

# Recommender Systems

## Harnessing the Power of Personalization

Michael Hahsler

Engineering Management, Information, and Systems (EMIS)  
Intelligent Data Analysis Lab (IDA@SMU)  
Bobby B. Lyle School of Engineering, Southern Methodist University

Southwest Airlines EDGe Analyst Community Meeting  
November 18, 2015



SMU | BOBBY B. LYLE  
SCHOOL OF ENGINEERING

**Mission:** At IDA we create novel techniques inspired by knowledge discovery, data mining, machine learning, artificial intelligence and statistical analysis to work with data from various sources. We currently focus on:

- Order modeling for massive data streams with applications in meteorology (hurricane intensity prediction) and personalized medicine (efficient classification and analysis of metagenomic data)
- Visual analytics using optimized reordering
- Simulation data analytics
- Recommender systems

**Team:** 3 faculty, 7 students, 2 collaborators

**Director:** M. Hahsler <mhahsler@lyle.smu.edu>

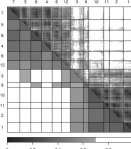
Supported by



National Human  
Genome Research  
Institute



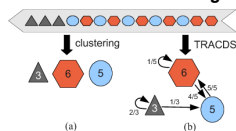
## Visual Analytics



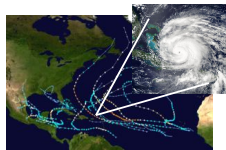
## Simulation Data Analytics



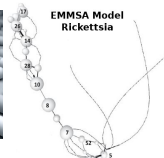
## Order Modeling Data Stream Clustering



## Meteorology



## Personalized Medicine



# Reproducible Research Using R



R has been consistently voted one of the most important tools for data mining and analytics, and being able to use R is one of the highest paying analytics skills.

Our team has developed and maintains several popular R packages:

## Association Rule Mining

- **arules**: Mining association rules and frequent itemsets.
- **arulesViz**: Visualizing association rules based on package arules.
- **arulesSequences**: Mine frequent sequences.

## Combinatorial Optimization

- **seriation**: Seriation/sequencing techniques to reorder matrices and dendrograms.
- **TSP**: Infrastructure and algorithms for the traveling salesperson problem.
- **DBSCAN**: Several density-based algorithms for spatial data.
- **QAP**: Heuristics for the Quadratic Assignment Problem (QAP).

## Data Stream Mining

- **stream**: Infrastructure for data stream mining.

## Recommender Systems

- **recommenderlab**: Infrastructure to test and develop recommender algorithms.

<http://michael.hahsler.net/#Software>

# A Message From the Department Chair

The Engineering Management, Information, and Systems Department program includes:

- Management Science
- Operations Research
- Analytics

We are looking for topics for [undergraduate senior design projects](#) in any of these areas.

Please contact Sila Centinkaya ([sila@lyle.smu.edu](mailto:sila@lyle.smu.edu)), Chair EMIS, with inquiries.

# Table of Contents

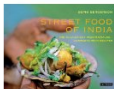
- 1 Motivation
- 2 Content-based Approach
- 3 Collaborative Filtering (CF)
  - Memory-based CF
  - Model-based CF
- 4 Further Considerations
  - Strategies for the Cold Start Problem
  - Recommender Systems and Revenue Management
- 5 Implementation
  - Open-Source Tools
  - An Example using recommenderlab
  - Deployment



## Kristina, Welcome to Your Amazon.com

## Today's Recommendations For You

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



[Street Food of India: The 50...](#)  
(Hardcover) by Sephil Bergerson

★★★★☆ (4) \$19.17

[Fix this recommendation](#)



[Lavazza Tierra! 100% Arabica Whole Bean Espresso...](#)

★★★★☆ (38) \$34.41

[Fix this recommendation](#)



[Entourage: The Complete Fou... DVD ~ Adrian Grenier](#)

★★★★☆ (44) \$16.49

[Fix this recommendation](#)

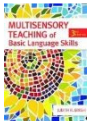
## New For You®



[The Race \(Isaac Bell\) Clive Cussler, Justin Scott Hardcover](#)

~~\$27.95~~ \$14.97

[Fix this recommendation](#)



[Multisensory Teaching of Basic... Language Skills Judith R. Birsh, Sally E. Shaywitz](#)

Hardcover  
~~\$79.95~~ \$44.99

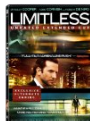
[Fix this recommendation](#)



[Kill Shot \(Mitch Rapp\) Vince Flynn Hardcover](#)

~~\$27.99~~ \$16.62

[Fix this recommendation](#)



[Limitless \(Unrated Extended Cut\) Bradley Cooper, Anna Friel, Abbie... DVD](#)

~~\$29.99~~ \$15.19

[Fix this recommendation](#)

Suggestions (1141)

Suggestions by Genre ▾

Rate Movies

Rate Genres

Movies You've Rated (262)

## Movies You'll Love

Suggestions based on your ratings

 You have 1141  
 Suggestions  
 from 262 ratings.

### New Suggestions for you

Based on your recent ratings



Play + All


 Not Interested

 Dexter  
 1 (4)  
 Beco  
 enjo  
 Lost  
 Battl  
 Gala  
 Seas  
 Rom

### The Fugitive (1993)

Wrongfully convicted of murdering his wife, Dr. Richard Kimble (Harrison Ford) escapes custody after a ferocious train accident (one of the most thrilling wrecks ever filmed). While Kimble tries to find the true murderer, gung-ho U.S. Marshal Samuel Gerard (Tommy Lee Jones, in an Oscar-winning performance) is hot on Kimble's trail, pulling out all stops to put him back behind bars.

