



10-601 Introduction to Machine Learning

Machine Learning Department
School of Computer Science
Carnegie Mellon University

Ensemble Methods + Recommender Systems

Matt Gormley
Lecture 28
Apr. 29, 2019

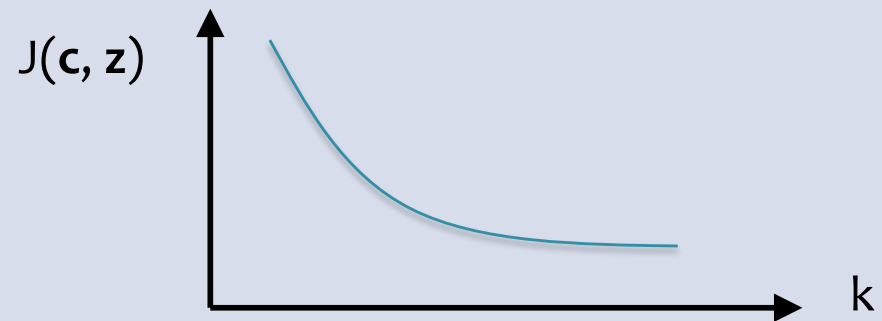
Reminders

- **Homework 9: Learning Paradigms**
 - Out: Wed, Apr 24
 - Due: Wed, May 1 at 11:59pm
 - Can only be submitted up to 3 days late, so we can return grades before final exam
- **Today's In-Class Poll**
 - <http://p28.mlcourse.org>

Q&A

Q: In k-Means, since we don't have a validation set, how do we pick k ?

A: Look at the training objective function as a function of k and pick the value at the “elbo” of the curve.



Q: What if our random initialization for k-Means gives us poor performance?

A: Do **random restarts**: that is, run k-means from scratch, say, 10 times and pick the run that gives the lowest training objective function value.

The objective function is **nonconvex**, so we're just looking for the best local minimum.

ML Big Picture

Learning Paradigms:

What data is available and when? What form of prediction?

- supervised learning
- unsupervised learning
- semi-supervised learning
- reinforcement learning
- active learning
- imitation learning
- domain adaptation
- online learning
- density estimation
- recommender systems
- feature learning
- manifold learning
- dimensionality reduction
- ensemble learning
- distant supervision
- hyperparameter optimization

Theoretical Foundations:

What principles guide learning?

- ☐ probabilistic
- ☐ information theoretic
- ☐ evolutionary search
- ☐ ML as optimization

Problem Formulation:

What is the structure of our output prediction?

boolean	Binary Classification
categorical	Multiclass Classification
ordinal	Ordinal Classification
real	Regression
ordering	Ranking
multiple discrete	Structured Prediction
multiple continuous	(e.g. dynamical systems)
both discrete & cont.	(e.g. mixed graphical models)

Facets of Building ML Systems:

How to build systems that are robust, efficient, adaptive, effective?

1. Data prep
2. Model selection
3. Training (optimization / search)
4. Hyperparameter tuning on validation data
5. (Blind) Assessment on test data

Big Ideas in ML:

Which are the ideas driving development of the field?

- inductive bias
- generalization / overfitting
- bias-variance decomposition
- generative vs. discriminative
- deep nets, graphical models
- PAC learning
- distant rewards

Application Areas

Key challenges?

NLP, Speech, Computer Vision, Robotics, Medicine, Search

Outline for Today

We'll talk about two distinct topics:

1. **Ensemble Methods:** combine or learn multiple classifiers into one
(i.e. a family of algorithms)
2. **Recommender Systems:** produce recommendations of what a user will like
(i.e. the solution to a particular type of task)

We'll use a prominent example of a recommender systems (the Netflix Prize) to motivate both topics...

RECOMMENDER SYSTEMS

Recommender Systems

A Common Challenge:

- Assume you're a company selling **items** of some sort: movies, songs, products, etc.
- Company collects millions of **ratings** from **users** of their **items**
- To maximize profit / user happiness, you want to **recommend** items that users are likely to want

Recommender Systems



You could be seeing useful stuff here!
Sign in to get your order status, balances and rewards.



Recommended for you, Matt



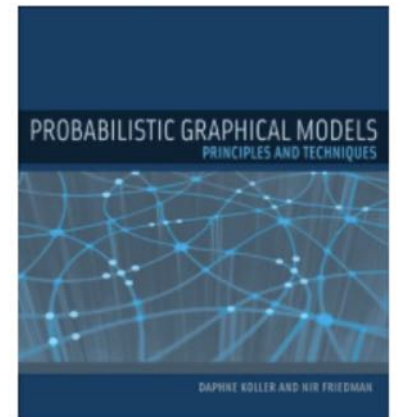
Buy It Again in Grocery
14 ITEMS



Buy It Again in Pets
6 ITEMS

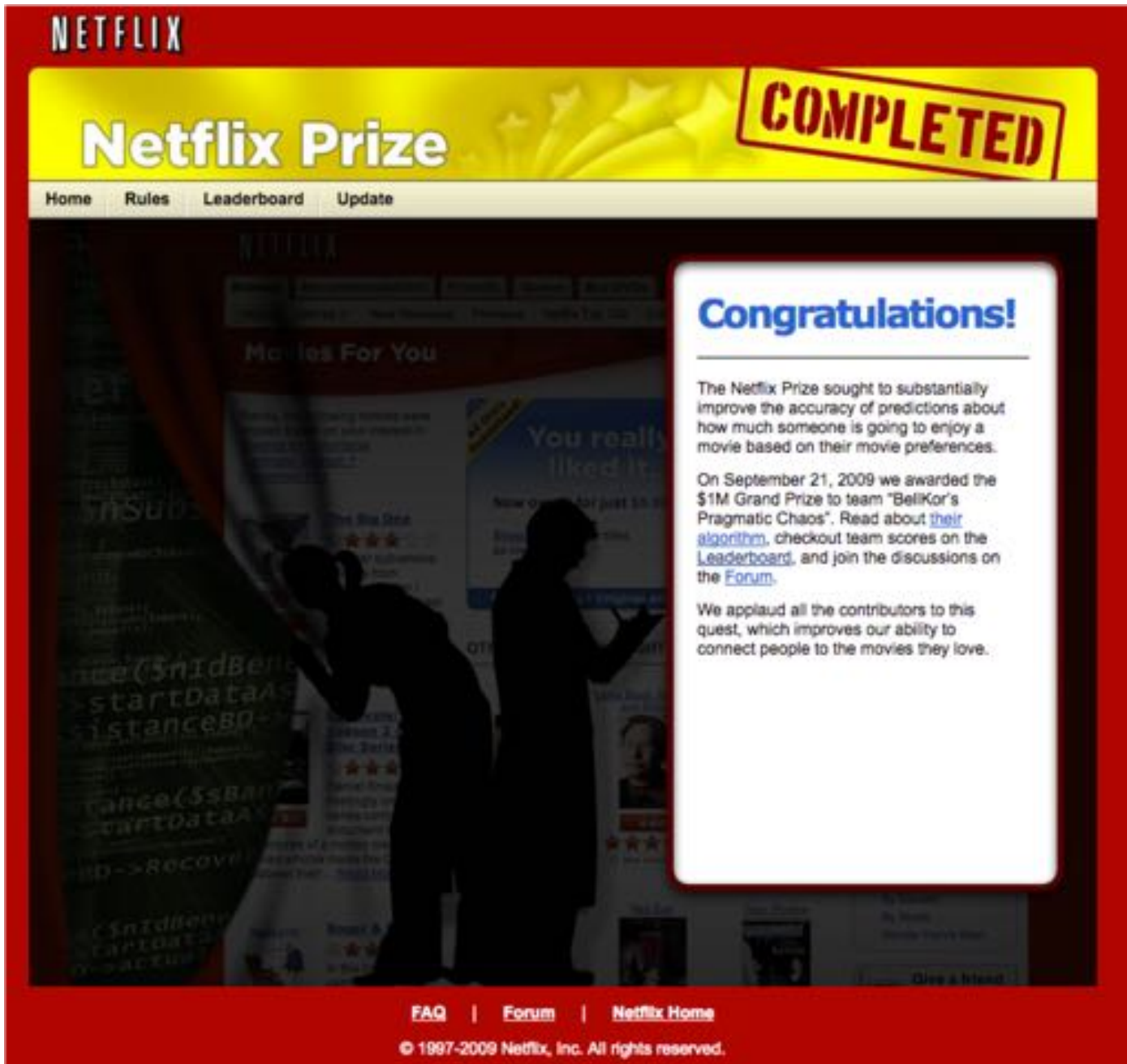


Buy It Again in Baby Products
5 ITEMS



Engineering Books
86 ITEMS

Recommender Systems



The image is a screenshot of the Netflix Prize announcement page. At the top, the Netflix logo is on the left, and a large yellow banner with the text "Netfix Prize" and a "COMPLETED" stamp is on the right. Below the banner is a navigation bar with links: Home, Rules, Leaderboard, and Update. The main content area features a large, dark, stylized image of two people silhouettes, one holding a large banner that reads "Netflix Prize". To the right of this image is a white box with the heading "Congratulations!" and text announcing the \$1M Grand Prize to team "BellKor's Pragmatic Chaos". The background of the page is a dark red color. At the bottom, there is a footer with links for FAQ, Forum, and Netflix Home, and a copyright notice for 1997-2009 Netflix, Inc.

NETFLIX

Netfix Prize

COMPLETED

Home Rules Leaderboard Update

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.

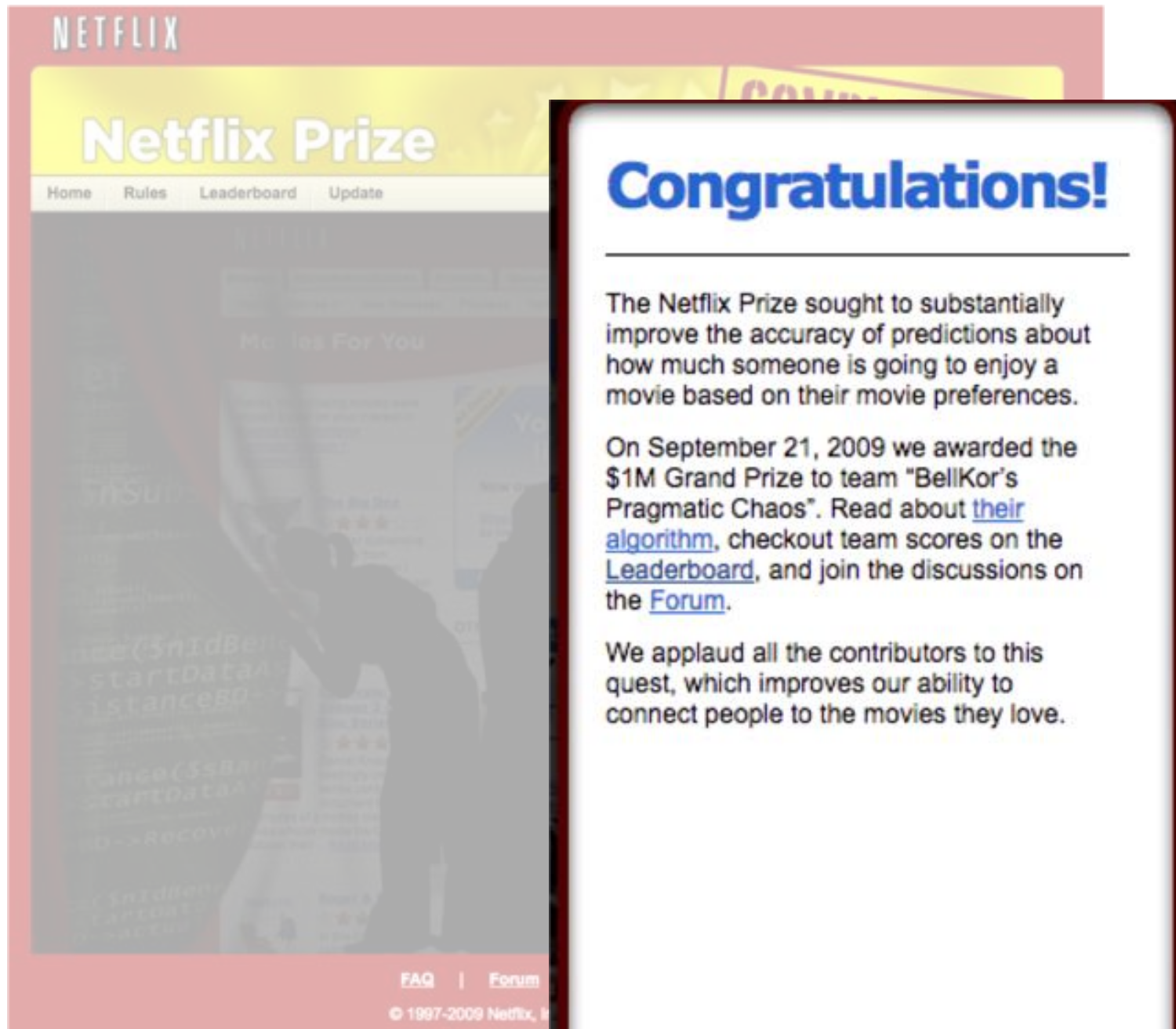
On September 21, 2009 we awarded the \$1M Grand Prize to team "BellKor's Pragmatic Chaos". Read about [their algorithm](#), checkout team scores on the [Leaderboard](#), and join the discussions on the [Forum](#).

