

# **INTRO to DATA SCIENCE**

## **LECTURE 14: RECOMMENDATION SYSTEMS**

### **LAST TIME:**

- ENSEMBLE METHODS**
- BAGGING**
- BOOSTING**
- RANDOM FORESTS**

### **QUESTIONS?**

**I. CONTENT-BASED FILTERING**

**II. COLLABORATIVE FILTERING**

**III. A SIMPLE MATRIX FACTORIZATION MODEL**

**IV. THE NETFLIX PRIZE**

**EXERCISE:**

**V. RECSYS IN PYTHON WITH MOVIELENS DATA**

*The purpose of a **recommendation system** is to predict a rating that a user will give an item that they have not yet rated.*

*The purpose of a **recommendation system** is to predict a rating that a user will give an item that they have not yet rated.*

*This rating is produced by analyzing other user/item ratings (and sometimes item characteristics) to provide personalized recommendations to users.*

*There are two general approaches to recsys design:*

*There are two general approaches to recsys design:*

**Content-based filtering:** *items are mapped into a feature space, and recommendations depend on item characteristics.*

**Collaborative filtering:** *recommendations are only based on user-ratings*

## TV Shows

Your **taste preferences**  
created this row.

TV Shows.

As well as your interest in...





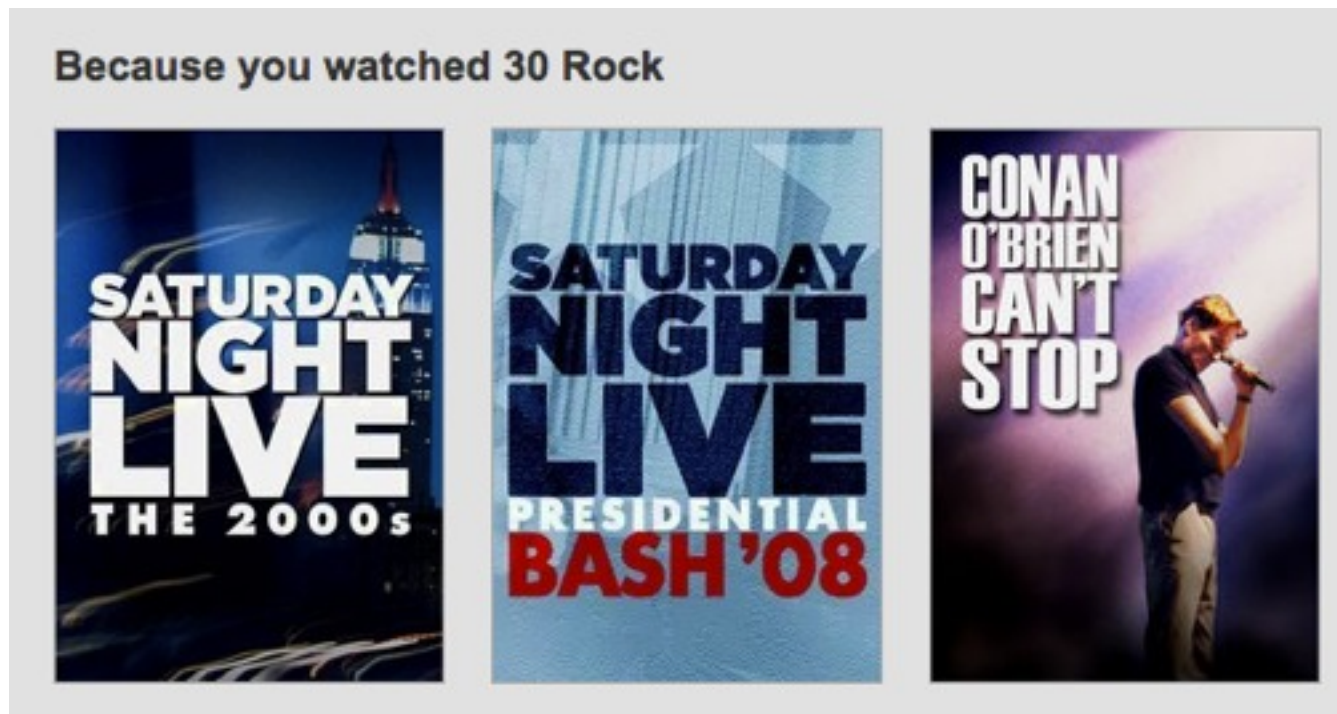
# EXAMPLES – NETFLIX: CONTENT-BASED FILTERING