**Starring:** Harrison Ford, Tommy Lee Jones

**Director:** Andrew Davis

**Genre:** Action & Adventure

**MPAA:** PG-13


4.7 Our best guess for Michael



4.1 Customer Average

Recommended based on 8 ratings

 nos:  
 e  
 u  
 ther  
 ther,  
 nder:


Add


 Not Interested

### The Fugitive

Because you enjoyed:

[Patriot Games](#)  
[Indiana Jones and the Last Crusade](#)  
[Die Hard](#)

★★★★★ SCI-FI &amp;



Incredib

[See all 26 >](#)


Spacehunter



RoboCop:

# PANDORA®

internet radio

[register](#) | [sign in](#)

Help

Register for FREE to save your stations and access Pandora anytime, anywhere.

[register now](#)

share



[Create a New Station...](#)

Your Stations

**Lady Gaga Radio**

[add variety...](#)

[options](#) ▾

QuickMix ▾

Alejandro

buy

by: Lady Gaga  
on: The Fame M...



menu



Evacuate The Dancefloor

buy

by: Cascada  
on: Evacuate Th...



menu



Toxic

buy

by: Britney Spears  
on: Greatest Hits...



menu





 Search Flights → Select Flights → Price → Purchase → Confirmed

## Select Departing Flight:

### Dallas (Love Field), TX to Washington (Reagan National), DC

 Modify Search  Round Trip  One-Way [Additional Search Options](#)

From:  To:  [+ Add another flight](#)

 **First 2 Bags Fly Free®.** Weight, size & excess limits apply.  Gov't taxes & fees now included

NOV	NOV	NOV	NOV	NOV	NOV	NOV	NOV	NOV	NOV	
13 FRI	14 SAT	15 SUN	16 MON	17 TUE	<b>18</b> WED	19 THU	20 FRI	21 SAT	22 SUN	23 MON

 **Flexible dates?**  
Search the low fare calendar.

#### Filter My Results

Nonstop  Direct (No plane change, with stops)

#### Show fares in

All fares are rounded up to the nearest dollar.

Depart	Arrive	Flight #	Routing	Travel Time	Business Select \$496 - \$505	Anytime \$474 - \$483	Wanna Get Away \$197 - \$369
6:00 AM	12:10 PM	1628 1001	1 stop Change Planes HOU	5h 10m	Sold Out	<input type="radio"/> \$483	<input type="radio"/> \$205
6:00 AM	11:35 AM	310	1 stop Change Planes	4h 35m	 \$505	<input type="radio"/> \$483	<input type="radio"/> \$205

#### Quick Air Links

- [Check In](#)
- [Change Flight](#)
- [Check Flight Status](#)

Account Login [Enroll Now!](#)


#### Username


#### Password

Remember Me

Need help logging in?

 [Manage Travel](#)

 [Shopping Cart](#)

 [Rapid Rewards](#)

# Dallas (Love Field), TX to Washington (Reagan National), DC

**Air** Total Price: **\$1008.96**

## ITINERARY

Travel Date	Flight Segments	Flight Summary
<b>DEPART</b> <b>NOV 18</b> <b>WED</b>	<b>06:00 AM</b> Depart <b>Dallas (Love Field), TX (DAL)</b> on Southwest Airlines <b>07:35 AM</b> Arrive in St. Louis, MO (STL) <b>08:45 AM</b> Change  to Southwest St. Louis, MO (STL) <b>11:35 AM</b> Arrive in <b>Washington (Reagan National), DC (DCA)</b>	<b>Flight #310</b> Southwest WiFi available <b>Wednesday, November 18, 2015</b>
<b>RETURN</b> <b>NOV 21</b> <b>SAT</b>	<b>06:00 AM</b> Depart <b>Washington (Reagan National), DC (DCA)</b> on Airlines <b>08:00 AM</b> Arrive in Atlanta, GA (ATL) <b>09:40 AM</b> Change  to Southwest Atlanta, GA (ATL) <b>11:00 AM</b> Arrive in <b>Dallas (Love Field) (DAL)</b>	

### What you need to know to travel:

- Check-in: Be sure to arrive at the departure time. Otherwise, your reserved spa

### Quick Air Links

- Check In
- Change Flight
- Check Flight Status

Account Login [Enroll Now!](#)

### Username

### Password

Remember Me

[Need help logging in?](#)

## Add a Hotel

- We'll keep an eye on your cart for you while you shop. Products not confirmed until purchase.



The George, a Kimpton Hotel

\$309/night



[View Details](#)

Search for hotels in Washington (11/18/2015 - 11/21/2015)

Close To (optional)

within

Show Only (optional)

## Add a Car

- We'll keep an eye on your cart for you while you shop. Products not confirmed until purchase.

Earn up to **2,400** Plus save up to **30%** Rapid Rewards<sup>®</sup> points. [Let's go!](#)

Type in any city or airport in the U.S., [Canada](#) or [Mexico](#)

Pickup Location

Pickup Date

Dropoff Date

Vehicle Type (optional)

Which Company? (optional)

Can I recommend anything else?

<http://www.kdnuggets.com/2015/10/big-data-recommendation-systems-change-lives.html/2>

# Recommender Systems

## Original Definition

Recommender systems apply statistical and knowledge discovery techniques to the problem of making product recommendations.

