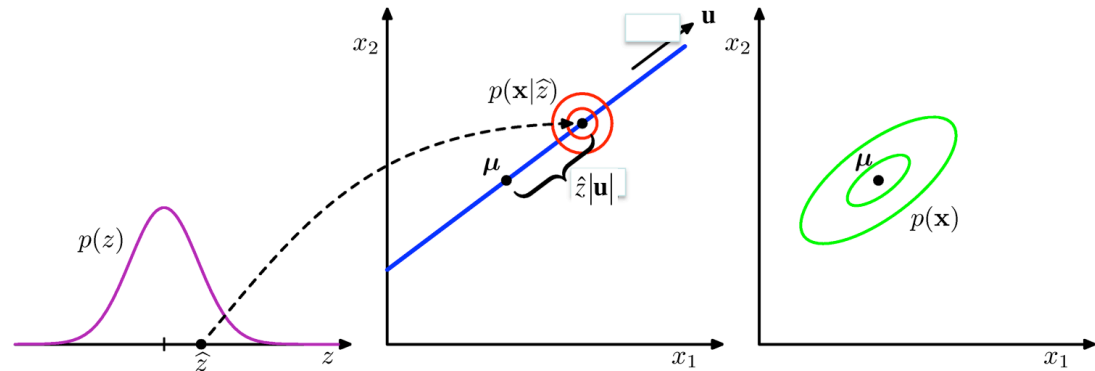
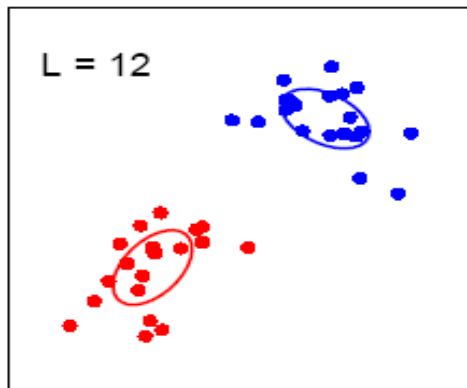
- Pick a *continuous* value z , which will be used to combine the “prototypes” \mathbf{u} in the model
- Pick the the point \mathbf{x} from a spherical Gaussian centered on $z\mathbf{u}$



PCA vs mixture of Gaussians



Comment: we can preprocess the data so that the mean is $\mathbf{0}$ to simplify the model



Finding the Principle Components

- There are different algorithms that can be used
 - EM (Roweis, NIPS 2007)
 - Can also be turned into an eigenvector computation (next)

Outline

- PCA
 - Example (Bishop, ch 12)
 - PCA as a mixture model variant
 - With a continuous latent variable
 - Breaking down PCA
 - Optimization problem
 - Solution
 - Intuition

The PCA Problem (vectors)

Start with a zero-mean dataset, where \mathbf{x}^t is a the t -th instance:

$$X = \begin{bmatrix} x_1^1 & x_2^1 & \dots & x_m^1 \\ \vdots & \ddots & & \vdots \\ x_1^n & \dots & & x_m^n \end{bmatrix} = \begin{bmatrix} \vdots \\ - \mathbf{x}^t - \\ \vdots \end{bmatrix}$$

We want to find small number of *orthogonal* prototypes $\mathbf{u}_1, \dots, \mathbf{u}_k$ and k weights z_1^t, \dots, z_k^t for each instance \mathbf{x}^t so that if we approximate \mathbf{x}^t by

$$\hat{\mathbf{x}}^t = \sum_{i=1}^k z_i^t \mathbf{u}_i$$

the approximation error will be small: we want to find \mathbf{u} 's and \mathbf{z} 's to minimize

$$J = \frac{1}{N} \sum_{t=1}^N \|\mathbf{x}^t - \hat{\mathbf{x}}^t\|^2$$

The PCA Problem (matrices)

Given a zero-mean dataset

$$X = \begin{bmatrix} x_1^1 & x_2^1 & \dots & x_m^1 \\ \vdots & \ddots & & \vdots \\ x_1^n & \dots & & x_m^n \end{bmatrix} = \begin{bmatrix} \vdots \\ - \mathbf{x}^t - \\ \vdots \end{bmatrix}$$

Find factors U and Z so that X is approximately their outer product:

$$\begin{bmatrix} \vdots \\ \dots \mathbf{z}^t \dots \\ \vdots \end{bmatrix} \begin{bmatrix} - & - & \mathbf{u}_k & - & - \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots \\ - \hat{\mathbf{x}}^t - \\ \vdots \end{bmatrix} = \hat{X} \quad \hat{\mathbf{x}}^t = \sum_{i=1}^k z_k^t \mathbf{u}_k$$

Specifically minimizing the square of the reconstruction error

$$J = \frac{1}{N} \sum_{t=1}^N \|\mathbf{x}^t - \hat{\mathbf{x}}^t\|^2$$

under the constraint that the rows of U are *orthogonal*.

A PCA Algorithm

Start with a zero-mean dataset, where

- \mathbf{x}^t is a the t -th instance
- \mathbf{f}_i is a column of feature values for the i -th feature.
- Compute the sample covariance matrix

i.e.,

$$C_X = X^T X$$
$$C_X(i, j) = \sum_t f_i^t f_j^t$$
$$X = \begin{bmatrix} x_1^1 & x_2^1 & \dots & x_m^1 \\ \vdots & \ddots & & \vdots \\ x_1^n & \dots & & x_m^n \end{bmatrix} = \begin{bmatrix} \vdots \\ - \mathbf{x}^t - \\ \vdots \end{bmatrix}$$
$$X = \begin{bmatrix} x_1^1 & x_2^1 & \dots & x_m^1 \\ \vdots & \ddots & & \vdots \\ x_1^n & \dots & & x_m^n \end{bmatrix} = \begin{bmatrix} \dots & | & \dots \\ \dots & \mathbf{f}_i & \dots \\ \dots & | & \dots \end{bmatrix}$$

- Find the largest k eigenvectors of C_X . These are the prototypes, U .
- Now find Z given X and U .

PCA Algorithm: Intuitions

Start with a zero-mean dataset, where

- \mathbf{x}^t is a the t -th instance
- \mathbf{f}_i is a column of feature values for the i -th feature.
- Compute the sample covariance matrix

$$C_X = X^T X$$

Some intuitions:

1. Suppose you wanted to predict feature i from feature j . Your best guess would be

$$f_i \text{ is predicted as } C_X(i, j) \cdot f_j$$

2. If you wanted to predict feature i from *all other* feature's j , a plausible guess is

$$f_i \text{ is predicted as } \frac{1}{n} \sum_{j \neq i} C_X(i, j) \cdot f_j$$

3. Any *eigenvector*, \mathbf{e} , of C_X leads to an *internally consistent** set of predictions

* up to a multiplier

$$\exists \lambda : \lambda \mathbf{e} = C_X \mathbf{e} \quad \longrightarrow \quad \forall i, \lambda e_i = \frac{1}{n} \sum_j C_X(i, j) e_j$$

PCA: Eigenfaces

Turk and Pentland, 1991



PCA: Eigenfaces

Turk and Pentland, 1991

Average face



Six eigenfaces (PC's)