9

## Taste Preferences

Select a category type: **Featured Categories** ▼

How often do you watch	Never	Sometimes	Often	
<b>Moods</b>				
Absurd	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Need some examples?
Campy	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Need some examples?
Cerebral	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Need some examples?
Chilling	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Need some examples?
Controversial	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Need some examples?



Trick 'r Treat (2008)  
Trick or Treat




**Add to DVD Queue**

This is like "You" →


⊘ ★ ★ ★ ★ ☆

Average of raters like you: 3.3 stars  
Average of 37,634 ratings: 3.5 stars


← This is like "All"




Recommended for you because you watched  
[Sugar Minott - Oh Mr De \(Studio One\)](#)



**Mikey Dread - Roots and Culture**  
by klaxonklaxon · 1,164,133 views  
Lyrics:  
Now here comes a special request  
To each and everyone



Recommended for you because you watched  
[Thelonious Monk Quartet - Monk in Denmark](#)



**Bill Evans Portrait in Jazz (Full Album)**  
by hansgy1 · 854,086 views  
Bill Evans Portrait in Jazz 1960  
1. Come Rain or Come Shine - 3.19 (0:00)  
2. Autumn Leaves - 5.23 (3:24)

amazon.com

## Recommended for You



### [MapReduce Design Patterns: Building Effective Algorithms and Analytics for Hadoop and Other Systems](#)

by Donald Miner (December 22, 2012)

In Stock

**List Price:** \$44.99

**Price:** \$33.18

[54 used & new](#) from \$23.77

[Add to Cart](#)

[Add to Wish List](#)

Rate this item



☐ I own it

☐ Not interested

## Because you purchased...



### [Programming Hive](#) (Paperback)

by Edward Capriolo (Author), et al.



☐ This was a gift

☐ Don't use for recommendations

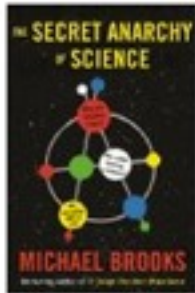
# EXAMPLES – AMAZON: WHAT ABOUT THIS ONE?

14

## Inspired by Your Wish List

You wished for

Customers who viewed this also viewed

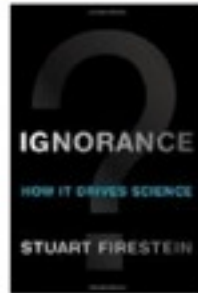


The Secret Anarchy of Science

► Michael Brooks

Paperback

★★★★☆ (6)



Ignorance: How It Drives Science

► Stuart Firestein

Hardcover

★★★★☆ (31)

~~\$21.95~~ **\$13.02**



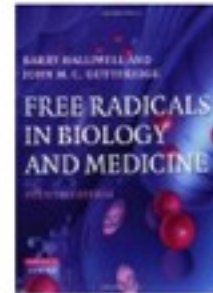
13 Things that Don't Make Sense: The...

► Michael Brooks

Paperback

★★★★☆ (65)

~~\$15.95~~ **\$12.49**



Free Radicals in Biology and Medicine

Barry Halliwell, John Gutteridge

Paperback

★★★★★ (6)

~~\$90.00~~ **\$75.78**



Nonsense on Stilts: How to Tell...