Sarwar *et al.* (2000)

## Advantages of recommender systems (Schafer *et al.*, 2001):

- Improve conversion rate: Help customers find a product she/he wants to buy.
- Cross-selling: Suggest additional products.
- Up-selling: Suggest premium products.
- Improve customer loyalty: Create a value-added relationship.



## A Significant Number Of Travelers Can Be Tempted With Up-Sell/Cross-Sell Offers

Percent of US travelers who will consider paying a reasonable premium for the following:



Receive better service: **48%**



Save time/reduce hassle: **47%**



Enjoy more comfort: **45%**

Base US online travelers

Source: Hudson Crossing's US Travel Online Study, Q1 2013

# A More General View of Recommender Systems

A recommender system is a **fully automatic system** to provide (near) **personalized decision support** given **limited information** while optimizing a set of potentially conflicting **objective functions**.

# A More General View of Recommender Systems

A recommender system is a **fully automatic system** to provide (near) **personalized decision support** given **limited information** while optimizing a set of potentially conflicting **objective functions**.

## Important aspects:

- Personalization
- Available information
- Incentive structure
- Trust
- Quality of recommendations
- Speed

# Common Approaches

- **Content-based filtering:** Consumer preferences for product attributes.
- **Collaborative filtering:** Mimics word-of-mouth based on analysis of rating/usage/sales data from many users.

(Ansari *et al.*, 2000)

- **Hybrid recommender systems:** Incorporate content, collaborative and expert information.



# Table of Contents

- 1 Motivation
- 2 Content-based Approach
- 3 Collaborative Filtering (CF)
  - Memory-based CF
  - Model-based CF
- 4 Further Considerations
  - Strategies for the Cold Start Problem
  - Recommender Systems and Revenue Management
- 5 Implementation
  - Open-Source Tools
  - An Example using recommenderlab
  - Deployment

# Content-based Approach

The screenshot shows the IMDb page for the movie 'The Social Network' (2010). The page features a search bar at the top, navigation tabs for Movies, TV, News, Videos, Community, and IMDb. The movie's poster is on the left, with the text 'YOU DON'T GET TO 500 MILLION FRIENDS WITHOUT MAKING A FEW ENEMIES'. The main content area includes the movie title, a '95' rating badge, a 'PG-13' rating, a 120-minute runtime, and genres 'Biography' and 'Drama'. It also shows the release date '1 October 2010 (USA)'. A star rating section displays a yellow star with '8.1' and a mouse cursor, with 'Your rating' and a star bar. Below this, it shows 'Ratings: 8.1/10 from 141,802 users', 'Metascore: 95/100', and 'Reviews: 515 user | 459 critic | 42 from Metacritic.com'. A synopsis follows: 'A chronicle of the founding of Facebook, the social-networking Web site.' Credits for Director (David Fincher), Writers (Aaron Sorkin, Ben Mezrich), and Stars (Jesse Eisenberg, Andrew Garfield, Justin Timberlake) are listed.

- 1 Analyze the objects (documents, video, music, etc.) and extract attributes/features (e.g., words, phrases, actors, genre).
- 2 Recommend objects with similar attributes to an object the user likes.



**Lady Gaga**

[Just Dance \(Remix Single\)](#)

Just Dance (Redone Remix F. Kardinal Offishall)

Play Sample

**PANDORA®**

---

#### Features Of This Track

electronica roots  
trip hop roots  
r&b influences  
funk influences  
beats made for dancing  
unsyncopated ensemble rhythms  
straight drum beats  
a female vocal  
clear pronunciation  
a rhythmic intro  
use of modal harmonies  
the use of chordal patterning  
melodic part writing  
use of strings  
subtle use of arpeggiated synths  
affected synths

**Create A Station**

**Bookmark This Track**

**Buy on iTunes**

**Buy CD From Amazon**

**Buy From Amazon MP3**

“The [Music Genome Project](#) is an effort to capture the essence of music at the fundamental level using almost 400 attributes to describe songs and a complex mathematical algorithm to organize them.”

[http://en.wikipedia.org/wiki/Music\\_Genome\\_Project](http://en.wikipedia.org/wiki/Music_Genome_Project)

**KAYAK** HOTELS FLIGHTS CARS PACKAGES TRIPS 1

**DFW ↔ VIE** | Dec 22 → Jan 6 | Economy | 3  
 672 of 1059 flights | Tuesday | Wednesday | cabin | travelers

Sort by: price (low to high) ▾

**See Deal** **JustFly, Up To 80% Off Flights**  
 Save big on flights to Vienna with up to 80% Off Flights. For A

[Select](#) **See Deal** 2+ stops [Select](#)

[www.justfly.com](http://www.justfly.com)

**Advice: BUY Confidence: 80%**  
 Prices may rise within 7 days ⓘ

[Create a price alert](#)

**Stops**

- nonstop
- 1 stop \$1502
- 2+ stops \$1363

**Times**

Content can be dynamic...

An issue with content based filtering?

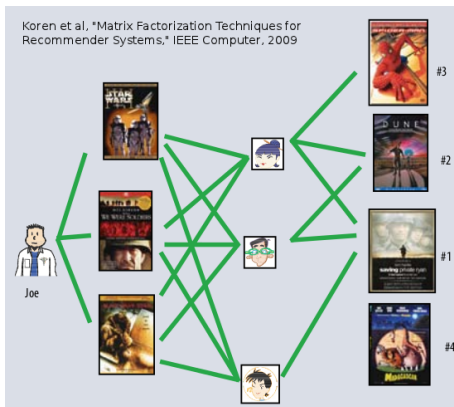
An issue with content based filtering?