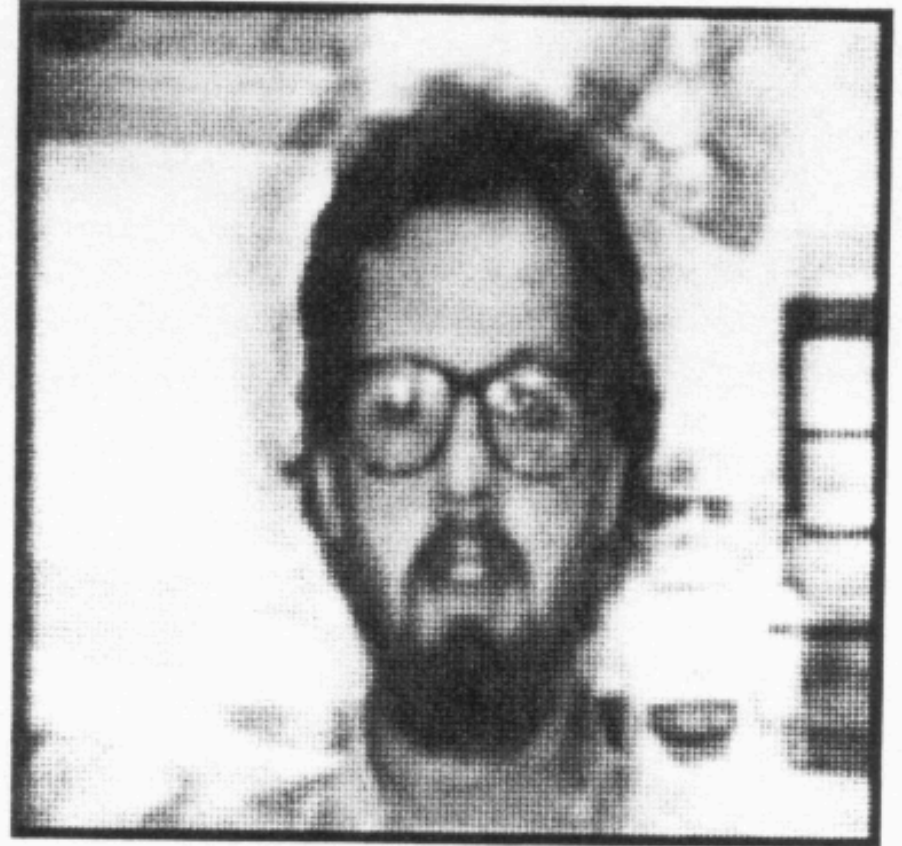
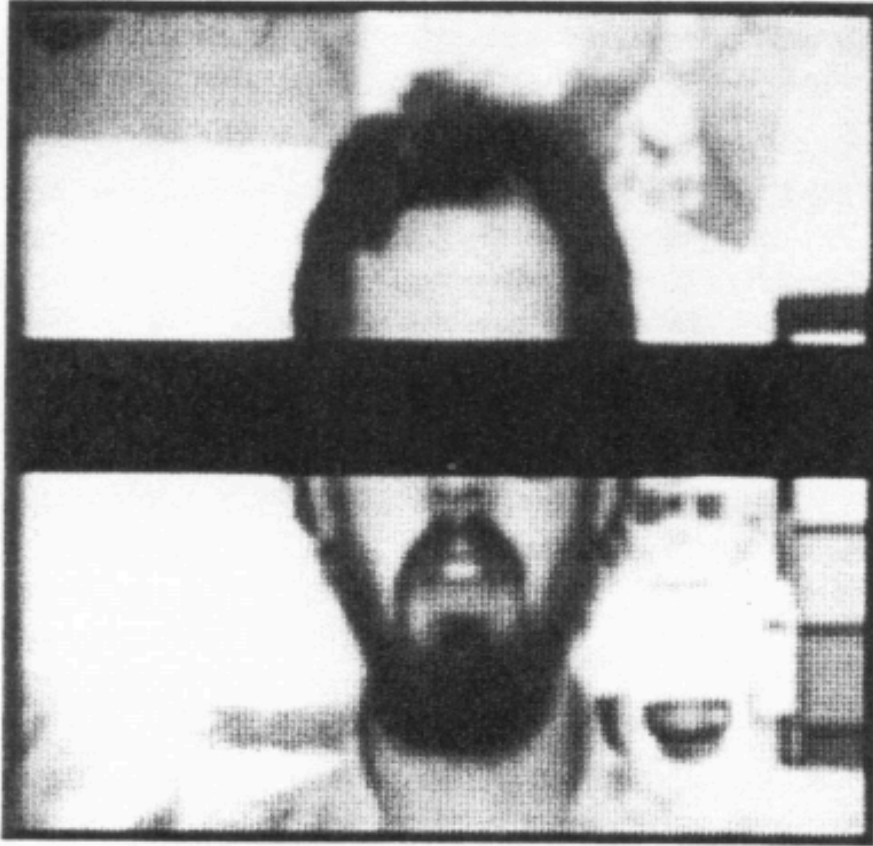


PCA: Eigenfaces

Turk and Pentland, 1991



PCA: Eigenfaces

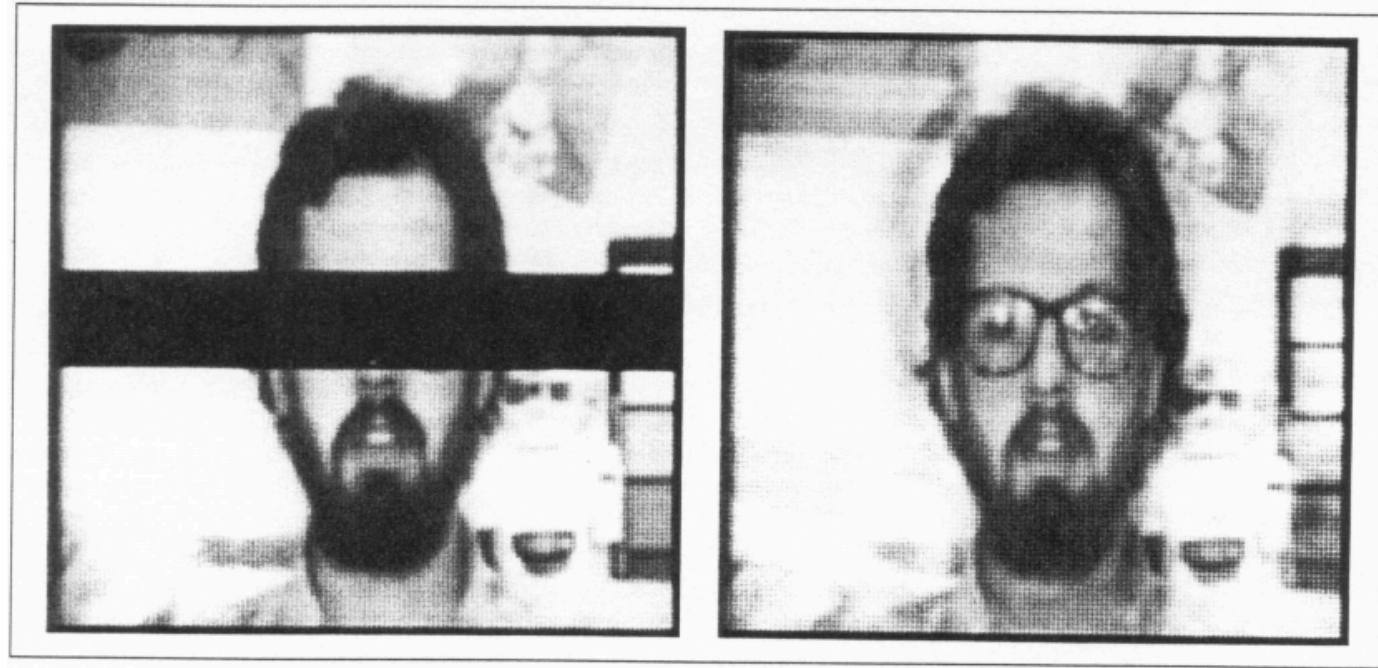


PCA: Eigenfaces

How is this done?

Simplest approach:

- Add the image with missing values to the data matrix
- Minimize reconstruction error over the non-missing values



$$\begin{bmatrix} \vdots \\ \dots & \mathbf{z}^t & \dots \\ \vdots \end{bmatrix} \begin{bmatrix} \vdots \\ - & - & \mathbf{u}_k & - & - \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots \\ - & \hat{\mathbf{x}}^t & - \\ - & ? & - \\ \vdots \end{bmatrix} = \hat{X}$$

Matrix completion for image denoising



Outline

- Principle Components Analysis (PCA)
- Other types of/applications of matrix factorization
 - Collaborative filtering/recommendation
 - Matrix factorization for CF using gradient descent

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?



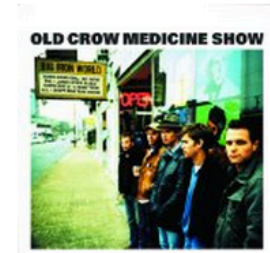
New Release
Let's Be Still
The Head And The Heart
★★★★☆ (21)
\$9.49
Why recommended?



Leaving Eden
Carolina Chocolate Drops
★★★★☆ (66)
\$10.49
Why recommended?



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

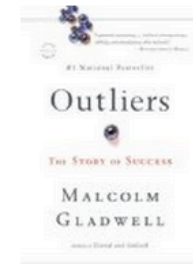
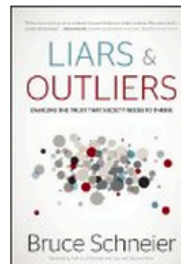
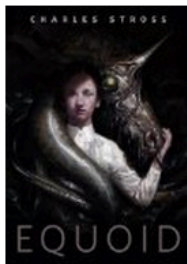


Big Iron World
Old Crow Medicine Show
★★★★☆ (39)
\$9.49
Why recommended?

▶ See all recommendations in MP3 Albums

Kindle eBooks

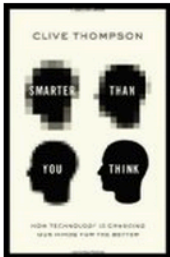
Page 1 of 20



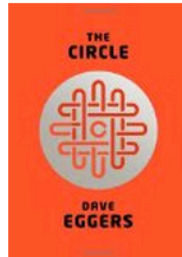
What is collaborative filtering?

Books

Page 1 of 20



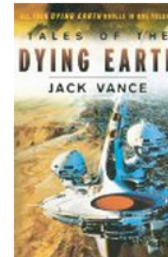
New Release
Smarter Than You ...
▶ Clive Thompson
★★★★☆ (26)
\$27.95 **\$20.82**
Why recommended?