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

FAQ | Forum | Netflix Home

© 1997-2009 Netflix, Inc. All rights reserved.

Recommender Systems



The image shows a screenshot of the Netflix Prize website. The top navigation bar includes links for Home, Rules, Leaderboard, and Update. The main content area features a "Movies For You" section with a list of movie recommendations. A large, semi-transparent watermark of a person's silhouette is overlaid on the page. At the bottom, there are links for FAQ and Forum, and a copyright notice for 1997-2009 Netflix, Inc.

Netflix Prize

Home Rules Leaderboard Update

Movies For You

FAQ | Forum

© 1997-2009 Netflix, Inc.

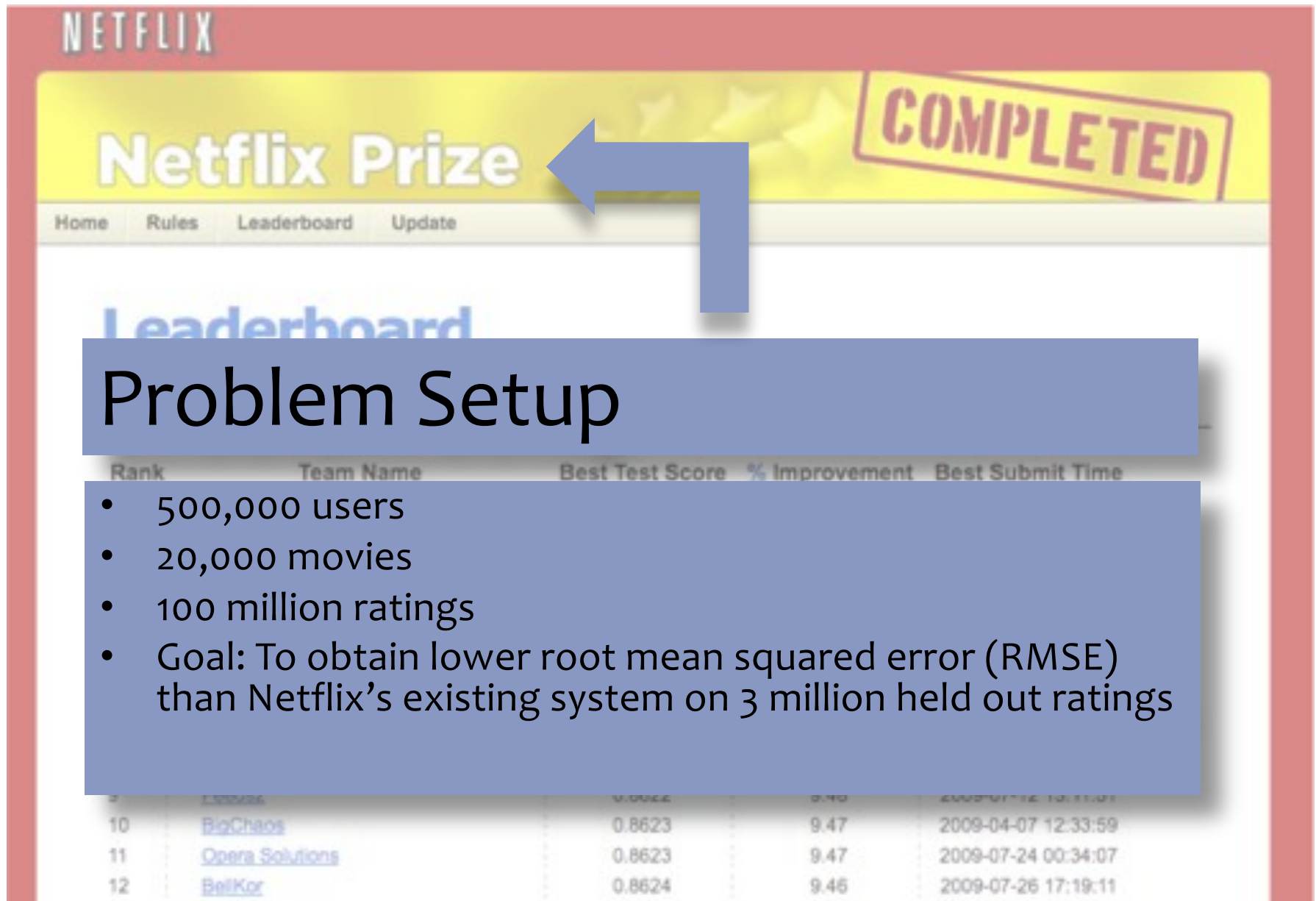
Congratulations!

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.

Recommender Systems



The image is a screenshot of the Netflix Prize website. At the top, the 'NETFLIX' logo is visible. Below it, a yellow banner features the text 'Netflix Prize' and a 'COMPLETED' stamp. A blue arrow points from the 'COMPLETED' stamp to the 'Netflix Prize' text. Below the banner is a navigation bar with links: 'Home', 'Rules', 'Leaderboard', and 'Update'. The 'Leaderboard' link is highlighted. Below the navigation bar, the word 'Leaderboard' is written in large blue letters. A large blue box with the text 'Problem Setup' is overlaid on the page. Below this box, a table shows the leaderboard data.

Problem Setup

- 500,000 users
- 20,000 movies
- 100 million ratings
- Goal: To obtain lower root mean squared error (RMSE) than Netflix's existing system on 3 million held out ratings

Rank	Team Name	Best Test Score	% Improvement	Best Submit Time
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

ENSEMBLE METHODS

Recommender Systems

Netflix Prize **COMPLETED**

Home Rules Leaderboard Update

Leaderboard

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

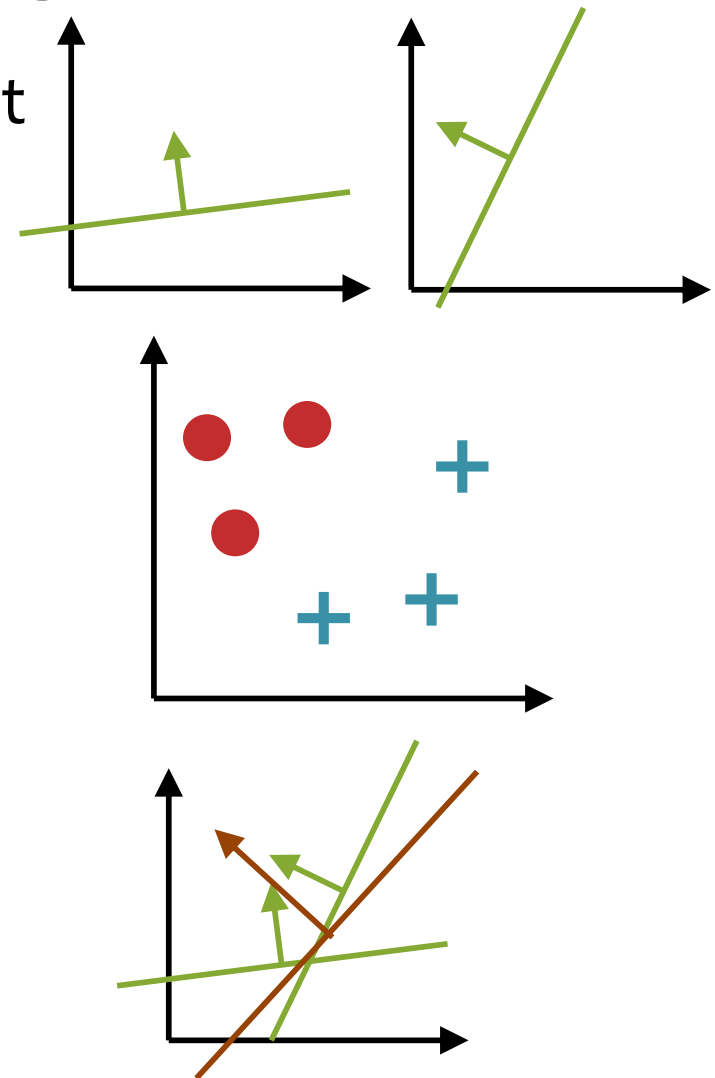
Rank	Team Name	Best Test Score	% Improvement	Best Submit Time
Grand Prize - RMSE = 0.8567 - Winning Team: BellKor's Pragmatic Chaos				
1	BellKor's Pragmatic Chaos	0.8567	10.06	2009-07-26 18:18:28
2	The Ensemble	0.8567	10.06	2009-07-26 18:38:22
3	Grand Prize Team	0.8582	9.90	2009-07-10 21:24:40
4	Opera Solutions and Vandelay United	0.8588	9.84	2009-07-10 01:12:31
5	Vandelay Industries I	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

Top performing systems were ensembles

Weighted Majority Algorithm

(Littlestone & Warmuth, 1994)

- **Given:** pool A of binary classifiers (that you know nothing about)
- **Data:** stream of examples (i.e. online learning setting)
- **Goal:** design a new learner that uses the predictions of the pool to make new predictions
- **Algorithm:**
 - Initially weight all classifiers equally
 - Receive a training example and predict the (weighted) majority vote of the classifiers in the pool
 - Down-weight classifiers that contribute to a mistake by a factor of β



Weighted Majority Algorithm

(Littlestone & Warmuth, 1994)

Suppose we have a pool of T binary classifiers $\mathcal{A} = \{h_1, \dots, h_T\}$ where $h_t : \mathbb{R}^M \rightarrow \{+1, -1\}$. Let α_t be the weight for classifier h_t .