Missing diversity!

# Table of Contents

- 1 Motivation
- 2 Content-based Approach
- 3 Collaborative Filtering (CF)
  - Memory-based CF
  - Model-based CF
- 4 Further Considerations
  - Strategies for the Cold Start Problem
  - Recommender Systems and Revenue Management
- 5 Implementation
  - Open-Source Tools
  - An Example using recommenderlab
  - Deployment

# Collaborative Filtering (CF)



Make automatic predictions (filtering) about the interests of a user by collecting preferences or taste information from many other users (collaboration).

**Assumption:** those who agreed in the past tend to agree again in the future.



# Data Collection

NETFLIX

Michael Hahsler | Your Account | Buy / Redeem Gift | Help

Browse DVDs | Watch Instantly | Your Queue | Movies You'll ♥ | Friends & Community | DVD Sale \$5.99

Movies, actors, directors, genres Search

Suggestions (1141) | Suggestions by Genre | Rate Movies | Rate Genres | Movies You've Rated (262)

## Rate Movies

You have 1141 Suggestions from 262 ratings.

**Rate More Movies in All Genres!**

Keep rating movies to get recommendations. You can rate movies you've seen in the theater as well as movies you've added from Netflix. Click the star that matches your opinion.

The Day After Tomorrow	National Treasure	Miss Congeniality	Pearl Harbor	The Longest Yard
Add	Add	Add	Add	Add
★ ★ ★ ★ ★ Not Interested	★ ★ ★ ★ ★ Not Interested	★ ★ ★ ★ ★ Not Interested	★ ★ ★ ★ ★ Not Interested	★ ★ ★ ★ ★ Not Interested
No Opinion	No Opinion	No Opinion	No Opinion	No Opinion

Con Air	Coyote Ugly	Annie Hall	Monster-in-Law	Mr. Deeds
Add	Add	Add	Add	Add
★ ★ ★ ★ ★ Not Interested	★ ★ ★ ★ ★ Not Interested	★ ★ ★ ★ ★ Not Interested	★ ★ ★ ★ ★ Not Interested	★ ★ ★ ★ ★ Not Interested
No Opinion	No Opinion	No Opinion	No Opinion	No Opinion

**Browse**

Favorite Genres: (Edit)

All Favorites

- Action & Adventure
- Drama
- Sci-Fi & Fantasy
- Comedy

Other Genres:

- All Genres
- Blu-ray
- Children & Family
- Classics
- Documentary
- Faith & Spirituality
- Foreign

## • Data sources:

- ▶ **Explicit:** ask the user for ratings, rankings, list of favorites, etc.
- ▶ **Observed behavior (Implicit):** clicks, page impressions, purchase, uses, downloads, posts, tweets, etc.

What is the incentive structure?

# Output of a Recommender System

The screenshot shows the Netflix interface with a red header and a yellow navigation bar. The main content area is titled "Movies You'll Love" and features a "New Suggestions for" section. A movie card for "The Fugitive (1993)" is highlighted with a red border. The card includes a description, a synopsis, and a predicted rating of 4.7 stars. Below the card, there are links for "Play" and "Add", and a "Not Interested" button. The background shows other movie suggestions and a "You have 1141 Suggestions from 262 ratings." notification.

**NETFLIX**

Suggestions (1141) | Suggestions by Genre | Rate Movies | Rate Genres | Movies You've Rated (262)

### Movies You'll Love

Suggestions based on your ratings

You have 1141 Suggestions from 262 ratings.

#### New Suggestions for

Based on your recent ratings

**The Fugitive (1993)**

Wrongfully convicted of murdering his wife, Dr. Richard Kimble (Harrison Ford) escapes custody after a ferocious train accident (one of the most thrilling wrecks ever filmed). While Kimble tries to find the true murderer, gung-ho U.S. Marshal Samuel Gerard (Tommy Lee Jones, in an Oscar-winning performance) is hot on Kimble's trail, pulling out all stops to put him back behind bars.

**Starring:** Harrison Ford, Tommy Lee Jones  
**Director:** Andrew Davis  
**Genre:** Action & Adventure  
**MPAA:** PG-13

★★★★★ 4.7 Our best guess for Michael  
★☆☆☆☆ 4.1 Customer Average

[Play](#) [+ All](#) [Not Interested](#)

[The Fugitive](#)

Because you enjoyed:  
[Patriot Games](#)  
[Indiana Jones and the Last Crusade](#)  
[Die Hard](#)

[Add](#) [Not Interested](#)

Recommended based on 8 ratings

[See all 26 >](#)

[Spacehunter](#) [ROBOCOP](#) [RoboCop](#)

- Predicted rating of unrated movies (Breese *et al.*, 1998)
- A top- $N$  list of unrated (unknown) movies ordered by predicted rating/score (Deshpande and Karypis, 2004)