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 Gutr 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**

- [Instant videos you've watched](#)
- [Items you've marked "I own it"](#)
- [Items you've rated](#)
- [Items you've liked](#)
- [Items you've marked "Not interested"](#)
- [Items you've marked as gifts](#)

EDIT YOUR PREFERENCES

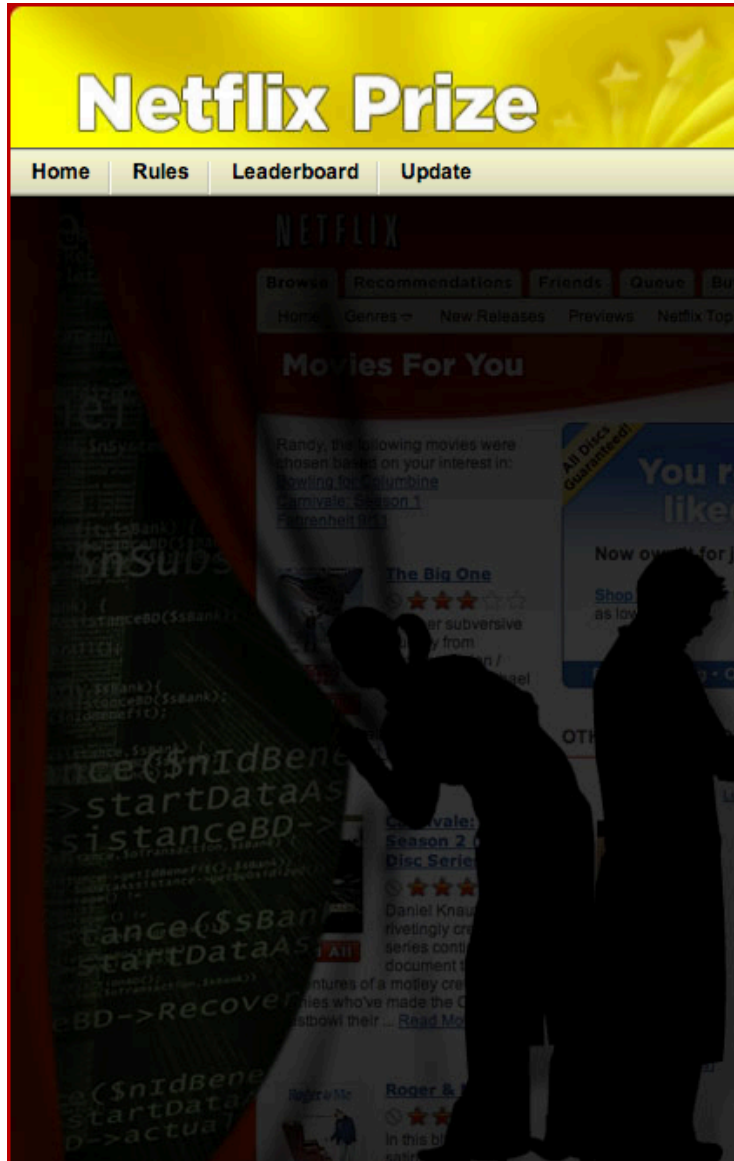
Show Amazon book recommendations as Kindle editions when possible.

Items you've purchased

- | | | Your Rating: |
|----|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------|
| 1. |  <p>Love Is Strange (A Paranormal Romance)
by Bruce Sterling
Your tags:
<input type="text"/> <input type="button" value="Add"/> (What's this?)
Click to Add: paranormal romance, nerd, futurist, science fiction romance, science fiction, technology, scifi, literature</p> | <input checked="" type="checkbox"/> ☆☆☆☆☆
<input type="checkbox"/> This was a gift
<input type="checkbox"/> Don't use for recommendations |
| 2. |  <p>Mad Magazine #1
by Harvey Kurtzman
Your tags:
<input type="text"/> <input type="button" value="Add"/> (What's this?)
Click to Add: harvey kurtzman, dc</p> | <input checked="" type="checkbox"/> ☆☆☆☆☆
<input type="checkbox"/> This was a gift
<input type="checkbox"/> Don't use for recommendations |
| 3. |  <p>Ahoy!
Punch Brothers Format: MP3 Music
Your tags:
<input type="text"/> <input type="button" value="Add"/> (What's this?)
Click to Add: bluegrass, music, punch brothers, singer-songwriters</p> | <input checked="" type="checkbox"/> ☆☆☆☆☆
<input type="checkbox"/> This was a gift
<input type="checkbox"/> Don't use for recommendations |

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

What is collabo



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

What is collaborative filtering?

A Framework for Optimizing Paper Matching

Laurent Charlin

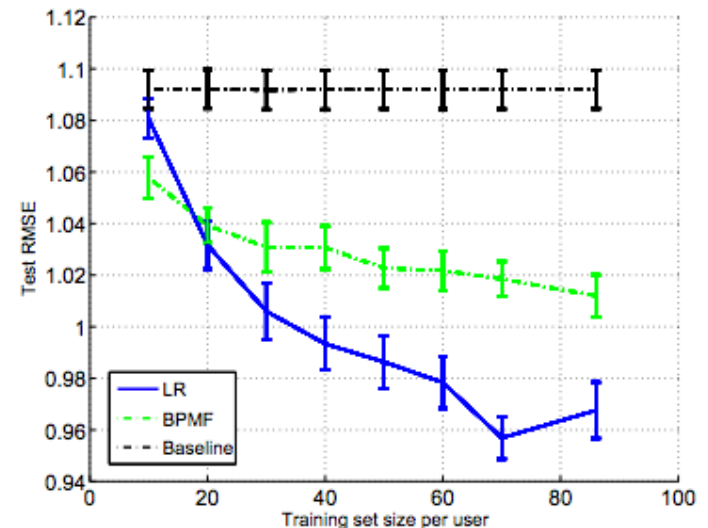
Department of Computer Science
University of Toronto
Toronto, ON M5S 3H5
laurent@cs.toronto.edu

Richard Zemel

Department of Computer Science
University of Toronto
Toronto, ON M5S 3H5
zemel@cs.toronto.edu

Craig Boutilier

Department of Computer Science
University of Toronto
Toronto, ON M5S 3H5
cebly@cs.toronto.edu



Other examples of social filtering....

[Web](#) [Images](#) [Maps](#) [Shopping](#) [News](#) [More ▾](#) [Search tools](#)


About 24,700,000 results (0.24 seconds)

[Carnegie Mellon University | CMU](#)

www.cmu.edu/ ▾

Carnegie Mellon University (CMU) is a global research university recognized for world-class arts and technology programs, collaboration across disciplines and ...

4.7 ★★★★★ 105 Google reviews · [Write a review](#)

 5000 Forbes Ave Pittsburgh, PA 15213
(412) 268-2000

[Careers](#)

Connect to Careers@CarnegieMellon to view and ...

[Academics](#)

Schools & Colleges - Graduate Admission - College of Fine Arts

[Carnegie Mellon Athletics](#)

Fitness & Recreation - Staff Directory - Men's Track and Field

[Prospective Students](#)

Admissions - Graduate Education - International Students - ...

[Contact Us](#)

Skip navigation and jump directly to page content. Contact Us ...

[Directory](#)

Public information for all university computing accounts.

[News for cmu](#)

[CMU president Subra Suresh elected to Institute of Medicine](#)

[Pittsburgh Post Gazette](#) - 1 day ago

Mr. Suresh, who became president of **CMU** in July, was elected to the Institute of Medicine today in recognition of his research into cell ...

[CMU ex-trustee's money laundering trial begins](#)

[Pittsburgh Post Gazette](#) - 1 day ago

[CMU student on treatment and recovery from ...](#)

[Carnegie Mellon University - Wikipedia, the free encyclopedia](#)

en.wikipedia.org/wiki/Carnegie_Mellon_University ▾

Carnegie Mellon has seven colleges and independent schools: the Carnegie Institute of Technology (engineering), College of Fine Arts, Dietrich College of ...

[Central Michigan University: Home](#)

www.cmich.edu/ ▾

Students are offered educational experiences in the arts, humanities, and natural and social sciences, in addition to educational depth in at least one academic ...

[Complete Music Update](#)

www.completemusicupdate.com/ ▾

An UnLimited Media Website, **Complete Music Update** | ThisWeek London | ThreeWeeks Edinburgh ... NEXT **CMU INSIGHTS COURSE** | NOW BOOKING

People also search for



University of Pittsburgh
Pittsburgh



Cornell University
Ithaca



Massach...
Institute of
Technology
Cambridge



Carnegie
Mellon
School of...
Pittsburgh



Stanford
University
Stanford

Other examples of social filtering....

October 26, 2013

COMBINED PRINT & E-BOOK FICTION

1. **STORM FRONT**, by John Sandford
2. **DOING HARD TIME**, by Stuart Woods
3. **DOCTOR SLEEP**, by Stephen King
4. **THE HUSBAND'S SECRET**, by Liane Moriarty
5. **THE LONGEST RIDE**, by Nicholas Sparks

[Complete List »](#)

COMBINED PRINT & E-BOOK NONFICTION

1. **KILLING JESUS**, by Bill O'Reilly and Martin Dugard
2. **MY STORY**, by Elizabeth Smart with Chris Stewart
3. **DAVID AND GOLIATH**, by Malcolm Gladwell
4. **I AM MALALA**, by Malala Yousafzai with Christina Lamb
5. **THE REASON I JUMP**, by Naoki Higashida

[Complete List »](#)

HARDCOVER FICTION

1. **STORM FRONT**, by John Sandford
2. **DOCTOR SLEEP**, by Stephen King
3. **THE LONGEST RIDE**, by Nicholas Sparks
4. **GONE**, by James Patterson and Michael Ledwidge
5. **DOG SONGS**, by Mary Oliver

[Complete List »](#)

HARDCOVER NONFICTION

1. **KILLING JESUS**, by Bill O'Reilly and Martin Dugard
2. **DAVID AND GOLIATH**, by Malcolm Gladwell
3. **THE REASON I JUMP**, by Naoki Higashida
4. **I AM MALALA**, by Malala Yousafzai with Christina Lamb
5. **MY STORY**, by Elizabeth Smart with Chris Stewart

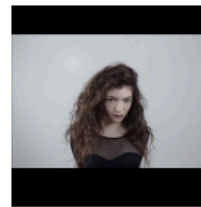
f 80.916 k t 19.815 k g+ 1.1 k

The Hot 100

The week's most popular current songs across ...