Algorithm 1 Weighted Majority Algorithm

- 1: **procedure** WEIGHTEDMAJORITY(\mathcal{A}, β)
- 2: Initialize classifier weights $\alpha_t = 1, \forall t \in \{1, \dots, T\}$
- 3: **for** each training example (\mathbf{x}, y) **do**
- 4: Predict majority vote class (splitting ties randomly)

$$\hat{h}(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$$

- 5: **if** a mistake is made $\hat{h}(x) \neq y$ **then**
 - 6: **for** each classifier $t \in \{1, \dots, T\}$ **do**
 - 7: **if** $h_t(x) \neq y$, then $\alpha_t \leftarrow \beta \alpha_t$
-

Weighted Majority Algorithm

(Littlestone & Warmuth, 1994)

Theorem 0.1 (Littlestone & Warmuth, 1994). *If the Weighted Majority Algorithm is applied to a pool \mathcal{A} of classifiers, and if each algorithm makes at most m mistakes on the sequence of examples, then the total number of mistakes is upper bounded by $2.4(\log |\mathcal{A}| + m)$.*



This is a “mistake bound” of the variety we saw for the Perceptron algorithm

ADABOOST

Comparison

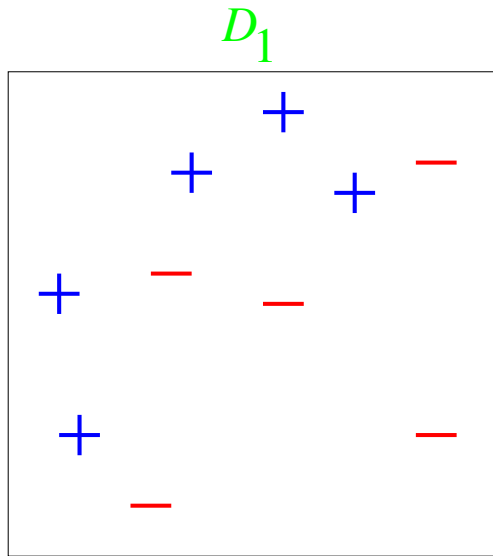
Weighted Majority Algorithm

- an example of an ensemble method
- assumes the classifiers are learned ahead of time
- only learns (majority vote) weight for each classifiers

AdaBoost

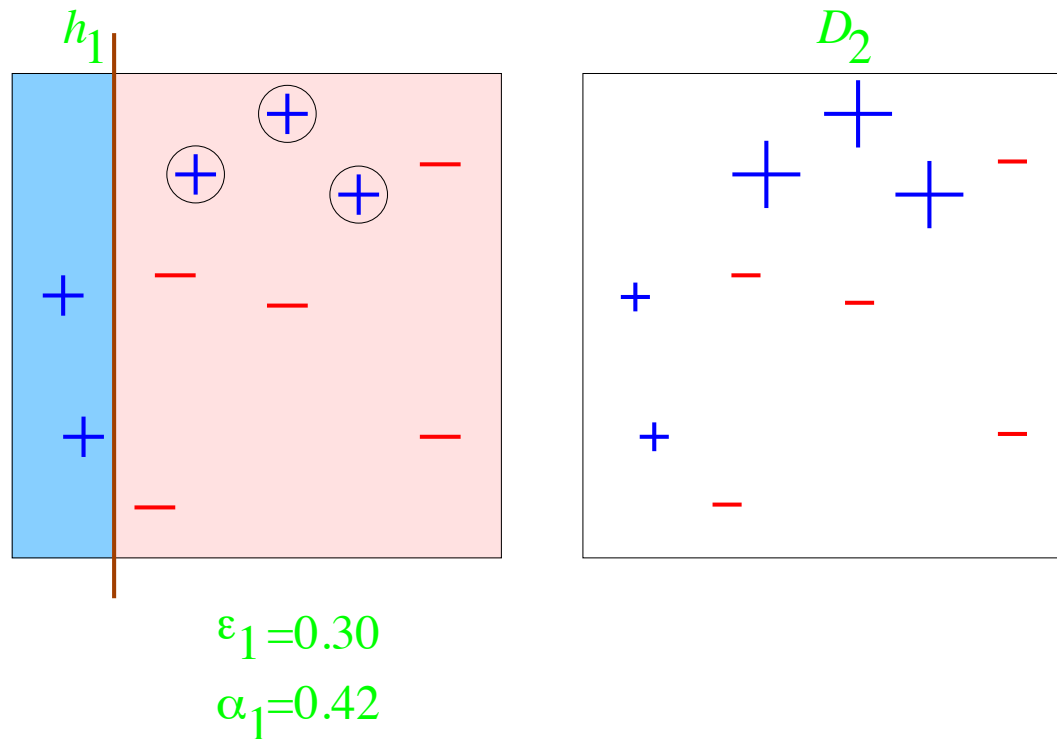
- an example of a boosting method
- simultaneously learns:
 - the classifiers themselves
 - (majority vote) weight for each classifiers