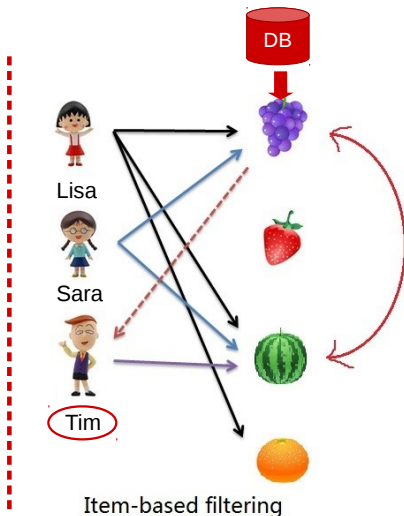
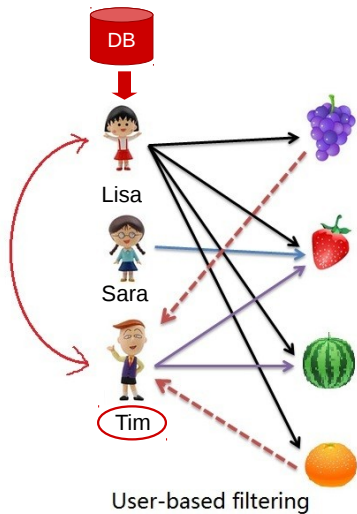
# Types of CF Algorithms

- **Memory-based:** Find similar users (user-based CF) or items (item-based CF) to predict missing ratings.
- **Model-based:** Build a model from the rating data (clustering, latent structure, etc.) and then use this model to predict missing ratings.

# Table of Contents

- 1 Motivation
- 2 Content-based Approach
- 3 Collaborative Filtering (CF)
  - Memory-based CF
  - Model-based CF
- 4 Further Considerations
  - Strategies for the Cold Start Problem
  - Recommender Systems and Revenue Management
- 5 Implementation
  - Open-Source Tools
  - An Example using recommenderlab
  - Deployment

# User-based vs. Item-based CF



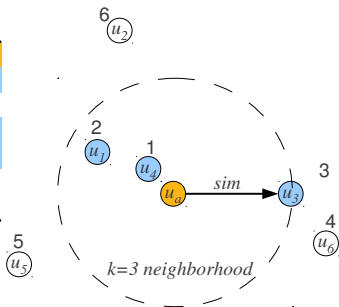
Source: <http://cuihelei.blogspot.com/2012/09/the-difference-among-three.html>

# User-based CF

Produce recommendations based on the preferences of similar users (Goldberg *et al.*, 1992; Resnick *et al.*, 1994; Mild and Reutterer, 2001).

	$i_1$	$i_2$	$i_3$	$i_4$	$i_5$	$i_6$
$u_a$	?	?	4.0	3.0	?	1.0
$u_1$	?	4.0	4.0	2.0	1.0	2.0
$u_2$	3.0	?	?	?	5.0	1.0
$u_3$	3.0	?	?	3.0	2.0	2.0
$u_4$	4.0	?	?	2.0	1.0	1.0
$u_5$	1.0	1.0	?	?	?	?
$u_6$	?	1.0	?	?	1.0	1.0
	3.5	4.0		1.3		

Recommendations:  $i_2, i_1$



- 1 Find  $k$  nearest neighbors for the user in the user-item matrix.
- 2 Generate recommendation based on the items liked by the  $k$  nearest neighbors. E.g., average ratings or use a weighting scheme.

## User-based CF II

- Pearson correlation coefficient:

$$\text{sim}_{\text{Pearson}}(\mathbf{x}, \mathbf{y}) = \frac{\sum_{i \in I} x_i y_i - I \bar{x} \bar{y}}{(I-1) s_x s_y}$$

- Cosine similarity:

$$\text{sim}_{\text{Cosine}}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2}$$

- Jaccard index (only binary data):

$$\text{sim}_{\text{Jaccard}}(X, Y) = \frac{|X \cap Y|}{|X \cup Y|}$$

where  $\mathbf{x} = b_{u_x, \cdot}$  and  $\mathbf{y} = b_{u_y, \cdot}$  represent the user's profile vectors and  $X$  and  $Y$  are the sets of the items with a 1 in the respective profile.

### Problem

Memory-based. Expensive online similarity computation.

# Item-based CF

Produce recommendations based on the relationship between items in the user-item matrix (Kitts *et al.*, 2000; Sarwar *et al.*, 2001)

S	$i_1$	$i_2$	$i_3$	$i_4$	$i_5$	$i_6$	$i_7$	$i_8$	$k=3$
$i_1$	-	0.1	0	<b>0.3</b>	<b>0.2</b>	<b>0.4</b>	0	0.1	$u_a = \{i_1, i_5, i_8\}$
$i_2$	0.1	-	<b>0.8</b>	<b>0.9</b>	0	<b>0.2</b>	0.1	0	$r_{ua} = \{2, ?, ?, ?, 4, ?, ?, 5\}$
$i_3$	0	<b>0.8</b>	-	0	<b>0.4</b>	0.1	0.3	<b>0.5</b>	
$i_4$	<b>0.3</b>	<b>0.9</b>	0	-	0	<b>0.3</b>	0	0.1	
$i_5$	<b>0.2</b>	0	<b>0.7</b>	0	-	<b>0.2</b>	0.1	0	
$i_6$	<b>0.4</b>	<b>0.2</b>	0.1	<b>0.3</b>	0.1	-	0	0.1	
$i_7$	0	<b>0.1</b>	<b>0.3</b>	0	0	0	-	0	
$i_8$	<b>0.1</b>	0	<b>0.9</b>	<b>0.1</b>	0	0.1	0	-	
	-	0	4.56	2.75	-	2.67	0	-	Recommendation: $i_3$

- 1 Calculate similarities between items and keep for each item only the values for the  $k$  most similar items.
- 2 Use the similarities to calculate a weighted sum of the user's ratings for related items.

$$\hat{r}_{ui} = \sum_{j \in s_i} s_{ij} r_{uj} / \sum_{j \in s_i} |s_{ij}|$$

Regression can also be used to create the prediction.