1 - 10 11 - 20 21 - 30 31 - 40 ... 91 - 100 Next chart

[Play entire chart](#)



Royals

Lorde
Love Club

▲ 1 📧 1 🕒 15

[Listen](#)

[Watch](#)

[Buy](#)

[Send Ringtone](#)

f 0 t 0 g+ 0 🗨 0

Lorde's 'Royals' Rules Hot 100 For Third Week

Lorde's "Royals" commands the Billboard Hot 100 for a third week, Ylvis' "The Fox" darts to the Digital Songs top 10 and Justin Bieber and Eminem soar into the Hot 100's top 20.

As we do every Wednesday, let's break down the Hot 100's top 10 and more.



Wrecking Ball

Miley Cyrus
Banaerz

[Listen](#)

[Watch](#)

[Buy](#)

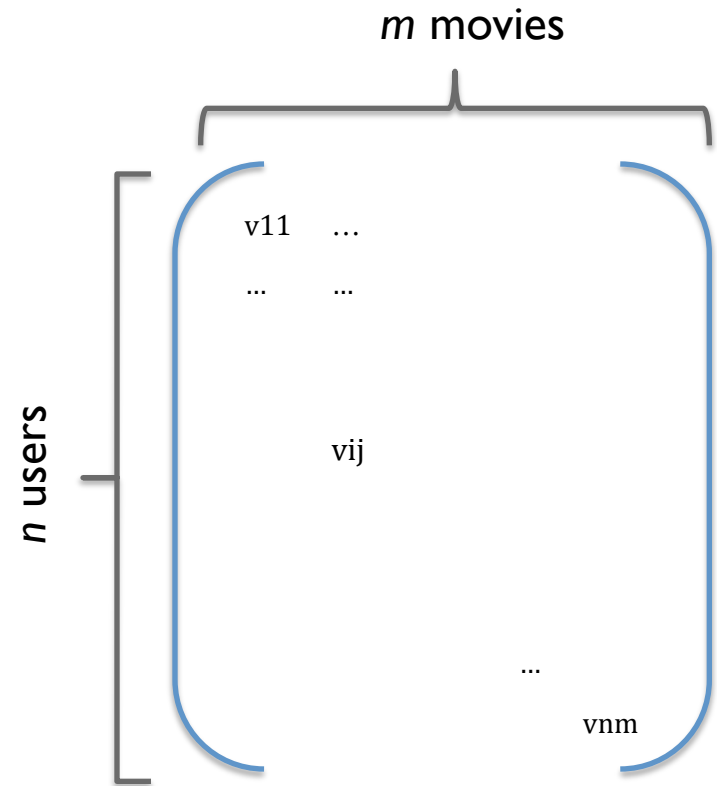
Everyday Examples of Collaborative Filtering...

- Bestseller lists
- Top 40 music lists
- The “recent returns” shelf at the library
- Unmarked but well-used paths thru the woods
- The printer room at work
- “Read any good books lately?”
-
- **Common insight:** personal tastes are *correlated*:
 - If Alice and Bob both like X and Alice likes Y then Bob is more likely to like Y
 - especially (perhaps) if Bob knows Alice

Outline

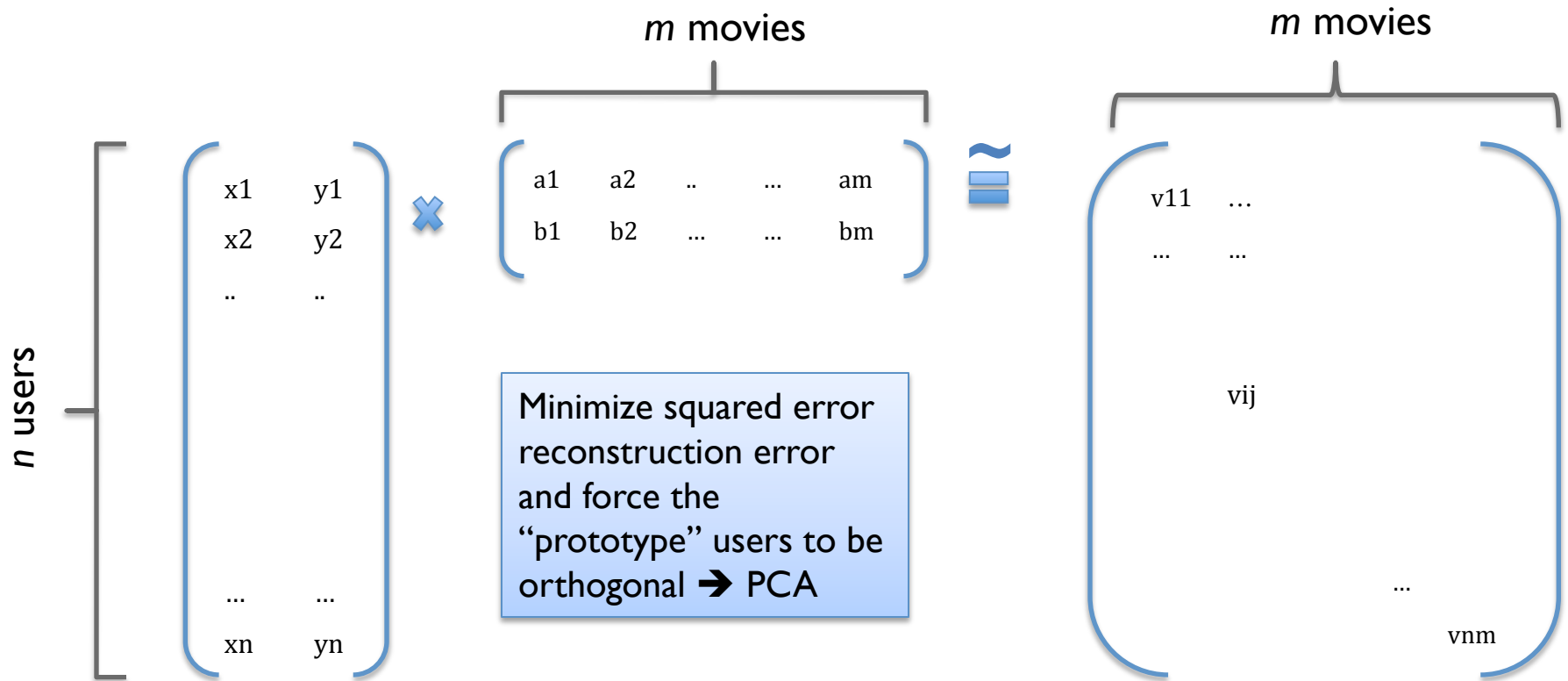
- Principle Components Analysis (PCA)
- Other types of/applications of matrix factorization
 - Collaborative filtering/recommendation
 - Algorithms:
 - K-NN type methods
 - Classification-base methods
 - ...
 - Matrix factorization