► Massimo Pigliucci

Paperback

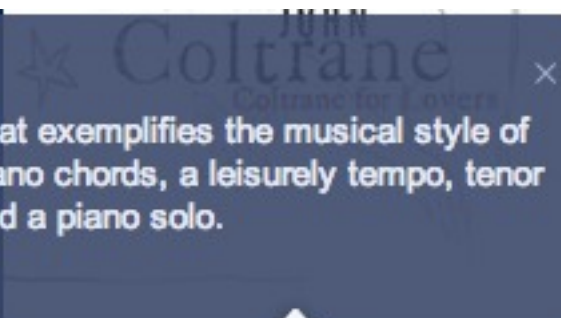
★★★★☆ (35)

~~\$20.00~~ **\$11.94**

## John Coltrane Radio

To start things off, we'll play a song that exemplifies the musical style of John Coltrane which features block piano chords, a leisurely tempo, tenor sax head, a melodic tenor sax solo and a piano solo.

That's not what I wanted, delete this station





## MOST E-MAILED

## RECOMMENDED FOR YOU

1. **How Big Data Is Playing Recruiter for Specialized Workers**
2. SLIPSTREAM  
**When Your Data Wanders to Places You've Never Been**
3. MOTHERLODE  
**The Play Date Gun Debate**
4. **For Indonesian Atheists, a Community of Support Amid Constant Fear**
5. **Justice Breyer Has Shoulder Surgery**
6. BILL KELLER  
**Erasing History**



### 8. How do you determine my Most Read Topics?

[Back to top](#) ▲

Each NYTimes.com article is assigned topic tags that reflect the content of the article. As you read articles, we use these tags to determine your most-read topics.

To search for additional articles on one of your most-read topics, click that topic on your personalized Recommendations page. To learn more about topic tags, visit [Times Topics](#).

# **I. CONTENT-BASED FILTERING**

**Content-based filtering** *begins by mapping each item into a feature space. Both users and items are represented by vectors in this space.*

**Content-based filtering** *begins by mapping each item into a feature space. Both users and items are represented by vectors in this space.*

*Item vectors measure the degree to which the item is described by each feature, and user vectors measure a user's preferences for each feature.*

**Content-based filtering** *begins by mapping each item into a feature space. Both users and items are represented by vectors in this space.*

*Item vectors measure the degree to which the item is described by each feature, and user vectors measure a user's preferences for each feature.*

*Ratings are generated by taking dot products of user & item vectors.*

**Content-based filtering** *begins by mapping each item into a feature space. Both users and items are represented by vectors in the*

**NOTE**

The idea is that users like items that are *similar* to other items they've consumed.

*Item vectors measure the degree to which the item is described by a particular feature, and user vectors measure a user's preferences for a particular feature.*

*Ratings are generated by taking dot products of user & item vectors.*

<i>Movies</i>	<i>  big box office  </i>	<i>aimed at kids  </i>	<i>famous actors</i>
<i>Finding Nemo</i>	<i>5</i>	<i>5</i>	<i>2</i>
<i>Mission Impossible</i>	<i>3</i>	<i>-5</i>	<i>5</i>
<i>Jiro Dreams of Sushi</i>	<i>-4</i>	<i>-5</i>	<i>-5</i>

<i>Movies</i>	<i>  big box office</i>	<i>  aimed at kids</i>	<i>  famous actors</i>
<i>Finding Nemo</i>	5	5	2
<i>Mission Impossible</i>	3	-5	5
<i>Jiro Dreams of Sushi</i>	-4	-5	-5

<i>User</i>	<i>  Big box office</i>	<i>  Aimed at kids</i>	<i>  Famous actors</i>
<i>Jayne</i>	-3	2	-2
<i>Bob</i>	5	-4	5



<i>Movies</i>	<i>  big box office</i>	<i>  aimed at kids</i>	<i>  famous actors</i>
<i>Finding Nemo</i>	$5 \times -3$	$5 \times 2$	$2 \times -2 = -9$
<i>Mission Impossible</i>	$3 \times -3$	$-5 \times 2$	$5 \times -2 = -29$
<i>Jiro Dreams of Sushi</i>	$-4 \times -3$	$-5 \times 2$	$-5 \times -2 = +12$