AdaBoost: Toy Example

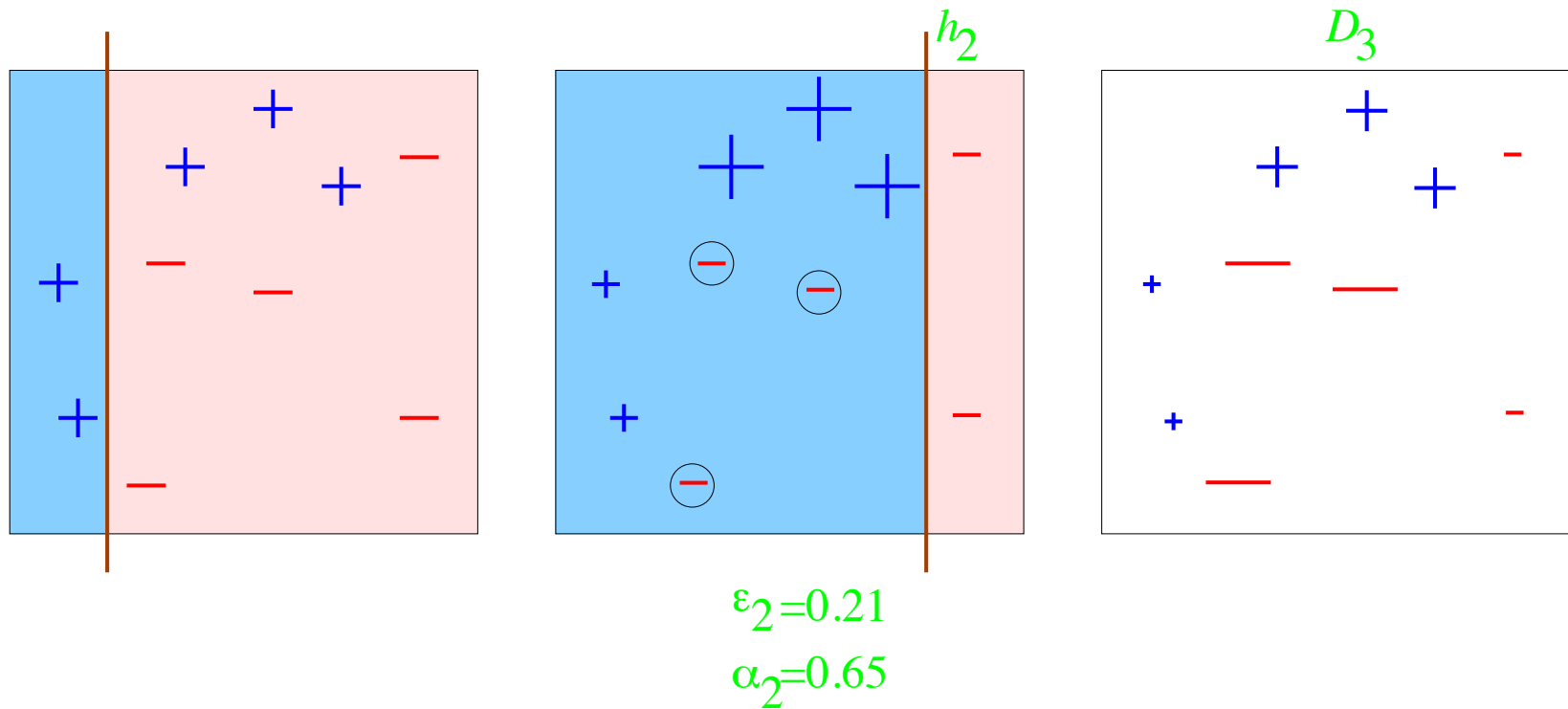


weak classifiers = vertical or horizontal half-planes

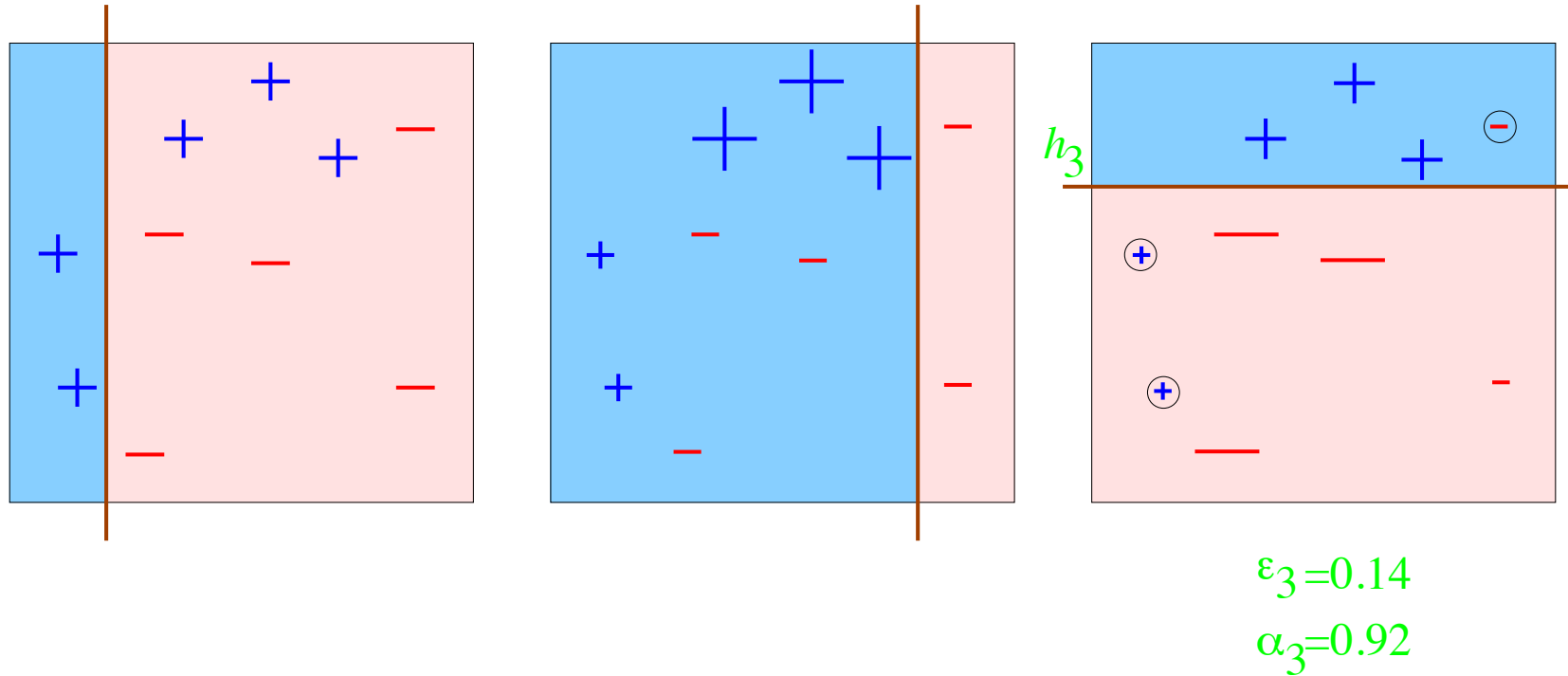
AdaBoost: Toy Example



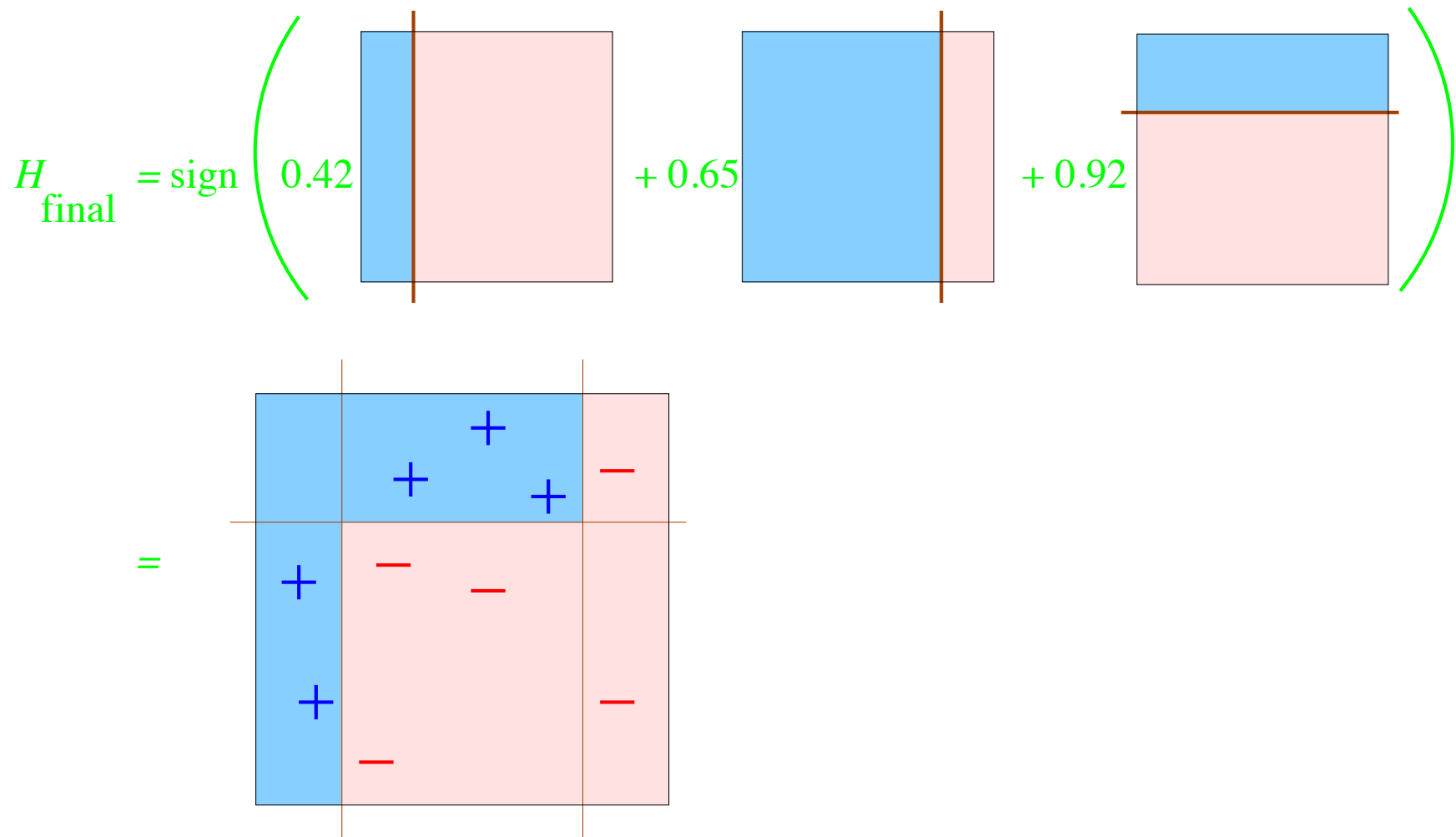
AdaBoost: Toy Example



AdaBoost: Toy Example



AdaBoost: Toy Example



AdaBoost

Given: $(x_1, y_1), \dots, (x_m, y_m)$ where $x_i \in X, y_i \in Y = \{-1, +1\}$

Initialize $D_1(i) = 1/m$.

For $t = 1, \dots, T$:

- Train weak learner using distribution D_t .
- Get weak hypothesis $h_t : X \rightarrow \{-1, +1\}$ with error

$$\epsilon_t = \Pr_{i \sim D_t} [h_t(x_i) \neq y_i].$$

- Choose $\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$.
- Update:

$$\begin{aligned} D_{t+1}(i) &= \frac{D_t(i)}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{if } h_t(x_i) = y_i \\ e^{\alpha_t} & \text{if } h_t(x_i) \neq y_i \end{cases} \\ &= \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t} \end{aligned}$$

where Z_t is a normalization factor (chosen so that D_{t+1} will be a distribution).

Output the final hypothesis:

$$H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right).$$

AdaBoost

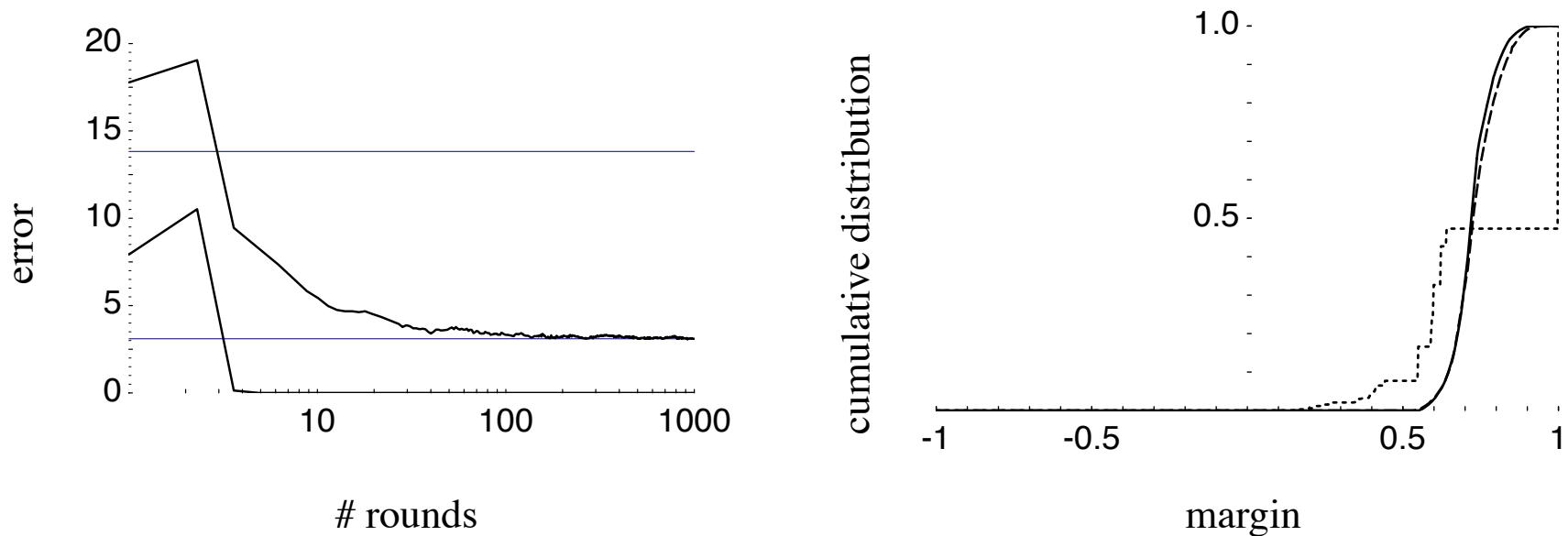


Figure 2: Error curves and the margin distribution graph for boosting C4.5 on the letter dataset as reported by Schapire et al. [41]. *Left*: the training and test error curves (lower and upper curves, respectively) of the combined classifier as a function of the number of rounds of boosting. The horizontal lines indicate the test error rate of the base classifier as well as the test error of the final combined classifier. *Right*: The cumulative distribution of margins of the training examples after 5, 100 and 1000 iterations, indicated by short-dashed, long-dashed (mostly hidden) and solid curves, respectively.

Learning Objectives

Ensemble Methods / Boosting

You should be able to...

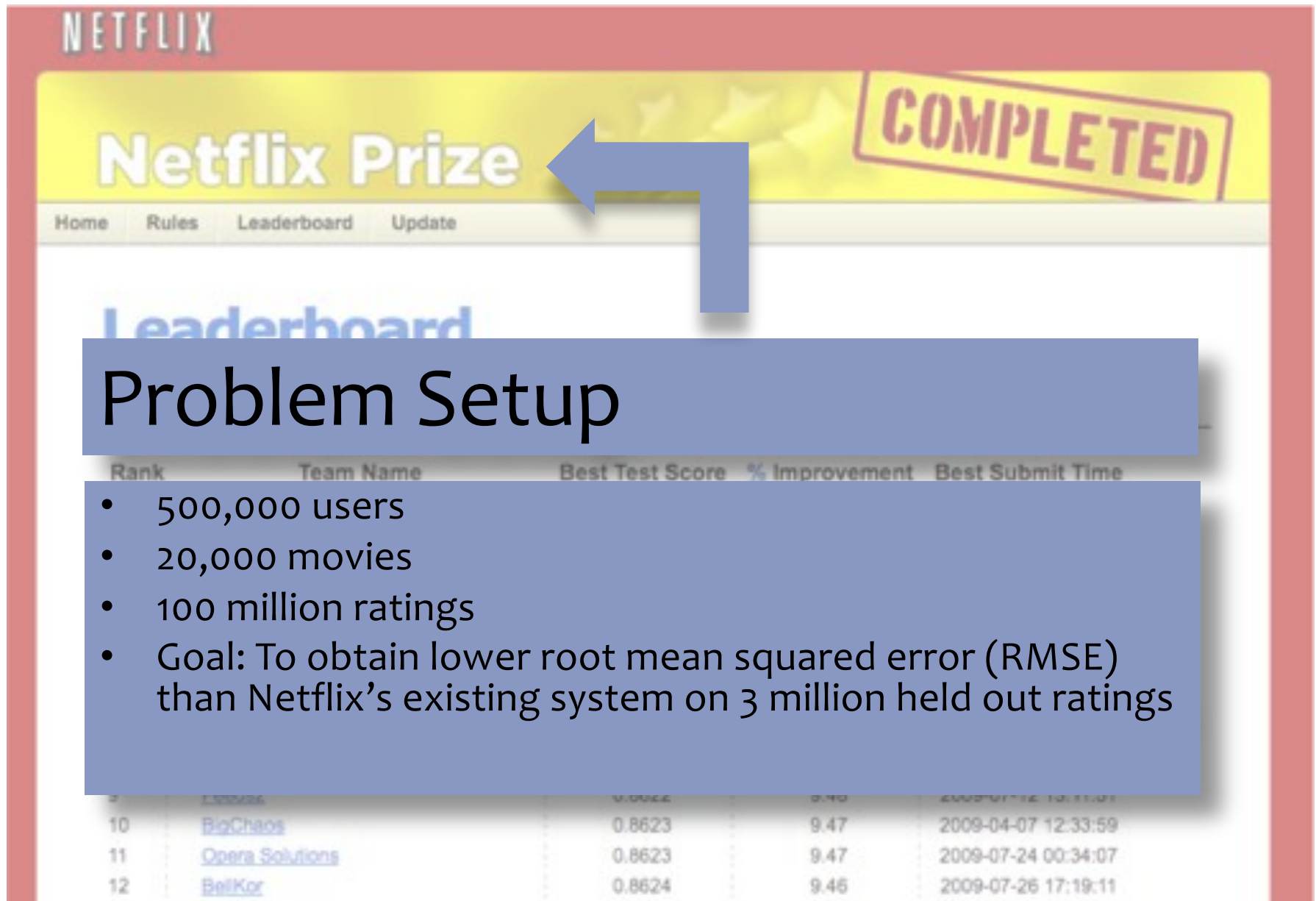
1. Implement the Weighted Majority Algorithm
2. Implement AdaBoost
3. Distinguish what is learned in the Weighted Majority Algorithm vs. Adaboost
4. Contrast the theoretical result for the Weighted Majority Algorithm to that of Perceptron
5. Explain a surprisingly common empirical result regarding Adaboost train/test curves

Outline

- **Recommender Systems**
 - Content Filtering
 - Collaborative Filtering (CF)
 - CF: Neighborhood Methods
 - CF: Latent Factor Methods
- **Matrix Factorization**
 - Background: Low-rank Factorizations
 - Residual matrix
 - Unconstrained Matrix Factorization
 - Optimization problem
 - Gradient Descent, SGD, Alternating Least Squares
 - User/item bias terms (matrix trick)
 - Singular Value Decomposition (SVD)
 - Non-negative Matrix Factorization

RECOMMENDER SYSTEMS

Recommender Systems



The screenshot shows the Netflix Prize website interface. At the top, the 'NETFLIX' logo is on the left, and a yellow banner with 'Netflix Prize' and a 'COMPLETED' stamp is on the right. A blue arrow points from the 'COMPLETED' stamp to the 'Netflix Prize' text. Below the banner is a navigation bar with 'Home', 'Rules', 'Leaderboard', and 'Update'. The 'Leaderboard' section is visible, showing a table of team rankings. A large blue box with the text 'Problem Setup' is overlaid on the lower part of the screenshot, containing a bulleted list of details about the competition.