Recovering latent factors in a matrix



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

Recovering latent factors in a matrix



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

Large-Scale Matrix Factorization with Distributed Stochastic Gradient Descent

Rainer Gemulla



talk pilfered from →

Peter J. Haas



Yannis Sismanis



Erik Nijkamp



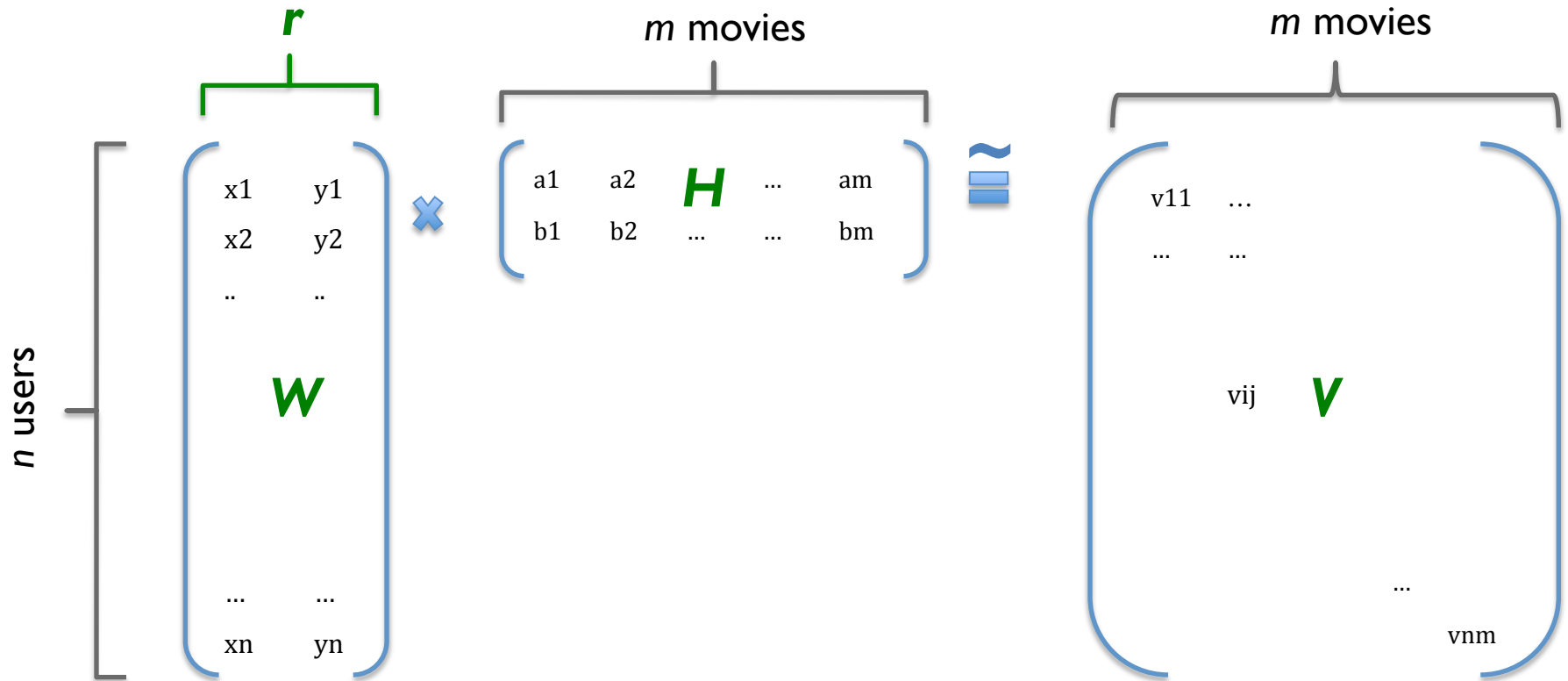
Collaborative Filtering

- ▶ 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

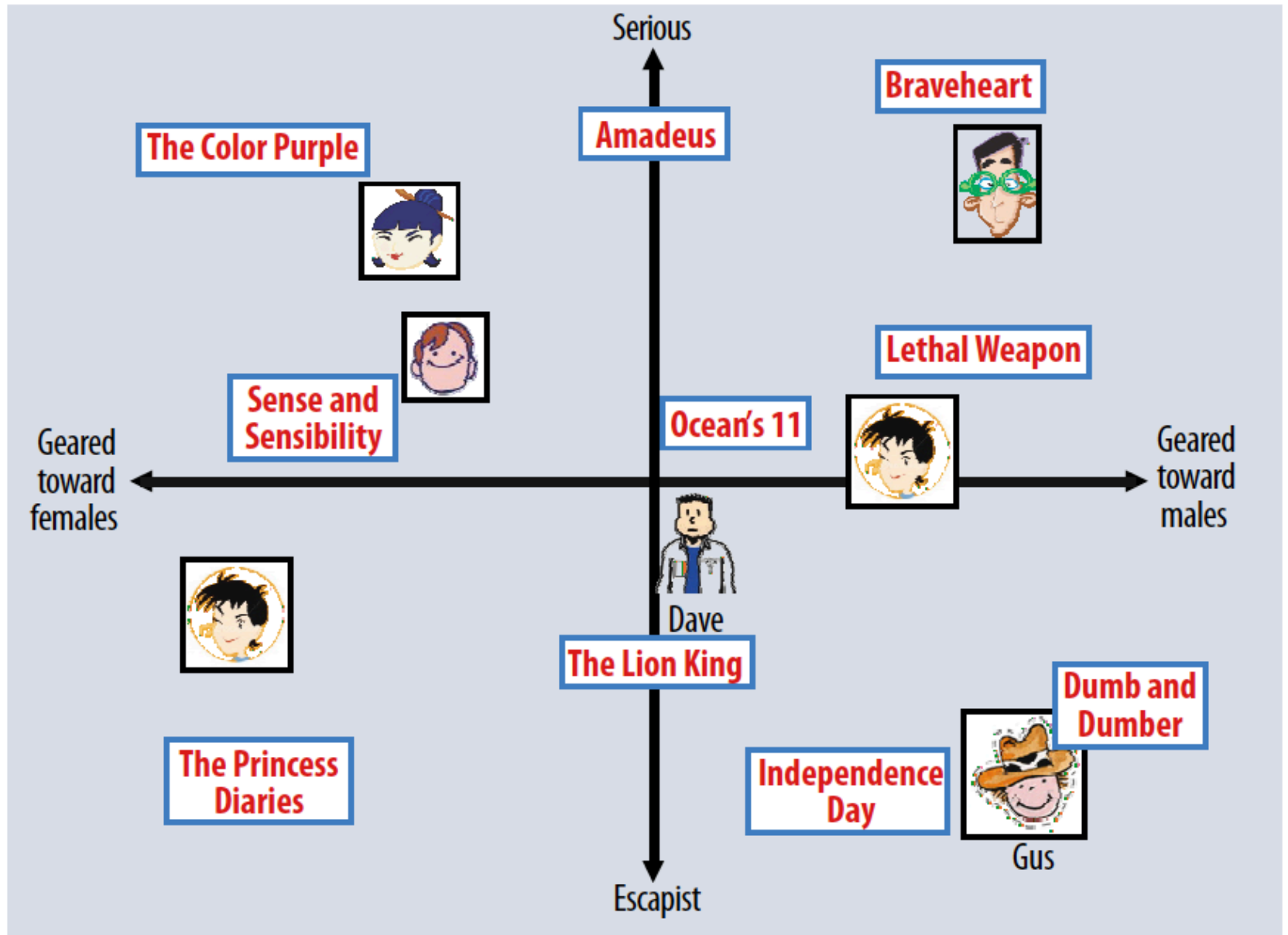
- ▶ Netflix competition: 500k users, 20k movies, 100M movie ratings, 3M question marks

Recovering latent factors in a matrix



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

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)	? (4.4)	4 (3.8)	2 (2.3)
Bob (1.21)	3 (2.7)	2 (2.3)	? (1.4)
Charlie (2.30)	5 (5.2)	? (4.4)	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

user-specific bias term

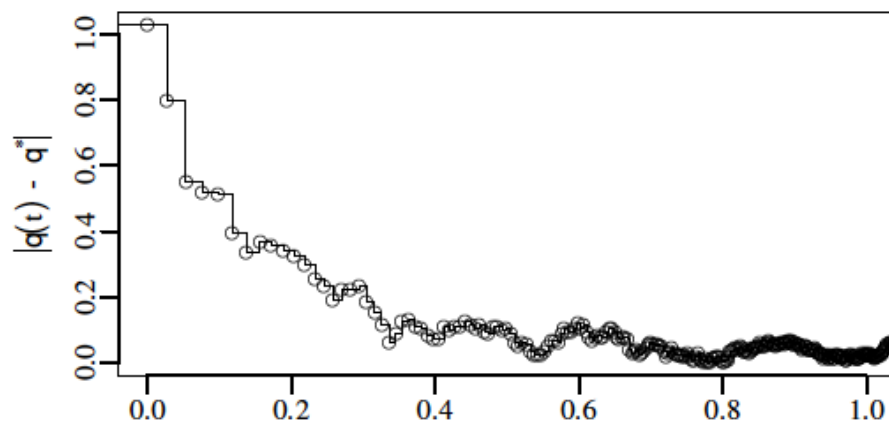
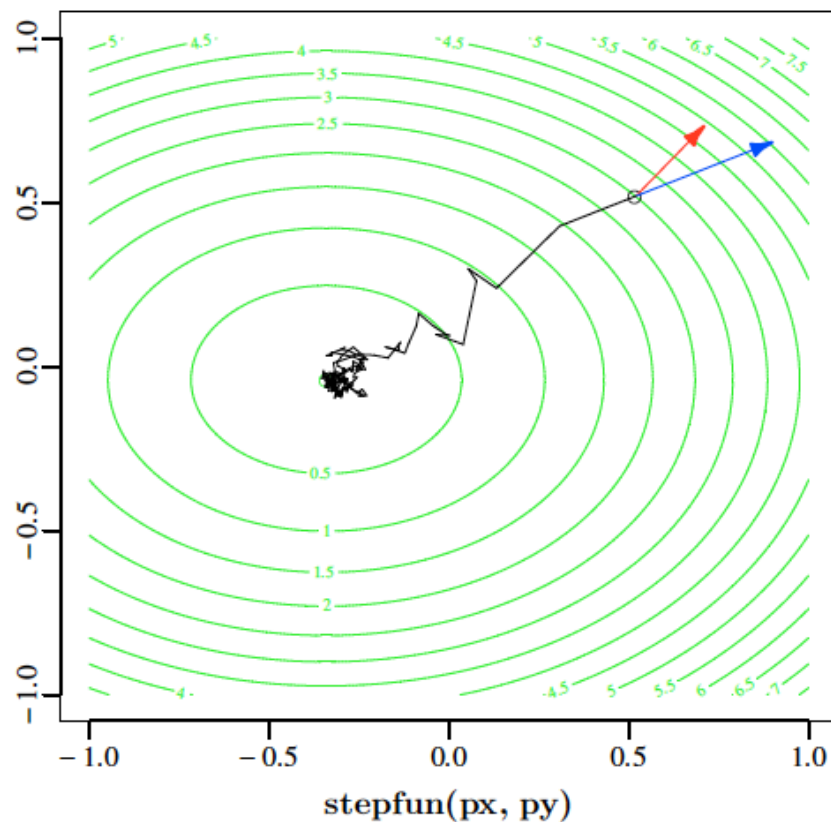
movie-specific bias term