<i>User</i>	<i>  Big box office</i>	<i>  Aimed at kids</i>	<i>  Famous actors</i>
<i>Jayne</i>	<i>-3</i>	<i>2</i>	<i>-2</i>
<i>Bob</i>	<i>5</i>	<i>-4</i>	<i>5</i>

<i>Movies</i>	<i>  big box office  </i>	<i>aimed at kids  </i>	<i>famous actors</i>
<i>Finding Nemo</i>	$5 \times 5$	$5 \times -4$	$2 \times 5 = +15$
<i>Mission Impossible</i>	$3 \times 5$	$-5 \times -4$	$5 \times 5 = +60$
<i>Jiro Dreams of Sushi</i>	$-4 \times 5$	$-5 \times -4$	$-5 \times 5 = -65$

<i>User</i>	<i>  Big box office  </i>	<i>Aimed at kids  </i>	<i>Famous actors</i>
<i>Jayne</i>	-3	2	-2
<i>Bob</i>	5	-4	5

*One notable example of content-based filtering is Pandora, which maps songs into a feature space using features (or “genes”) designed by the Music Genome Project.*

*Using song vectors that depend on these features, Pandora can create a station with music having similar properties to a song the user selects.*

## About The Music Genome Project®

We believe that each individual has a unique relationship with music – no one else has tastes exactly like yours. So delivering a great radio experience to each and every listener requires an incredibly broad and deep understanding of music. That's why Pandora is based on the Music Genome Project, the most sophisticated taxonomy of musical information ever collected. It represents over ten years of analysis by our trained team of musicologists, and spans everything from this past Tuesday's new releases all the way back to the Renaissance and Classical music.

Each song in the Music Genome Project is analyzed using up to 450 distinct musical characteristics by a trained music analyst. These attributes capture not only the musical identity of a song, but also the many significant qualities that are relevant to understanding the musical preferences of listeners. The typical music analyst working on the Music Genome Project has a four-year degree in music theory, composition or performance, has passed through a selective screening process and has completed intensive training in the Music Genome's rigorous and precise methodology. To qualify for the work, analysts must have a firm grounding in music theory, including familiarity with a wide range of styles and sounds.

*source: <http://www.pandora.com/about/mgp>*

*Content-based filtering has some difficulties:*

*Content-based filtering has some challenges:*

- need to map each item into a feature space (usually by hand!)*
- recs limited in scope (items must be similar to each other)*
- hard to create cross-content recommendations (eg books/music films...this would require comparing elements from different feature spaces!)*

# II. COLLABORATIVE FILTERING

*Instead of thinking about users and items in terms of a feature space, we are only interested in the existing user-item ratings themselves.*



*Instead of thinking about users and items in terms of a feature space, we are only interested in the existing user-item ratings themselves.*

*Our dataset is a ratings matrix of items and users.*

*Instead of thinking about users and items in terms of a feature space, we are only interested in the existing user-item ratings themselves.*

*Our dataset is a ratings matrix of items and users.*

*Assumption: users get value from recommendations based on other users with similar tastes*

x	1	1	x	...	x
x	x	x	5	...	x
x	x	3	x	...	x
x	4	3	x	...	2
...	x	x	x	...	x
x	5	x	1	...	x
x	x	3	3	...	x
x	1	x	x	...	2

**NOTE**

This matrix will always be *sparse*!

source: <http://www.eecs.berkeley.edu/~zhanghao/main/publications/subfolder/netflix.png>

*Collaborative filtering can be done in two different ways.*

*Collaborative filtering can be done in two different ways.*

***Memory-based:*** *uses a sample of the users that are most similar to a given user to predict ratings on unrated items.*

*Collaborative filtering can be done in two different ways.*

***Memory-based:*** *uses a sample of the users that are most similar to a given user to predict ratings on unrated items.*

***Model-based:*** *extracts complex patterns from the dataset, and uses that as a "model" to make recommendations without having to use the dataset every time*

*For a given user ( $u_j$ ), find users that are most similar (e.g. **vector similarity** or **Pearson correlation coefficient**)*