Problem Setup

- 500,000 users
- 20,000 movies
- 100 million ratings
- Goal: To obtain lower root mean squared error (RMSE) than Netflix's existing system on 3 million held out ratings

Rank	Team Name	Best Test Score	% Improvement	Best Submit Time
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

Recommender Systems

NETFLIX

Netfix Prize

COMPLETED

Home Rules Leaderboard Update

Leaderboard

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

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 I	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

Recommender Systems

- **Setup:**

- **Items:**

- movies, songs, products, etc.
(often many thousands)

- **Users:**

- watchers, listeners, purchasers, etc.
(often many millions)

- **Feedback:**

- 5-star ratings, not-clicking ‘next’,
purchases, etc.

- **Key Assumptions:**

- Can represent ratings numerically
as a user/item matrix

- Users only rate a small number of
items (the matrix is sparse)

	Doctor Strange	Star Trek: Beyond	Zootopia
Alice	1		5
Bob	3	4	
Charlie	3	5	2

Two Types of Recommender Systems

Content Filtering

- *Example:* Pandora.com music recommendations (Music Genome Project)
- **Con:** Assumes access to **side information** about items (e.g. properties of a song)
- **Pro:** Got a **new item** to add? No problem, just be sure to include the side information

Collaborative Filtering

- *Example:* Netflix movie recommendations
- **Pro:** Does not assume access to **side information** about items (e.g. does not need to know about movie genres)
- **Con:** Does not work on **new items** that have no ratings

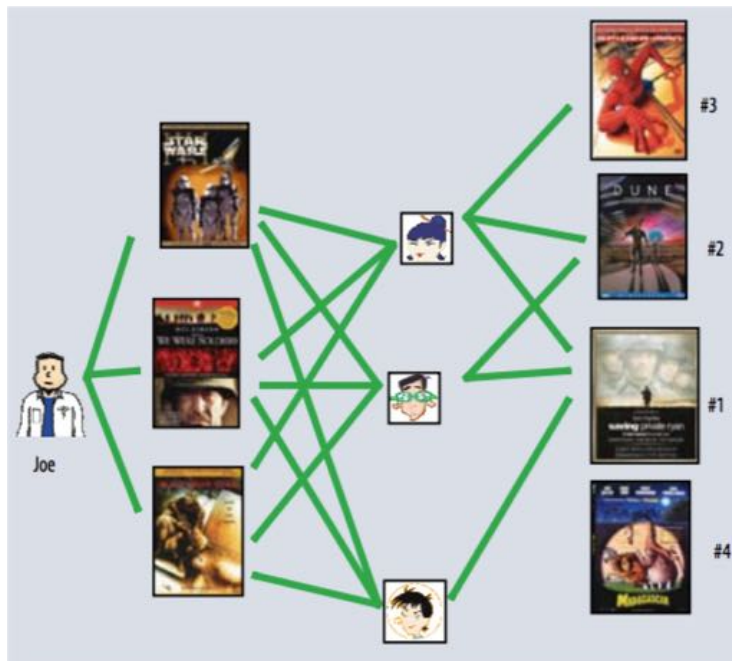
COLLABORATIVE FILTERING

Collaborative Filtering

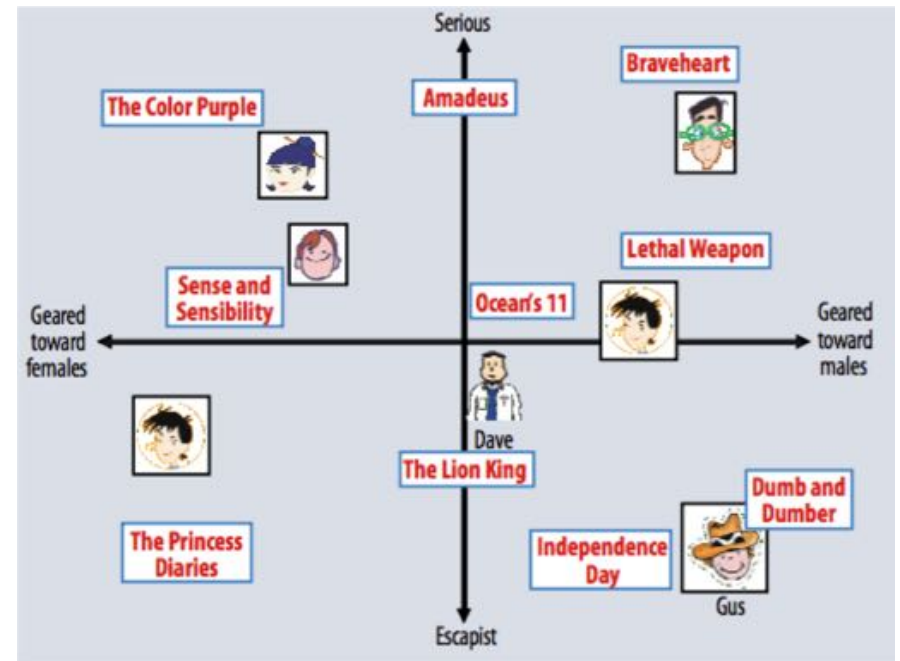
- **Everyday Examples of Collaborative Filtering...**
 - Bestseller lists
 - Top 40 music lists
 - The “recent returns” shelf at the library
 - Unmarked but well-used paths thru the woods
 - The printer room at work
 - “Read any good books lately?”
 - ...
- **Common insight:** personal tastes are correlated
 - 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

Two Types of Collaborative Filtering

1. Neighborhood Methods

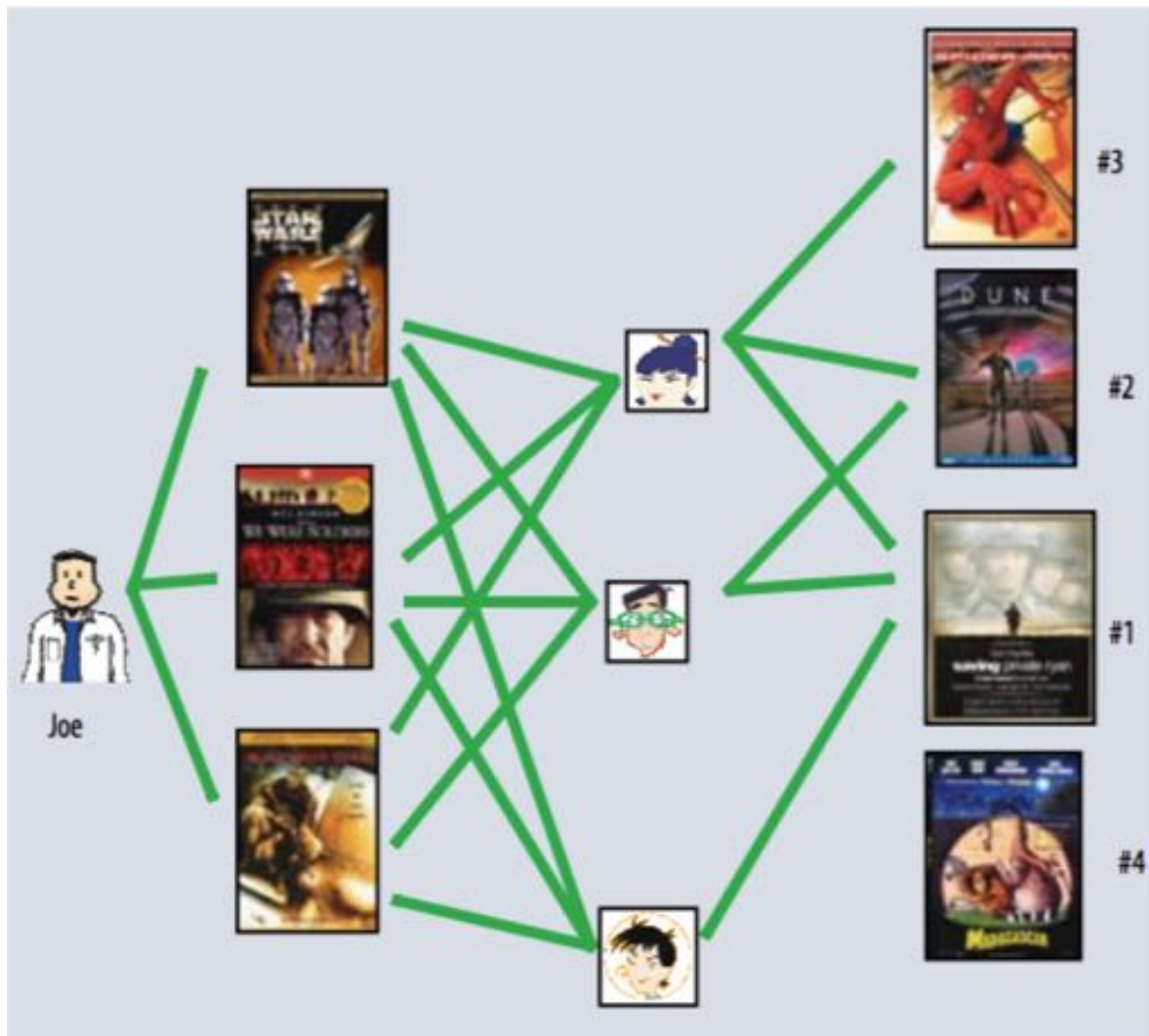


2. Latent Factor Methods



Two Types of Collaborative Filtering

1. Neighborhood Methods



In the figure, assume that a green line indicates the movie was **watched**

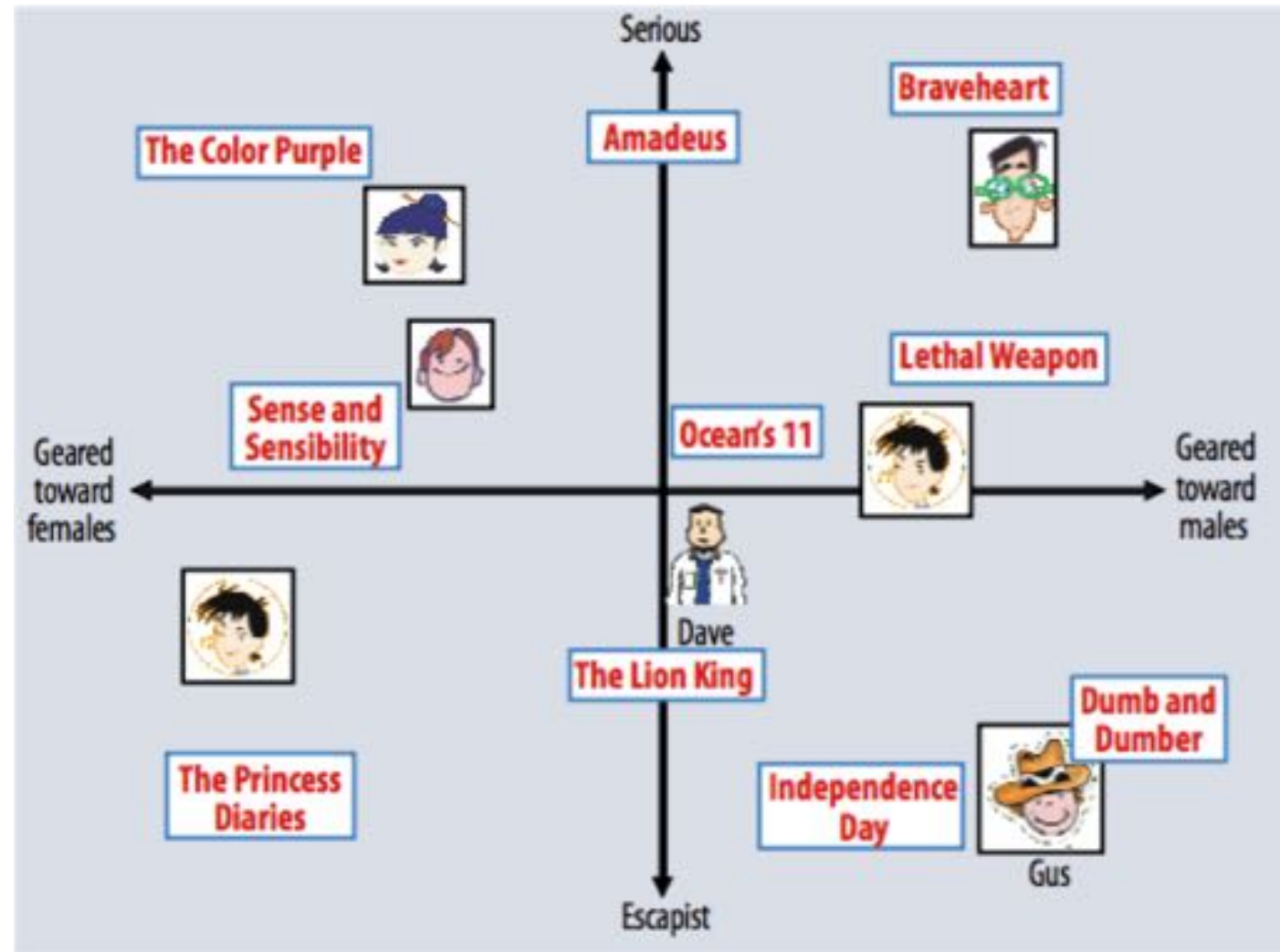
Algorithm:

1. **Find neighbors** based on similarity of movie preferences
2. **Recommend** movies that those neighbors watched

Two Types of Collaborative Filtering

2. Latent Factor Methods

- Assume that both movies and users live in some **low-dimensional space** describing their properties
- **Recommend** a movie based on its **proximity** to the user in the latent space
- **Example Algorithm:** Matrix Factorization



MATRIX FACTORIZATION

Recommending Movies

Question:

Applied to the Netflix Prize problem, which of the following methods *always* requires side information about the users and movies?