Stochastic Gradient Descent

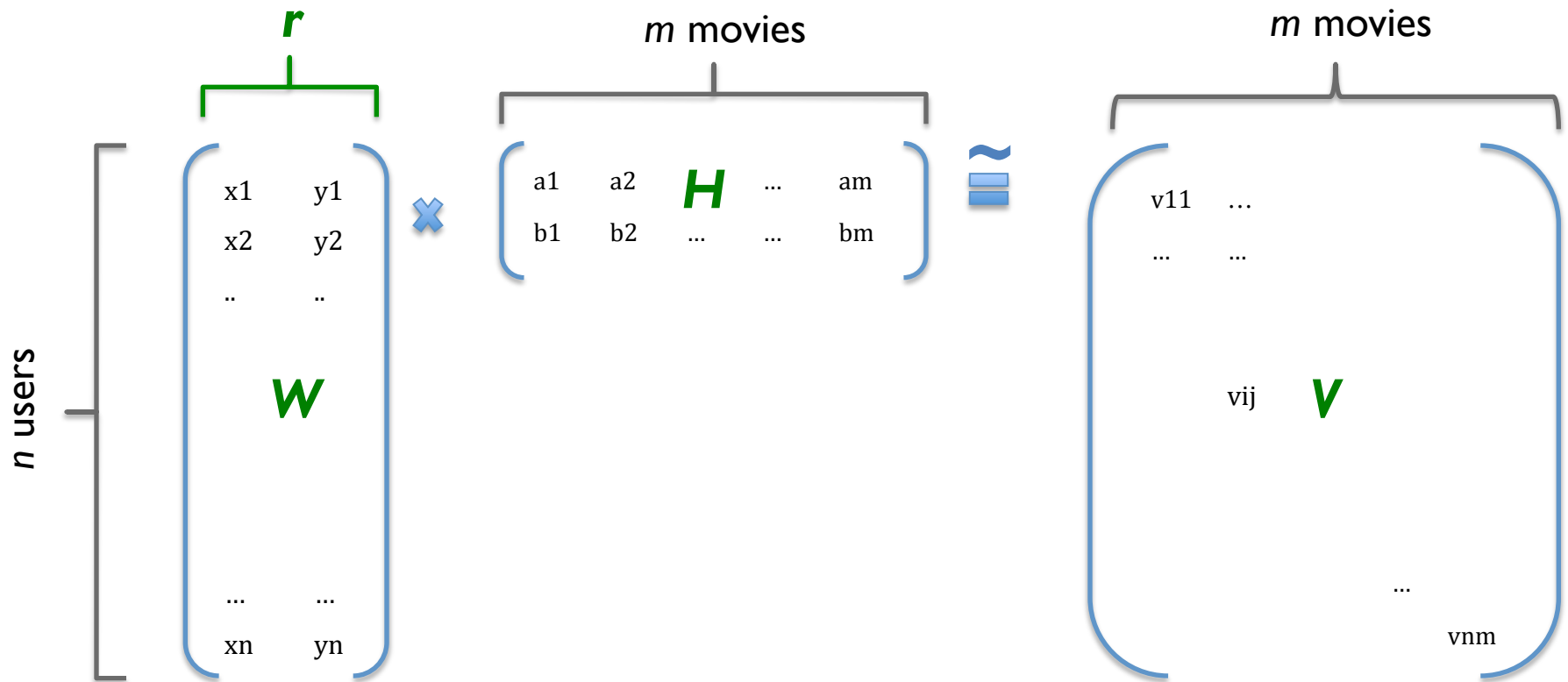
- ▶ Find minimum θ^* of function L
- ▶ Pick a starting point θ_0
- ▶ Approximate gradient $\hat{L}'(\theta_0)$
- ▶ Jump “approximately” downhill
- ▶ Stochastic difference equation

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

- ▶ Under certain conditions, asymptotically approximates (continuous) gradient descent