# Item-based CF II

## Similarity measures:

- Pearson correlation coefficient, cosine similarity, Jaccard index
- Conditional probability-based similarity (Deshpande and Karypis, 2004):

$$\text{sim}_{\text{Conditional}}(x, y) = \frac{\text{Freq}(xy)}{\text{Freq}(x)} = \hat{P}(y|x)$$

where  $x$  and  $y$  are two items,  $\text{Freq}(\cdot)$  is the number of users with the given item in their profile.

## Properties

- Model (reduced similarity matrix) is relatively small ( $N \times k$ ) and can be fully precomputed.
- Item-based CF was reported to only produce slightly inferior results compared to user-based CF (Deshpande and Karypis, 2004).
- Higher order models which take the joint distribution of sets of items into account are possible (Deshpande and Karypis, 2004).
- Successful application in large scale systems (e.g., Amazon.com)

# Table of Contents

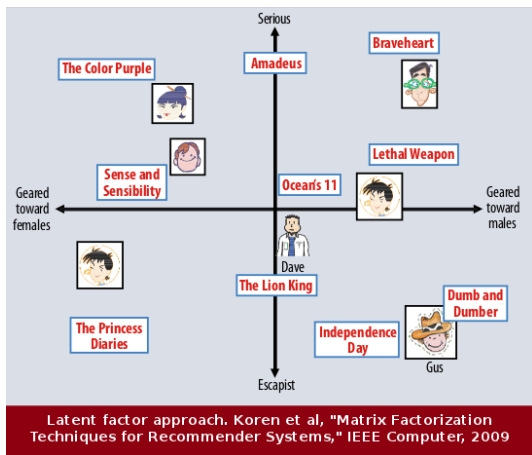
- 1 Motivation
- 2 Content-based Approach
- 3 Collaborative Filtering (CF)**
  - Memory-based CF
  - Model-based CF**
- 4 Further Considerations
  - Strategies for the Cold Start Problem
  - Recommender Systems and Revenue Management
- 5 Implementation
  - Open-Source Tools
  - An Example using recommenderlab
  - Deployment

# Different Model-based CF Techniques

There are many techniques:

- **Cluster users** (i.e., customer segmentation) and then recommend items the users in the cluster closest to the active user like.
- Mine **association rules** (if-then rules) and then use the rules to recommend items.
- Define a null-model (a stochastic process which models usage of independent items) and then find **significant deviation from the null-model**.
- **Learning to rank**: Logistic regression, neural networks (deep learning) and many other machine learning methods.
- Learn a **latent factor model** from the data and then use the discovered factors to find items with high expected ratings.

# Latent Factor Approach



Latent semantic indexing (LSI) developed by the IR community (late 80s) addresses sparsity, scalability and can handle synonyms  
⇒ Dimensionality reduction.

# Matrix Factorization

Given a user-item (rating) matrix  $M = (r_{ui})$ , map users and items on a joint latent factor space of dimensionality  $k$ .

- Each item  $i$  is modeled by a vector  $q_i \in \mathbb{R}^k$ .
- Each user  $u$  is modeled by a vector  $p_u \in \mathbb{R}^k$ .

such that a value close to the actual rating  $r_{ui}$  can be computed (e.g., by the dot product also known as the cosine similarity)

$$r_{ui} \approx \hat{r}_{ui} = q_i^T p_u$$

**The hard part is to find a suitable latent factor space!**

# Singular Value Decomposition (Matrix Fact.)

Linear algebra: Singular Value Decomposition (SVD) to factorizes  $M$

$$M = U\Sigma V^T$$

$M$  is the  $m \times n$  (users  $\times$  items) rating matrix of rank  $r$ .

Columns of  $U$  and  $V$  are the left and right singular vectors.

Diagonal of  $\Sigma$  contains the  $r$  singular values.

# Singular Value Decomposition (Matrix Fact.)

Linear algebra: Singular Value Decomposition (SVD) to factorizes  $M$

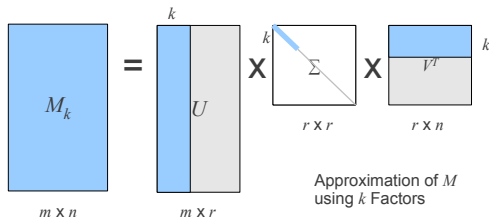
$$M = U\Sigma V^T$$

$M$  is the  $m \times n$  (users  $\times$  items) rating matrix of rank  $r$ .

Columns of  $U$  and  $V$  are the left and right singular vectors.

Diagonal of  $\Sigma$  contains the  $r$  singular values.

A low-rank approximation of  $M$  using only  $k$  factors is straight forward.



The approximation minimizes approx. error  $\|M - M_k\|_F$  (Frobenius norm).

# Challenges (Matrix Fact.)

SVD is  $O(m^3)$  and missing values are a problem.

- 1 Use **Incremental SVD** to add new users/items without recomputing the whole SVD (Sarwar *et al.*, 2002).
- 2 To avoid overfitting minimize the regularized square error on **only known ratings**:

$$\operatorname{argmin}_{p^*, q^*} \sum_{(u, i) \in \kappa} (r_{ui} - q_i^T p_u)^2 + \lambda(\|q_i\|^2 + \|p_u\|^2)$$

where  $\kappa$  are the  $(u, i)$  pairs for which  $r$  is known.