Select all that apply

- A. collaborative filtering
- B. latent factor methods
- C. ensemble methods
- D. content filtering
- E. neighborhood methods
- F. recommender systems

Answer:

Matrix Factorization

- Many different ways of factorizing a matrix
- We'll consider three:
 1. Unconstrained Matrix Factorization
 2. Singular Value Decomposition
 3. Non-negative Matrix Factorization
- MF is just another example of a **common recipe**:
 1. define a model
 2. define an objective function
 3. optimize with SGD

Matrix Factorization

Whiteboard

- Background: Low-rank Factorizations
- Residual matrix

Example: MF for Netflix Problem

		NERO	JULIUS CAESAR	CLEOPATRA	SLEEPLESS IN SEATTLE	PRETTY WOMAN	CASABLANCA
HISTORY	1	1	1	1	0	0	0
	2	1	1	1	0	0	0
	3	1	1	1	0	0	0
BOTH	4	1	1	1	1	1	1
	5	-1	-1	-1	1	1	1
ROMANCE	6	-1	-1	1	1	1	1
	7	-1	-1	-1	1	1	1

 \approx

		HISTORY	ROMANCE
1	1	1	0
2	1	1	0
3	1	1	0
4	1	1	1
5	-1	1	1
6	-1	1	1
7	-1	1	1

 \times

		NERO	JULIUS CAESAR	CLEOPATRA	SLEEPLESS IN SEATTLE	PRETTY WOMAN	CASABLANCA
HISTORY	1	1	1	1	0	0	0
ROMANCE	0	0	1	1	1	1	1

 V^T

(a) Example of rank-2 matrix factorization

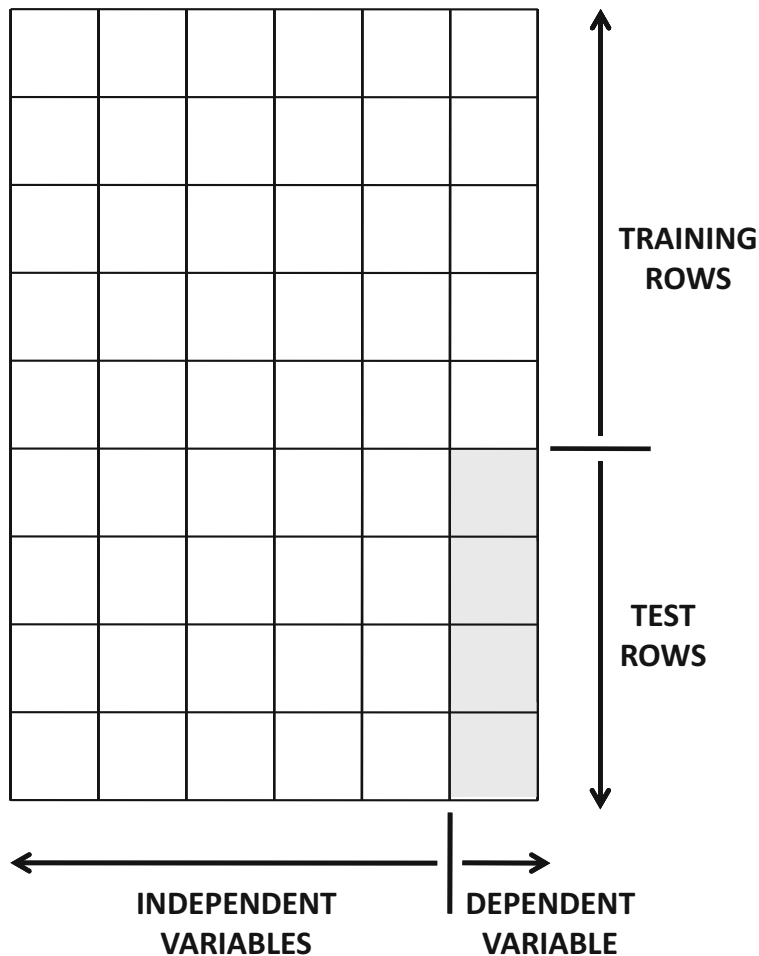
		NERO	JULIUS CAESAR	CLEOPATRA	SLEEPLESS IN SEATTLE	PRETTY WOMAN	CASABLANCA
HISTORY	1	0	0	0	0	0	0
	2	0	0	0	0	0	0
	3	0	0	0	0	0	0
BOTH	4	0	0	-1	0	0	0
	5	0	0	-1	0	0	0
ROMANCE	6	0	0	1	0	0	0
	7	0	0	-1	0	0	0

 R

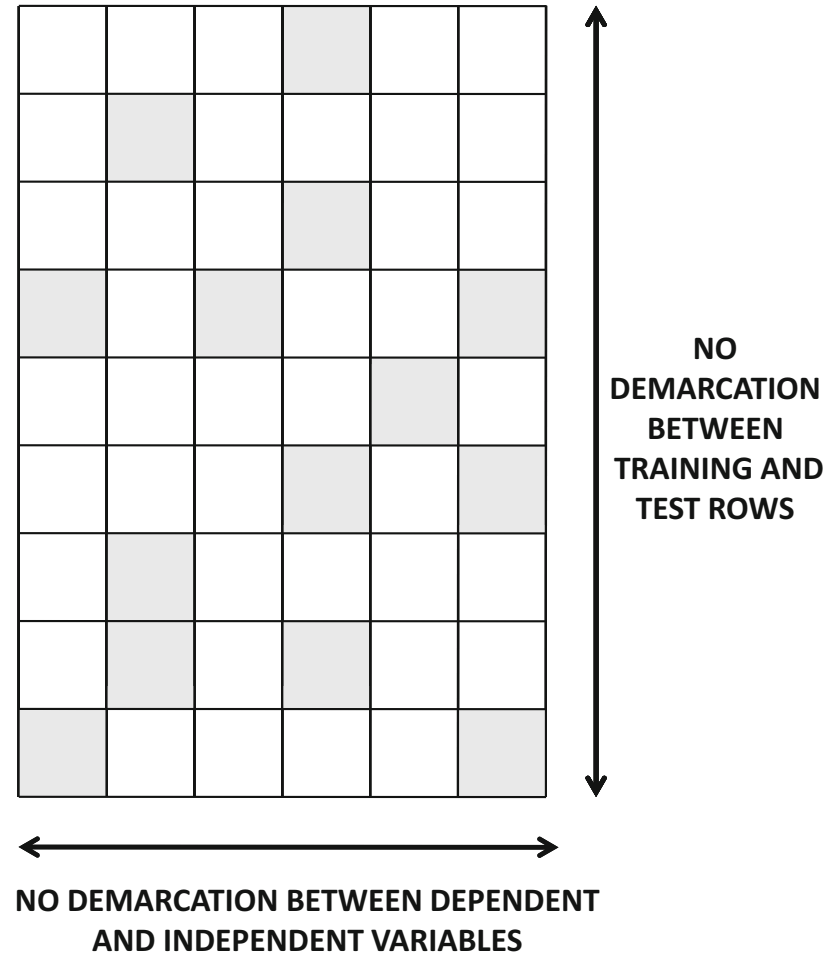
(b) Residual matrix

Regression vs. Collaborative Filtering

Regression



Collaborative Filtering



UNCONSTRAINED MATRIX FACTORIZATION

Unconstrained Matrix Factorization

Whiteboard

- Optimization problem
- SGD
- SGD with Regularization
- Alternating Least Squares
- User/item bias terms (matrix trick)

Unconstrained Matrix Factorization

In-Class Exercise

Derive a block coordinate descent algorithm for the Unconstrained Matrix Factorization problem.

- User vectors:

$$\mathbf{w}_u \in \mathbb{R}^r$$

- Item vectors:

$$\mathbf{h}_i \in \mathbb{R}^r$$

- Rating prediction:

$$v_{ui} = \mathbf{w}_u^T \mathbf{h}_i$$

- Set of non-missing entries

$$\mathcal{Z} = \{(u, i) : v_{ui} \text{ is observed}\}$$

- Objective:

$$\operatorname{argmin}_{\mathbf{w}, \mathbf{h}} \sum_{(u, i) \in \mathcal{Z}} (v_{ui} - \mathbf{w}_u^T \mathbf{h}_i)^2$$

Matrix Factorization (with matrices)

- User vectors:

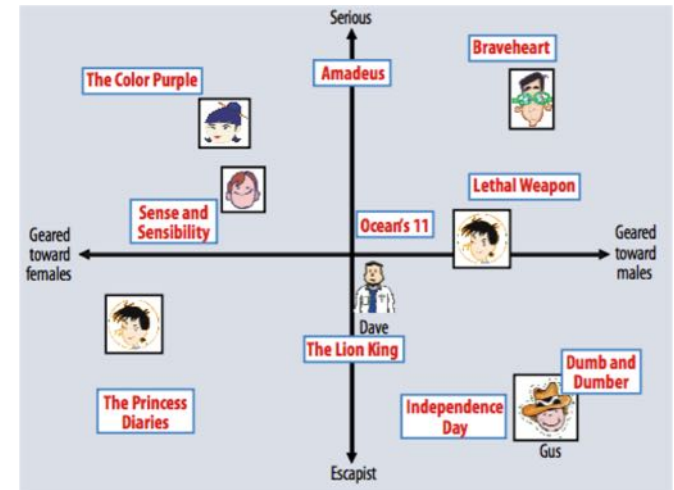
$$(W_{u*})^T \in \mathbb{R}^r$$

- Item vectors:

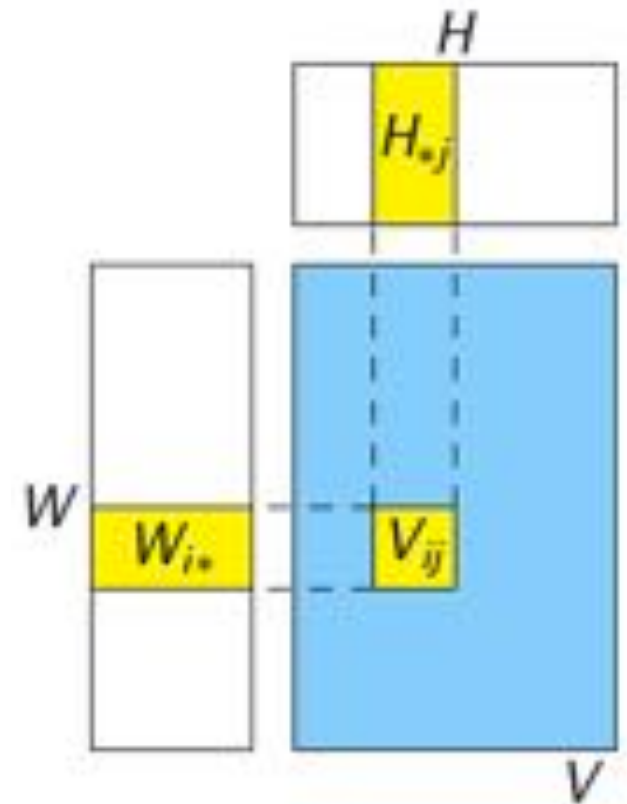
$$H_{*i} \in \mathbb{R}^r$$

- Rating prediction:

$$\begin{aligned} V_{ui} &= W_{u*} H_{*i} \\ &= [WH]_{ui} \end{aligned}$$



Figures from Koren et al. (2009)