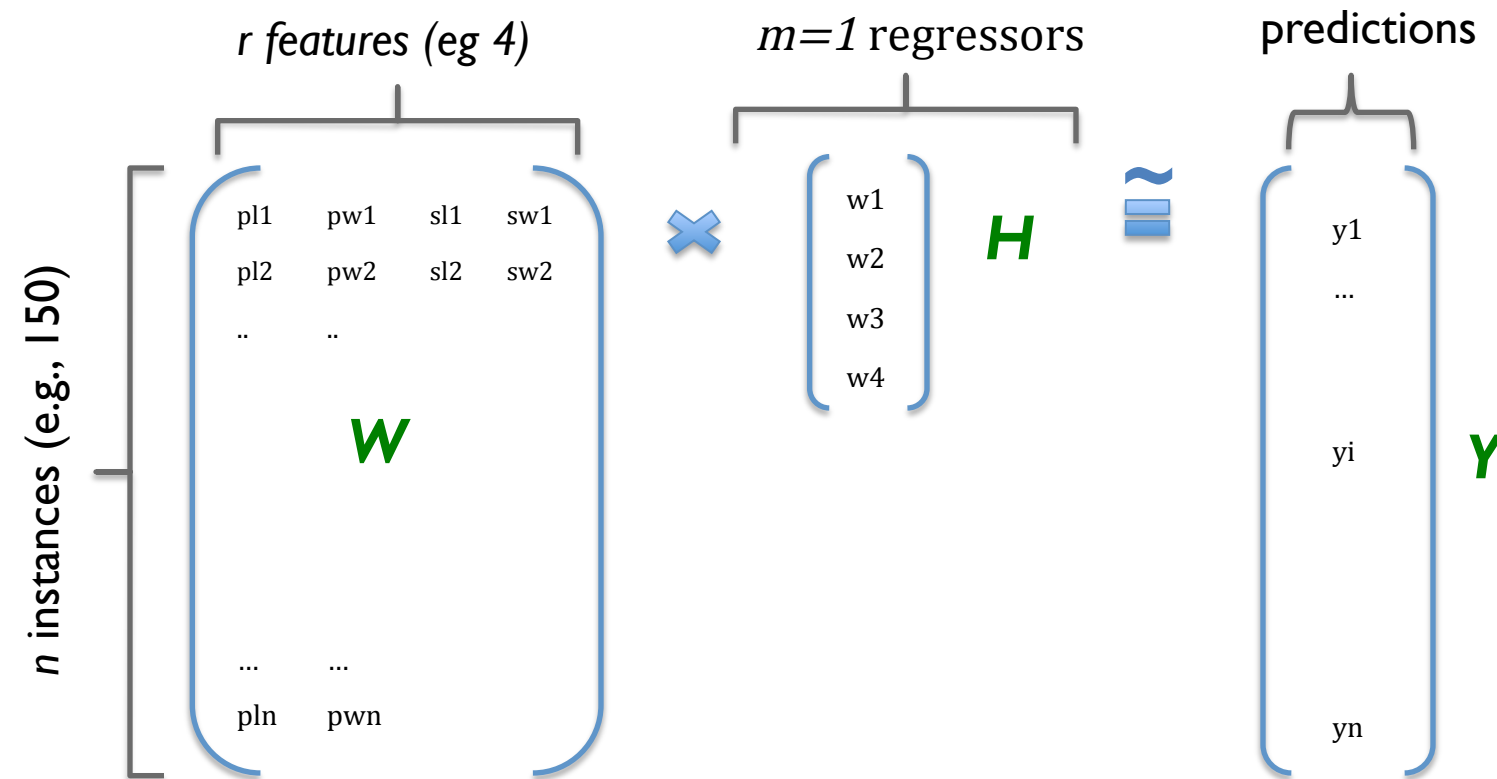


Recovering latent factors in a matrix



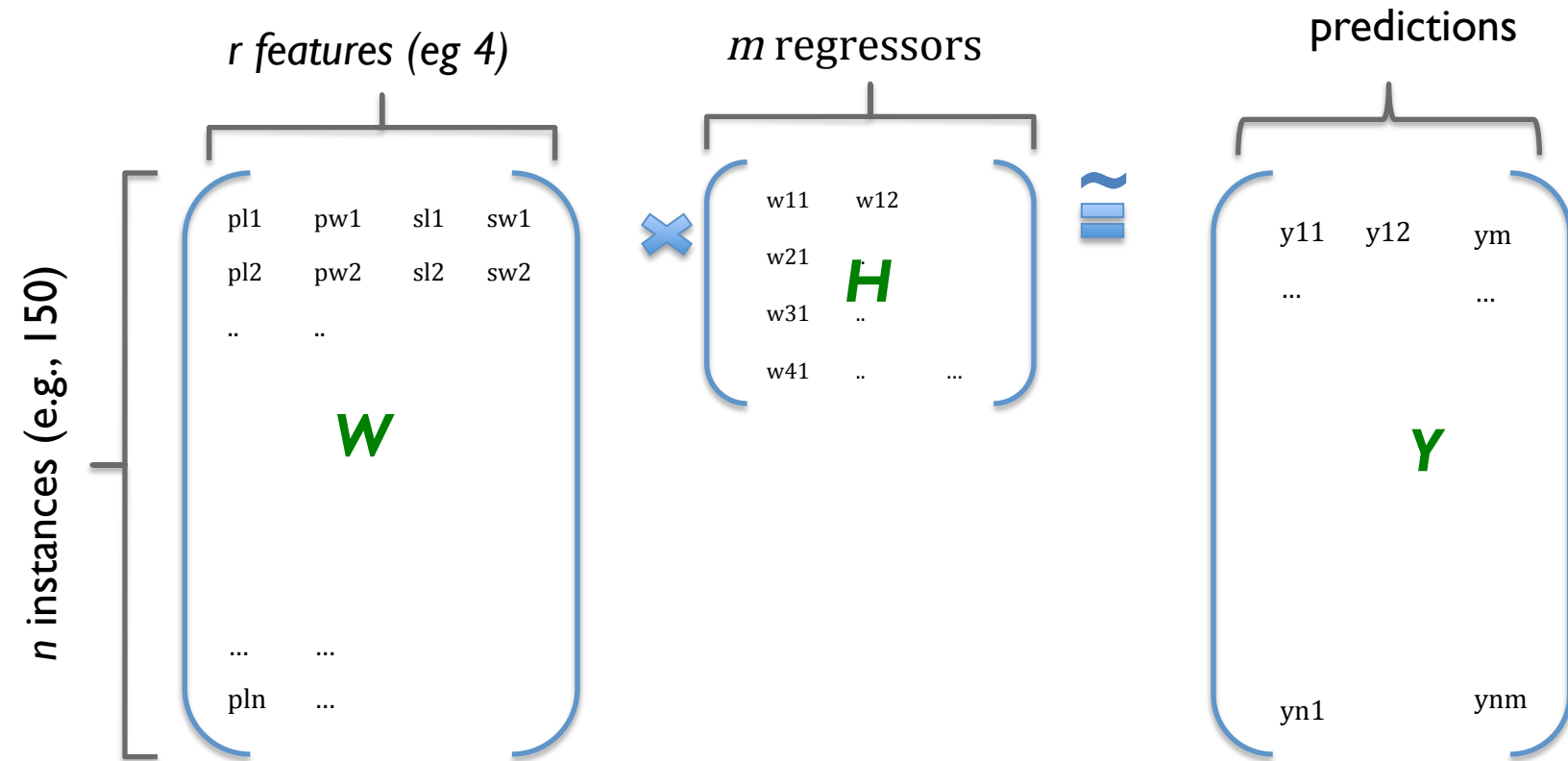
$V[i,j]$ = user i 's rating of movie j

... is like Linear Regression ...



$Y[i,1]$ = instance i 's prediction

.. for many outputs at once...



... where we also have to find the dataset!

$Y[l,j]$ = instance l 's prediction for regression task j

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



Matrix factorization as SGD - why does this work?

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

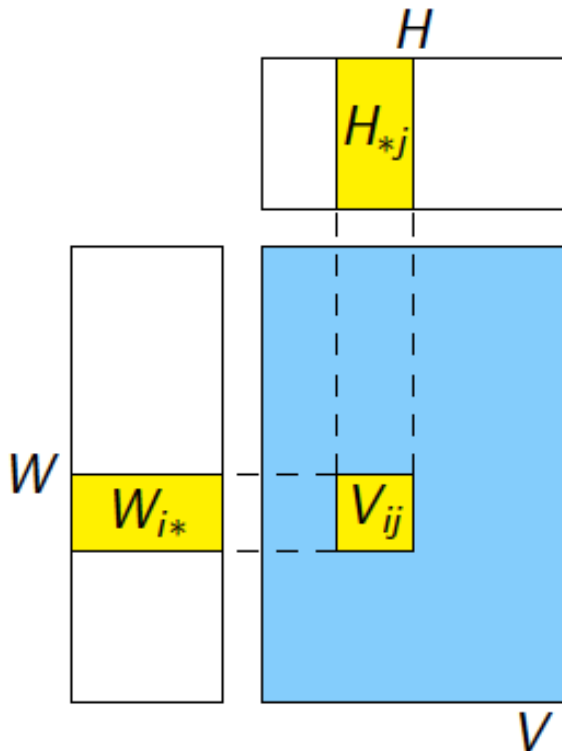
step size



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

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

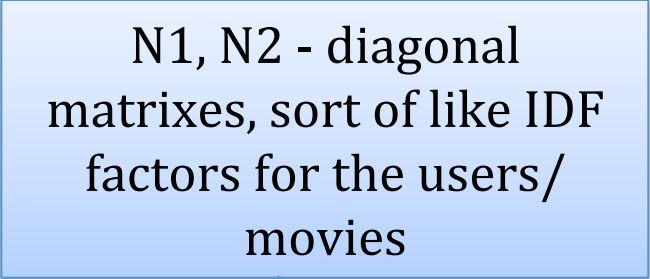
where $Z_{*j} = \{i : (i, j) \in Z\}$.

Think for SGD for logistic regression

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

What loss functions are possible?

N_1, N_2 - diagonal matrixes, sort of like IDF factors for the users/ movies



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

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

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

What loss functions are possible?

Loss Function	Definition and Derivatives
---------------	----------------------------

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

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

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

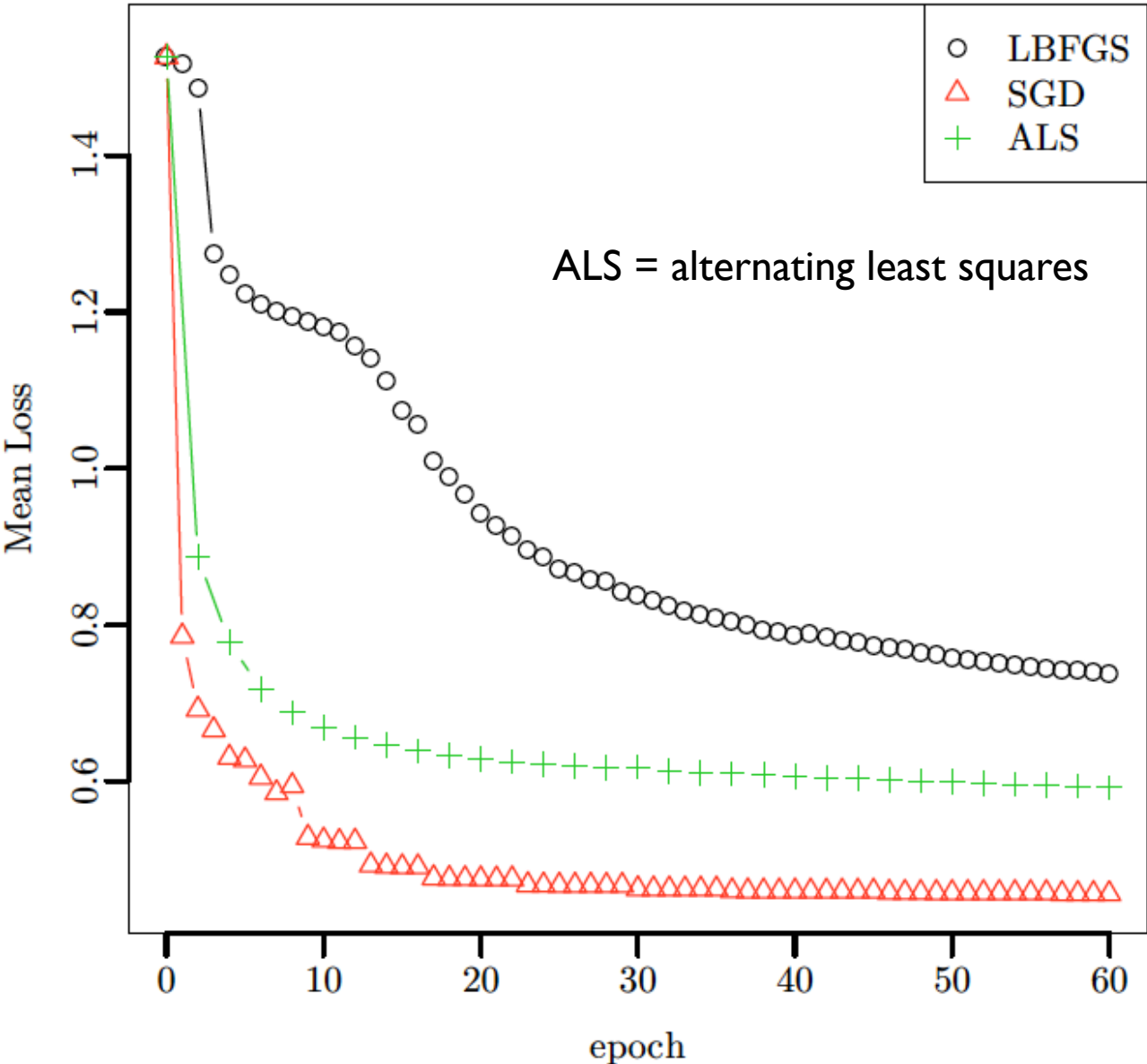
What loss functions are possible?