Good solutions can be found by **stochastic gradient descent** or **alternating least squares** (Koren *et al.*, 2009).



# Prediction (Matrix Fact.)

- 1 For new user (item) compute  $q_i$  ( $p_u$ ).
- 2 After all  $q_i$  and  $p_u$  are known, prediction is very fast:

$$\hat{r}_{ui} = q_i^T p_u$$

# Table of Contents

- 1 Motivation
- 2 Content-based Approach
- 3 Collaborative Filtering (CF)
  - Memory-based CF
  - Model-based CF
- 4 Further Considerations
  - Strategies for the Cold Start Problem
  - Recommender Systems and Revenue Management
- 5 Implementation
  - Open-Source Tools
  - An Example using recommenderlab
  - Deployment

# Table of Contents

- 1 Motivation
- 2 Content-based Approach
- 3 Collaborative Filtering (CF)
  - Memory-based CF
  - Model-based CF
- 4 Further Considerations**
  - Strategies for the Cold Start Problem**
  - Recommender Systems and Revenue Management
- 5 Implementation
  - Open-Source Tools
  - An Example using recommenderlab
  - Deployment

# Cold Start Problem

What do we recommend to new users for whom we have no ratings yet?

# Cold Start Problem

What do we recommend to new users for whom we have no ratings yet?

- Recommend popular items
- Have some start-up questions (e.g., "What are your 10 favorite movies?")
- Obtain/purchase personal information

# Cold Start Problem

What do we recommend to new users for whom we have no ratings yet?

- Recommend popular items
- Have some start-up questions (e.g., "What are your 10 favorite movies?")
- Obtain/purchase personal information

What do we do with new items?

# Cold Start Problem

What do we recommend to new users for whom we have no ratings yet?

- Recommend popular items
- Have some start-up questions (e.g., "What are your 10 favorite movies?")
- Obtain/purchase personal information

What do we do with new items?

- Content-based filtering techniques.
- Use expert/domain knowledge.
- Pay a focus group to rate new items.

# Table of Contents

- 1 Motivation
- 2 Content-based Approach
- 3 Collaborative Filtering (CF)
  - Memory-based CF
  - Model-based CF
- 4 Further Considerations**
  - Strategies for the Cold Start Problem
  - Recommender Systems and Revenue Management**
- 5 Implementation
  - Open-Source Tools
  - An Example using recommenderlab
  - Deployment



# Revenue Management

Recommender systems have the potential to increase revenue

- cross-selling
- up-selling

How about influencing which items are recommended using revenue considerations?

# Revenue Management

Recommender systems have the potential to increase revenue

- cross-selling
- up-selling

How about influencing which items are recommended using revenue considerations?

What about **trust + incentive to share information?**

# Table of Contents

- 1 Motivation
- 2 Content-based Approach
- 3 Collaborative Filtering (CF)
  - Memory-based CF
  - Model-based CF
- 4 Further Considerations
  - Strategies for the Cold Start Problem
  - Recommender Systems and Revenue Management
- 5 Implementation
  - Open-Source Tools
  - An Example using recommenderlab
  - Deployment

# Table of Contents

- 1 Motivation
- 2 Content-based Approach
- 3 Collaborative Filtering (CF)
  - Memory-based CF
  - Model-based CF
- 4 Further Considerations
  - Strategies for the Cold Start Problem
  - Recommender Systems and Revenue Management
- 5 Implementation
  - Open-Source Tools
  - An Example using recommenderlab
  - Deployment

# Open-Source Implementations

- **Apache Mahout**: ML library including collaborative filtering (Java)
- **C/Matlab Toolkit for Collaborative Filtering** (C/Matlab)
- **Cofi**: Collaborative Filtering Library (Java)
- **Crab**: Components for recommender systems (Python)
- **easyrec**: Recommender for Web pages (Java)
- **LensKit**: CF algorithms from GroupLens Research (Java)
- **MyMediaLite**: Recommender system algorithms. (C#/Mono)
- **RACOFI**: A rule-applying collaborative filtering system
- **Rating-based item-to-item recommender system** (PHP/SQL)
- **recommenderlab**: Infrastructure to test and develop recommender algorithms (R)

See <http://michael.hahsler.net/research/recommender/> for URLs.

# Table of Contents

- 1 Motivation
- 2 Content-based Approach
- 3 Collaborative Filtering (CF)
  - Memory-based CF
  - Model-based CF
- 4 Further Considerations
  - Strategies for the Cold Start Problem
  - Recommender Systems and Revenue Management
- 5 Implementation**
  - Open-Source Tools
  - An Example using recommenderlab**
  - Deployment

# recommenderlab: Reading Data

100k MovieLens ratings data set: The data was collected through `movielens.umn.edu` from 9/1997 to 4/1998. The data set contains about 100,000 ratings (1-5) from 943 users on 1664 movies.

```
R> library("recommenderlab")
R> data(MovieLens)
R> MovieLens
943 x 1664 rating matrix of class 'realRatingMatrix' with
99392 ratings.
R> train <- MovieLens[1:900]
R> u <- MovieLens[901]
R> u
1 x 1664 rating matrix of class 'realRatingMatrix' with 124
ratings.
R> as(u, "list")[[1]][1:5]
      Toy Story (1995)      Babe (1995)
                5                3
Usual Suspects, The (1995)  Mighty Aphrodite (1995)
                5                1
Mr. Holland's Opus (1995)
                5
```