Figures from Gemulla et al. (2011)₅₈

Matrix Factorization (with vectors)

- User vectors:

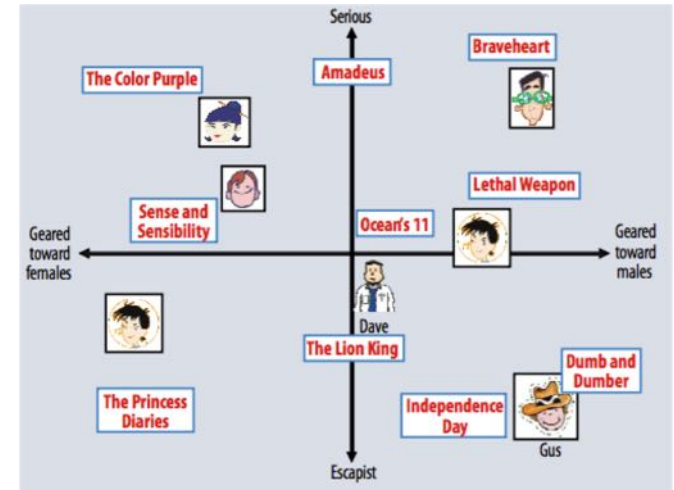
$$\mathbf{w}_u \in \mathbb{R}^r$$

- Item vectors:

$$\mathbf{h}_i \in \mathbb{R}^r$$

- Rating prediction:

$$v_{ui} = \mathbf{w}_u^T \mathbf{h}_i$$



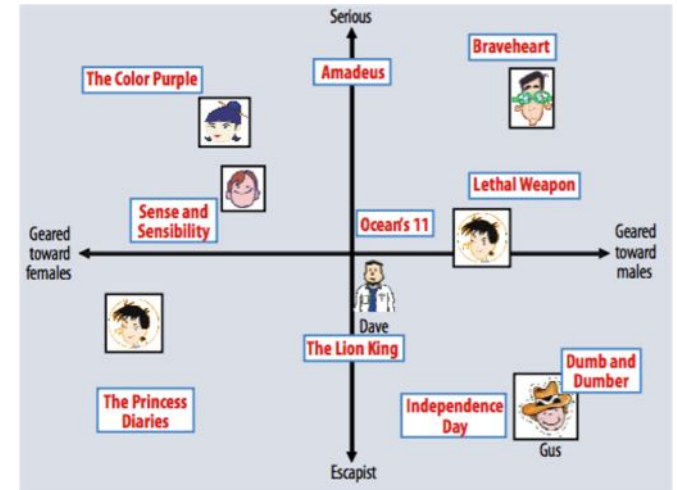
Figures from Koren et al. (2009)

Matrix Factorization (with vectors)

- Set of non-missing entries:
 $\mathcal{Z} = \{(u, i) : v_{ui} \text{ is observed}\}$

- Objective:

$$\operatorname{argmin}_{\mathbf{w}, \mathbf{h}} \sum_{(u, i) \in \mathcal{Z}} (v_{ui} - \mathbf{w}_u^T \mathbf{h}_i)^2$$

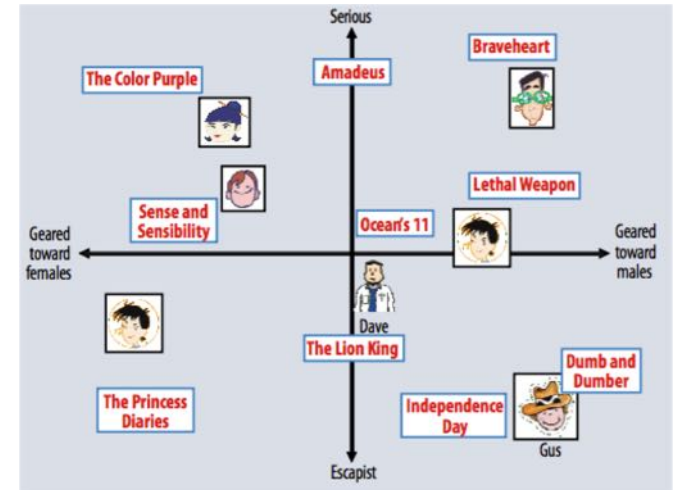


Figures from Koren et al. (2009)

Matrix Factorization (with vectors)

- Regularized Objective:

$$\begin{aligned} \operatorname{argmin}_{\mathbf{w}, \mathbf{h}} \sum_{(u, i) \in \mathcal{Z}} (v_{ui} - \mathbf{w}_u^T \mathbf{h}_i)^2 \\ + \lambda \left(\sum_i \|\mathbf{w}_i\|^2 + \sum_u \|\mathbf{h}_u\|^2 \right) \end{aligned}$$



Figures from Koren et al. (2009)

Matrix Factorization (with vectors)

- Regularized Objective:

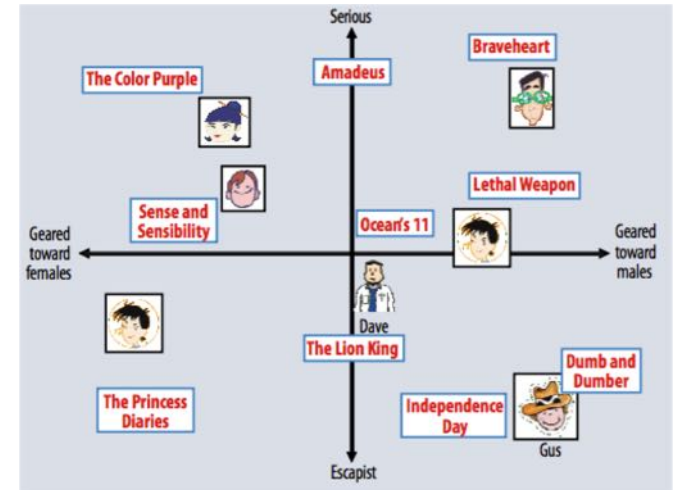
$$\begin{aligned} \operatorname{argmin}_{\mathbf{w}, \mathbf{h}} \sum_{(u,i) \in \mathcal{Z}} (v_{ui} - \mathbf{w}_u^T \mathbf{h}_i)^2 \\ + \lambda \left(\sum_i \|\mathbf{w}_i\|^2 + \sum_u \|\mathbf{h}_u\|^2 \right) \end{aligned}$$

- SGD update for random (u,i):

$$e_{ui} \leftarrow v_{ui} - \mathbf{w}_u^T \mathbf{h}_i$$

$$\mathbf{w}_u \leftarrow \mathbf{w}_u + \gamma(e_{ui} \mathbf{h}_i - \lambda \mathbf{w}_u)$$

$$\mathbf{h}_i \leftarrow \mathbf{h}_i + \gamma(e_{ui} \mathbf{w}_u - \lambda \mathbf{h}_i)$$



Figures from Koren et al. (2009)

Matrix Factorization (with matrices)

- User vectors:

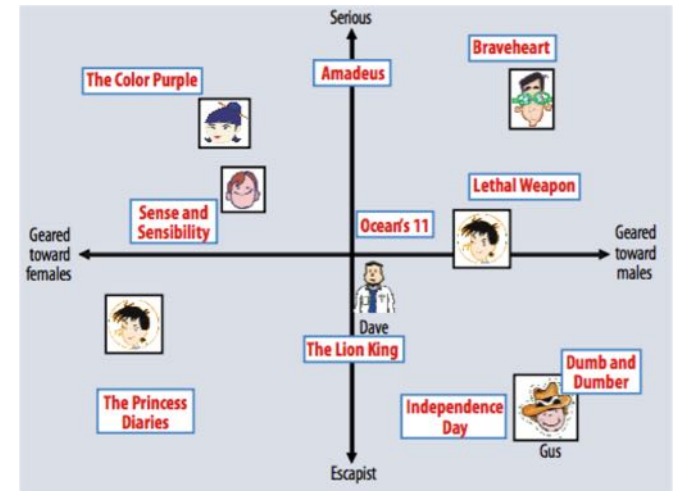
$$(W_{u*})^T \in \mathbb{R}^r$$

- Item vectors:

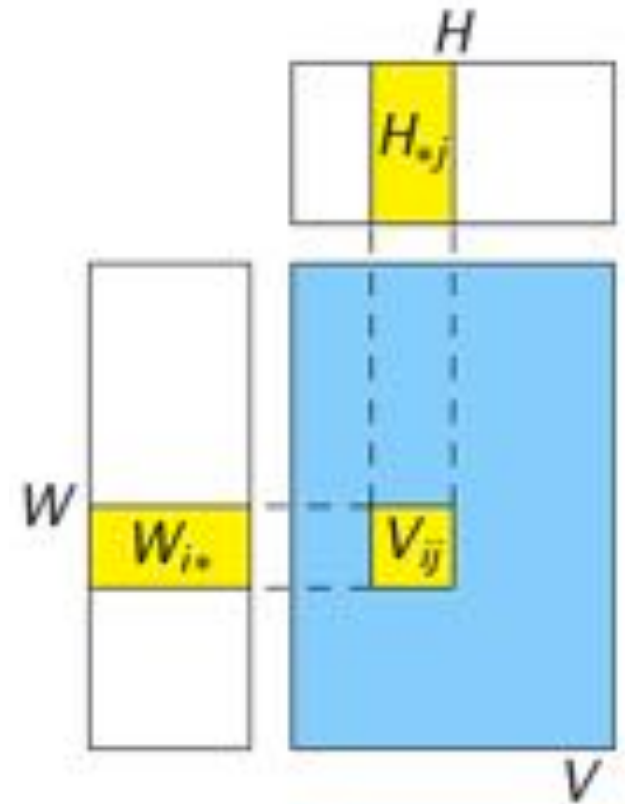
$$H_{*i} \in \mathbb{R}^r$$

- Rating prediction:

$$\begin{aligned} V_{ui} &= W_{u*} H_{*i} \\ &= [WH]_{ui} \end{aligned}$$



Figures from Koren et al. (2009)



Figures from Gemulla et al. (2011)

Matrix Factorization (with matrices)

- SGD

require that the loss can be written as

$$L = \sum_{(i,j) \in Z} l(V_{ij}, W_{i*}, H_{*j})$$

Algorithm 1 SGD for Matrix Factorization

Require: A training set Z , initial values W_0 and H_0

while not converged **do** {step}

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

$W'_{i*} \leftarrow W_{i*} - \epsilon_n N \frac{\partial}{\partial W_{i*}} l(V_{ij}, W_{i*}, H_{*j})$

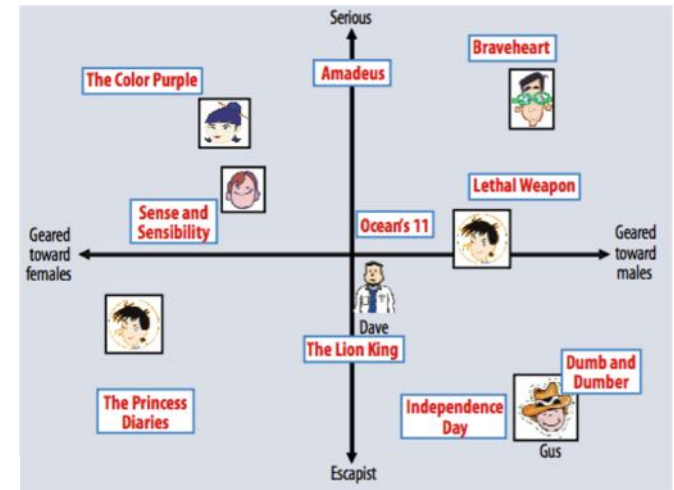
$H_{*j} \leftarrow H_{*j} - \epsilon_n N \frac{\partial}{\partial H_{*j}} l(V_{ij}, W_{i*}, H_{*j})$