Loss Function Definition and Derivatives

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

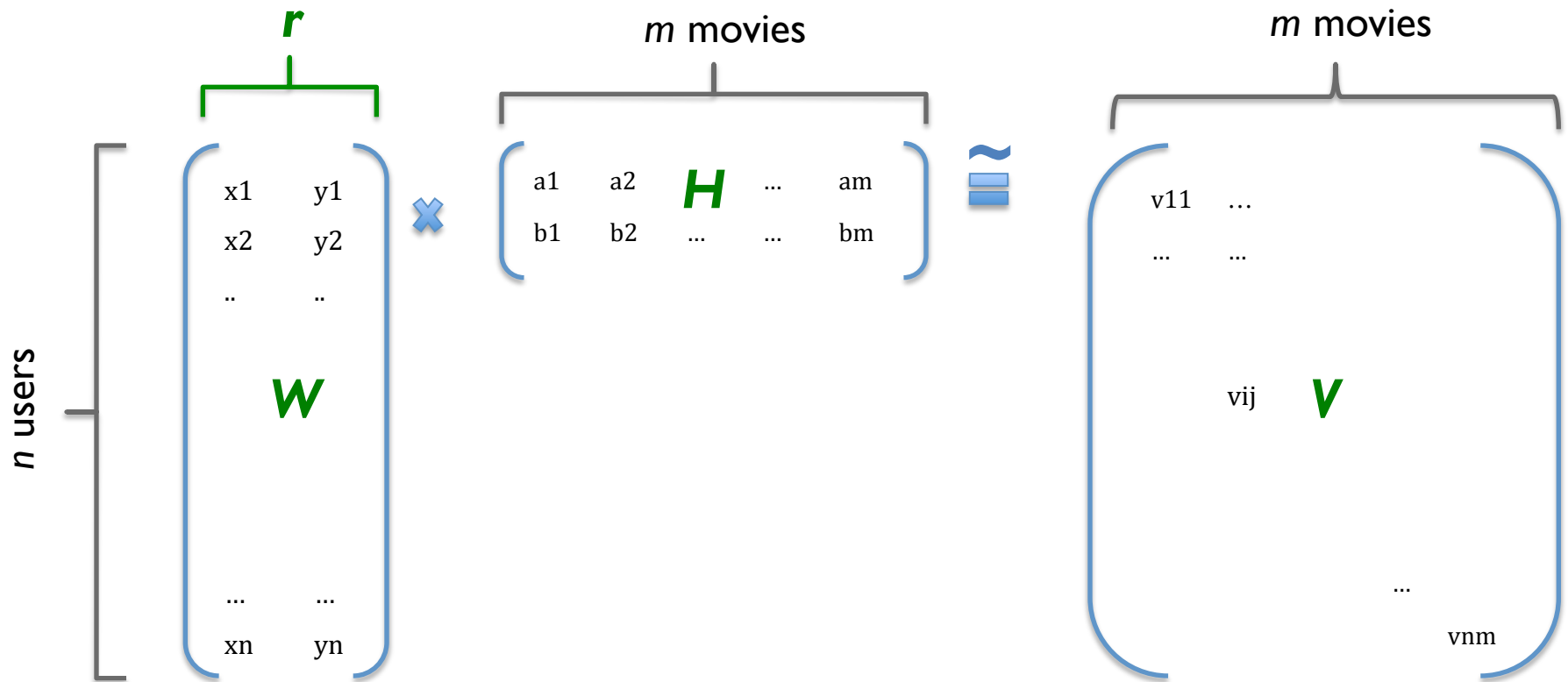
$$\frac{\partial}{\partial \mathbf{W}_{ik}} L_{ij} = -2(\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij})\mathbf{H}_{kj} + 2\lambda \frac{\mathbf{W}_{ik}}{N_{i*}}$$
$$\frac{\partial}{\partial \mathbf{H}_{kj}} L_{ij} = -2(\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij})\mathbf{W}_{ik} + 2\lambda \frac{\mathbf{H}_{kj}}{N_{*j}}$$

Stochastic Gradient Descent on Netflix Data



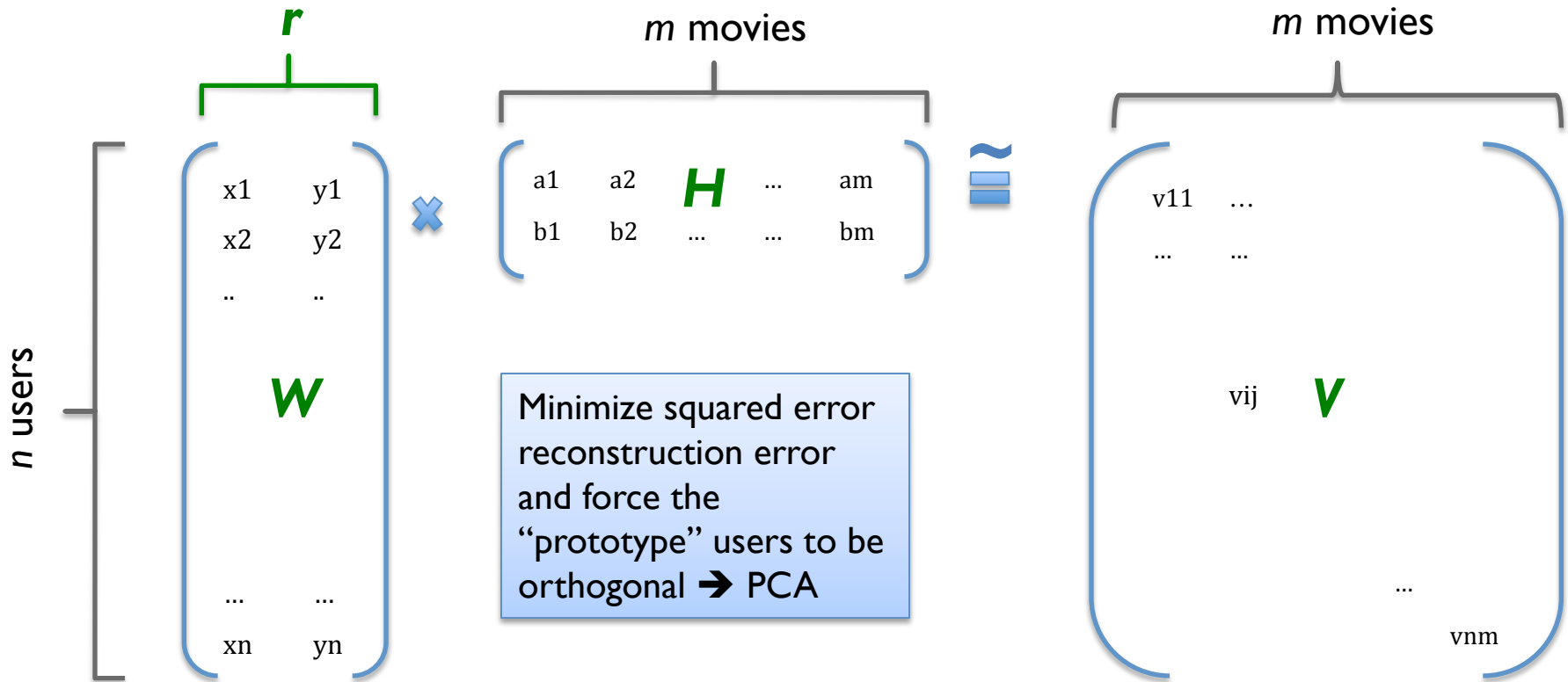
Wrapup: Matrix Multiplications in Machine Learning

Recovering latent factors in a matrix



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

... vs PCA

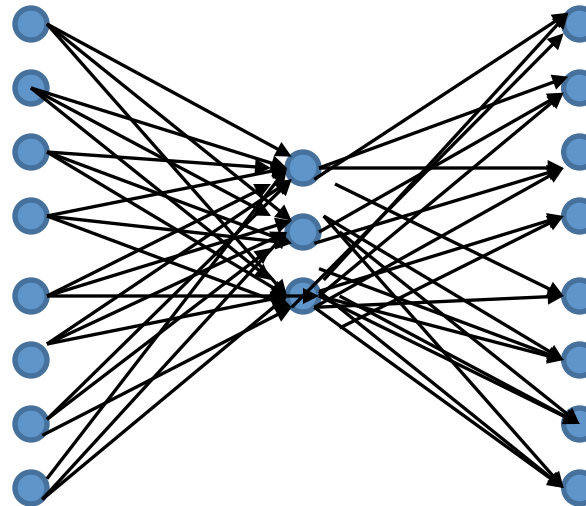


$V[i,j]$ = user i 's rating of movie j

... vs autoencoders & nonlinear PCA

- Assume we would like to learn the following (trivial?) output function:
- Using the following network:
- With *linear* hidden units, how do the weights match up to W and H ?

Input	Output
00000001	00000001
00000010	00000010
00000101	00000100
00001000	00001000
00010000	00010000
00100000	00100000
01000000	01000000
10000000	10000000

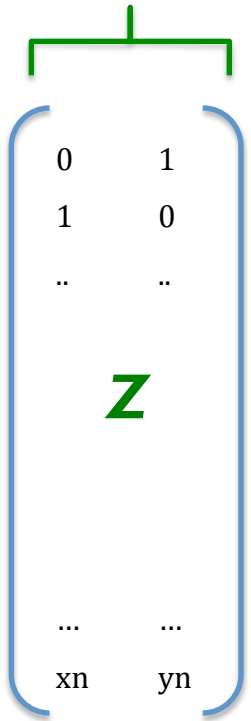


..... vs k-means

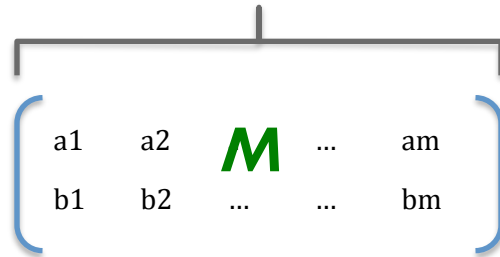
indicators for r clusters

clusters

n examples



cluster means



original data set