## recommenderlab: Creating Recommendations

```
R> r <- Recommender(train, method = "UBCF")
R> r
Recommender of type 'UBCF' for 'realRatingMatrix'
learned using 900 users.
R> recom <- predict(r, u, n = 5)
R> recom
Recommendations as 'topNList' with n = 5 for 1 users.
R> as(recom, "list")[[1]]
[1] "Fugitive, The (1993)"
[2] "Shawshank Redemption, The (1994)"
[3] "It's a Wonderful Life (1946)"
[4] "Princess Bride, The (1987)"
[5] "Alien (1979)"
```

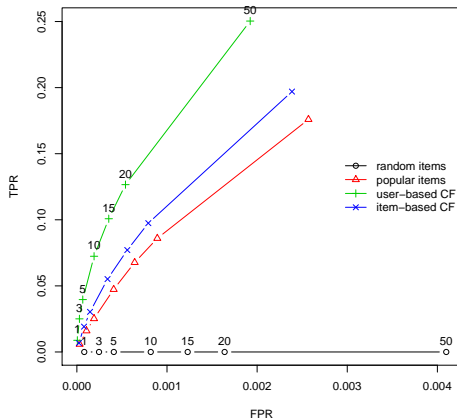


## recommenderlab: Compare Algorithms

```
R> scheme <- evaluationScheme(db, method = "cross", k = 4,  
+   given = 10)  
R> algorithms <- list(  
+ `random items` = list(name = "RANDOM", param = NULL),  
+ `popular items` = list(name = "POPULAR", param = NULL),  
+ `user-based CF` = list(name = "UBCF",  
+   param = list(method = "Cosine", nn = 50)),  
+ `item-based CF` = list(name = "IBCF",  
+   param = list(method = "Cosine", k = 50)))  
R> results <- evaluate(scheme, algorithms,  
+   n = c(1, 3, 5, 10, 15, 20, 50))
```

# recommenderlab: Compare Algorithms II

```
R> plot(results, annotate = c(1, 3), legend = "right")
```



# Table of Contents

- 1 Motivation
- 2 Content-based Approach
- 3 Collaborative Filtering (CF)
  - Memory-based CF
  - Model-based CF
- 4 Further Considerations
  - Strategies for the Cold Start Problem
  - Recommender Systems and Revenue Management
- 5 Implementation**
  - Open-Source Tools
  - An Example using recommenderlab
  - Deployment**

# Technology



<http://techblog.netflix.com>

# References I

- Asim Ansari, Skander Essegaier, and Rajeev Kohli. Internet recommendation systems. *Journal of Marketing Research*, 37:363–375, 2000.
- John S. Breese, David Heckerman, and Carl Kadie. Empirical analysis of predictive algorithms for collaborative filtering. In *Proceedings of the Fourteenth Annual Conference on Uncertainty in Artificial Intelligence*, pages 43–52, 1998.
- Mukund Deshpande and George Karypis. Item-based top-n recommendation algorithms. *ACM Transactions on Information Systems*, 22(1):143–177, 2004.
- David Goldberg, David Nichols, Brian M. Oki, and Douglas Terry. Using collaborative filtering to weave an information tapestry. *Communications of the ACM*, 35(12):61–70, 1992.
- Brendan Kitts, David Freed, and Martin Vrieze. Cross-sell: a fast promotion-tunable customer-item recommendation method based on conditionally independent probabilities. In *KDD '00: Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 437–446. ACM, 2000.
- Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42:30–37, August 2009.
- Andreas Mild and Thomas Reutterer. Collaborative filtering methods for binary market basket data analysis. In *AMT '01: Proceedings of the 6th International Computer Science Conference on Active Media Technology*, pages 302–313, London, UK, 2001. Springer-Verlag.
- Paul Resnick, Neophytos Iacovou, Mitesh Suchak, Peter Bergstrom, and John Riedl. Grouplens: an open architecture for collaborative filtering of netnews. In *CSCW '94: Proceedings of the 1994 ACM conference on Computer supported cooperative work*, pages 175–186. ACM, 1994.
- Badrul Sarwar, George Karypis, Joseph Konstan, and John Riedl. Analysis of recommendation algorithms for e-commerce. In *EC '00: Proceedings of the 2nd ACM conference on Electronic commerce*, pages 158–167. ACM, 2000.
- Badrul Sarwar, George Karypis, Joseph Konstan, and John Riedl. Item-based collaborative filtering recommendation algorithms. In *WWW '01: Proceedings of the 10th international conference on World Wide Web*, pages 285–295. ACM, 2001.
- Badrul Sarwar, George Karypis, Joseph Konstan, and John Riedl. Incremental singular value decomposition algorithms for highly scalable recommender systems. In *Fifth International Conference on Computer and Information Science*, pages 27–28, 2002.
- J. Ben Schafer, Joseph A. Konstan, and John Riedl. E-commerce recommendation applications. *Data Mining and Knowledge Discovery*, 5(1/2):115–153, 2001.

# Thank you!

This presentation can be downloaded from  
<http://michael.hahsler.net/> (under publications/talks)

For questions, please contact the author at [mhahsler@lyle.smu.edu](mailto:mhahsler@lyle.smu.edu)

**recommenderlab** is available in R from CRAN.

An introduction can be found at <https://cran.r-project.org/web/packages/recommenderlab/vignettes/recommenderlab.pdf>