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

end while

step size

Figure from Gemulla et al. (2011)



Figures from Koren et al. (2009)

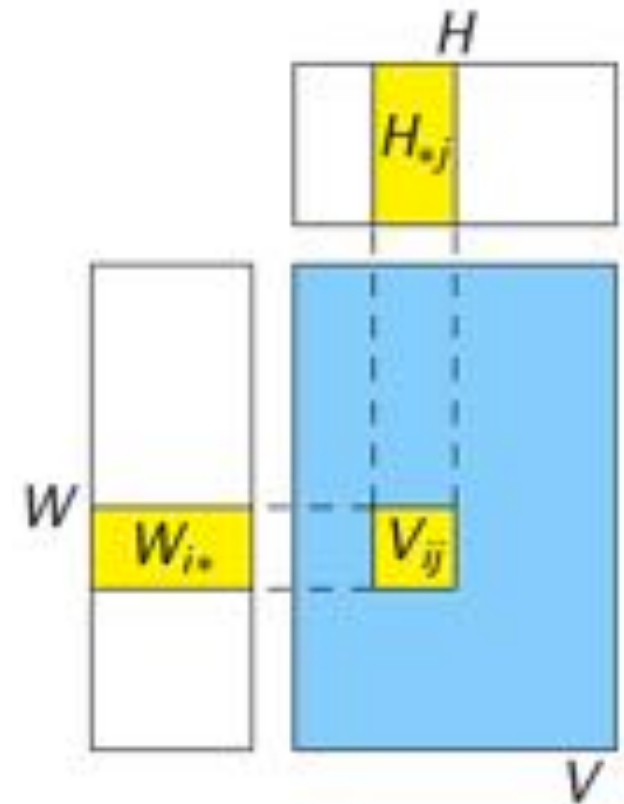


Figure from Gemulla et al. (2011)₆₄

Matrix Factorization

Example Factors

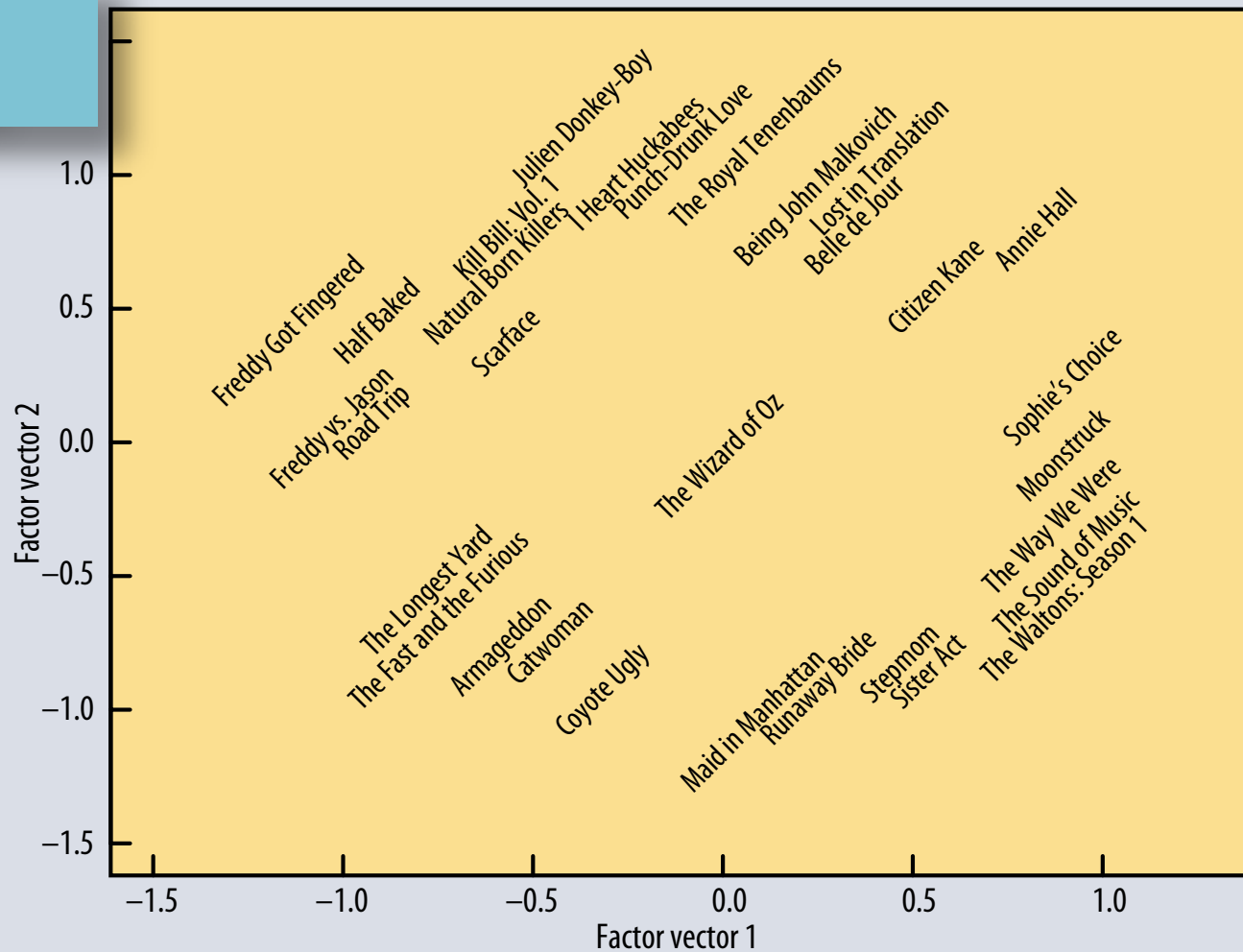
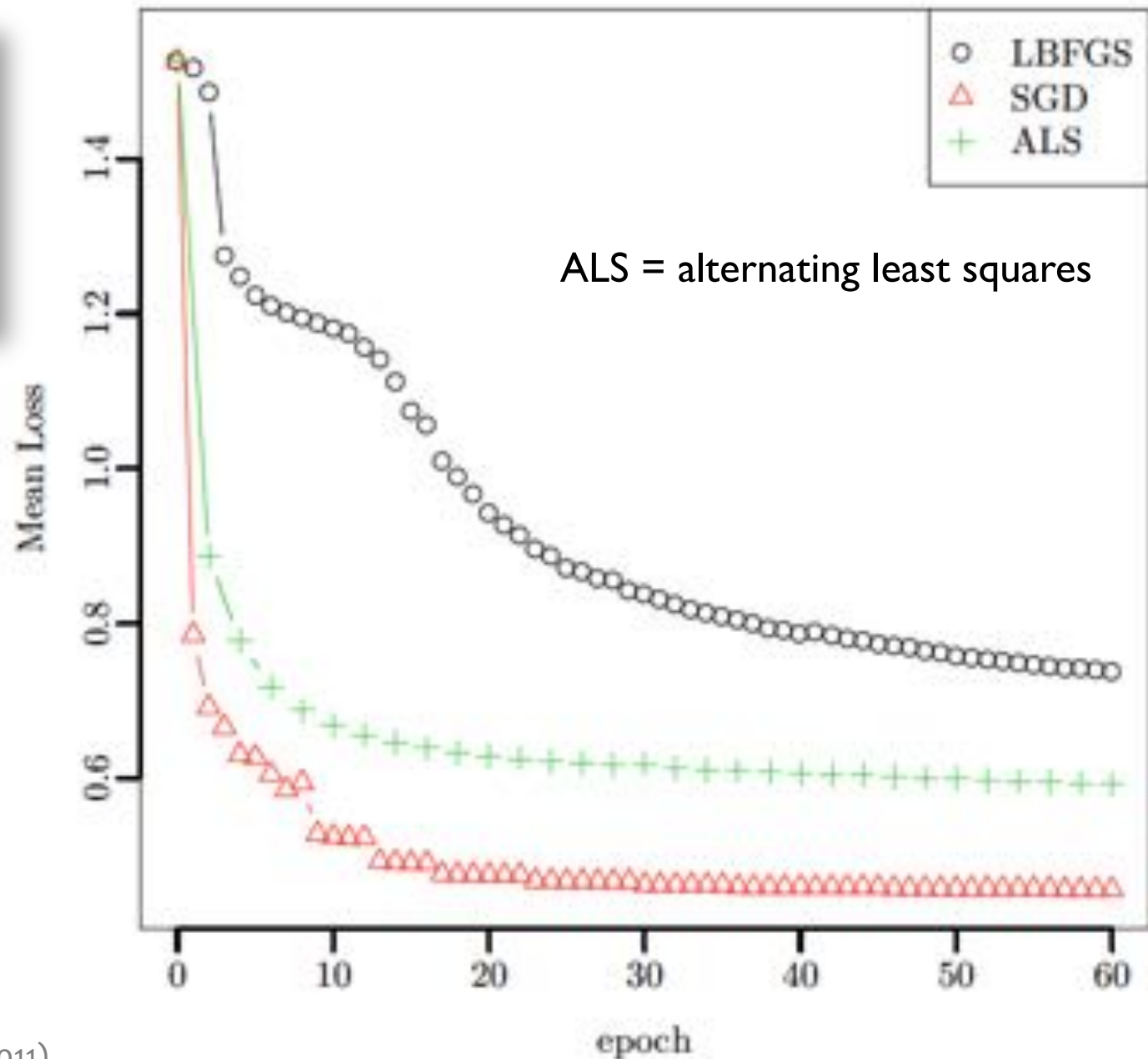


Figure 3. The first two vectors from a matrix decomposition of the Netflix Prize data. Selected movies are placed at the appropriate spot based on their factor vectors in two dimensions. The plot reveals distinct genres, including clusters of movies with strong female leads, fraternity humor, and quirky independent films.

Matrix Factorization

Comparison of Optimization Algorithms



SVD FOR COLLABORATIVE FILTERING

Singular Value Decomposition for Collaborative Filtering

For any arbitrary matrix \mathbf{A} , SVD gives a decomposition:

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$$

where $\mathbf{\Lambda}$ is a diagonal matrix, and \mathbf{U} and \mathbf{V} are orthogonal matrices.

Suppose we have the SVD of our ratings matrix

$$R = Q\Sigma P^T,$$

but then we truncate each of Q , Σ , and P s.t. Q and P have only k columns and Σ is $k \times k$:

$$R \approx Q_k \Sigma_k P_k^T$$

For collaborative filtering, let:

$$U \triangleq Q_k \Sigma_k$$

$$V \triangleq P_k$$

$$\Rightarrow U, V = \operatorname{argmin}_{U, V} \frac{1}{2} \|R - UV^T\|_2^2$$

s.t. columns of U are mutually orthogonal

s.t. columns of V are mutually orthogonal

Theorem: If R fully observed and no regularization, the optimal UV^T from SVD equals the optimal UV^T from Unconstrained MF

NON-NEGATIVE MATRIX FACTORIZATION

Implicit Feedback Datasets

- What information does a five-star rating contain?



- Implicit Feedback Datasets:
 - In many settings, users don't have a way of expressing *dislike* for an item (e.g. can't provide negative ratings)
 - The only mechanism for feedback is to “like” something
- Examples:
 - Facebook has a “Like” button, but no “Dislike” button
 - Google's “+1” button
 - Pinterest pins
 - Purchasing an item on Amazon indicates a preference for it, but there are many reasons you might *not* purchase an item (besides dislike)
 - Search engines collect click data but don't have a clear mechanism for observing dislike of a webpage

Non-negative Matrix Factorization

Constrained Optimization Problem:

$$U, V = \operatorname{argmin}_{U, V} \frac{1}{2} \|R - UV^T\|_2^2$$

$$\text{s.t. } U_{ij} \geq 0$$

$$\text{s.t. } V_{ij} \geq 0$$

Multiplicative Updates: simple iterative algorithm for solving just involves multiplying a few entries together

Summary

- Recommender systems solve many **real-world** (*large-scale) **problems**
- Collaborative filtering by Matrix Factorization (MF) is an **efficient** and **effective** approach
- MF is just another example of a **common recipe**:
 1. define a model
 2. define an objective function
 3. optimize with SGD