*For a given item  $i_j$ , the predicted rating of  $u_j$  is the average of the known ratings of  $i_j$  within the group of similar users.*

*We just talked about user-based CF, but the same process can be done for item-based CF.*

*Item-based CF is more commonly used than user-based CF, why do you think that is?*



*Model-based collaborative filtering abandons the neighbor approach and applies other techniques to the ratings matrix.*

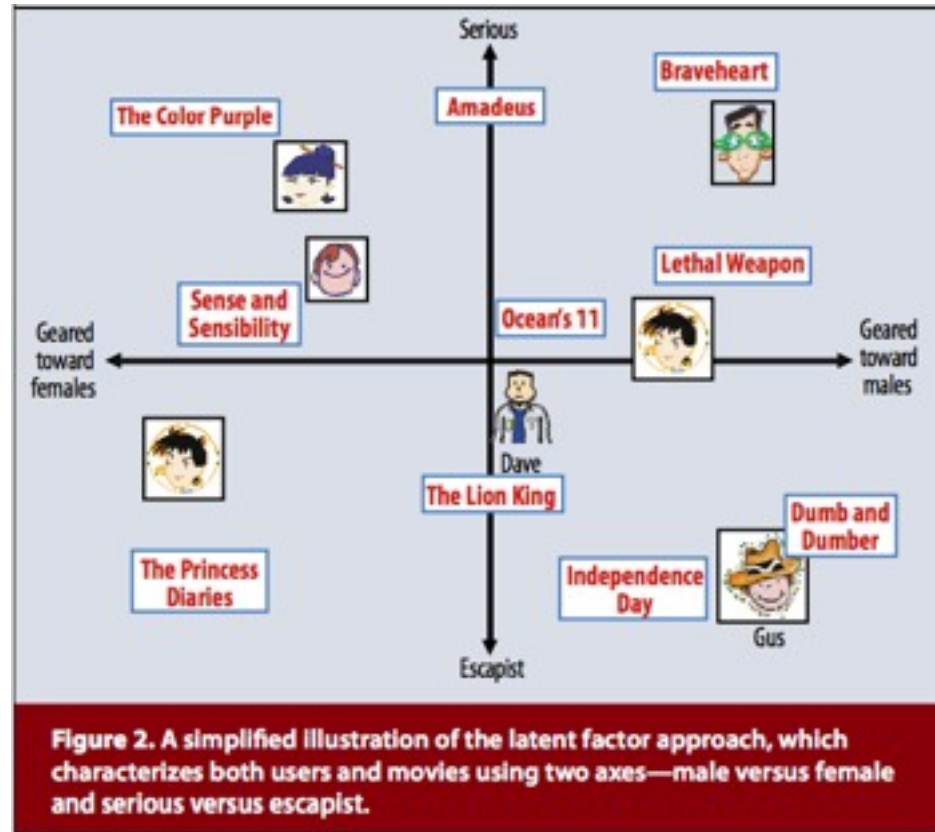
*Model-based collaborative filtering abandons the neighbor approach and applies other techniques to the ratings matrix.*

*The most popular model-based CF techniques use matrix decomposition techniques to find deeper structure in the ratings data.*

*Model-based collaborative filtering abandons the neighbor approach and applies other techniques to the ratings matrix.*

*The most popular model-based CF techniques use matrix decomposition techniques to find deeper structure in the ratings data.*

*For example, we could decompose the ratings matrix via SVD to reduce the dimensionality and extract latent variables.*



source: <http://www2.research.att.com/~volinsky/papers/ieeecomputer.pdf>

*Once we identify the latent variables in the ratings matrix, we can express both users and items in terms of these latent variables.*

*Once we identify the latent variables in the ratings matrix, we can express both users and items in terms of these latent variables.*

*As before, values in the item vectors represent the degree to which an item exhibits a given feature, and values in the user vectors represent user preferences for a given feature.*

*Once we identify the latent variables in the ratings matrix, we can express both users and items in terms of these latent variables.*

*As before, values in the item vectors represent the degree to which an item exhibits a given feature, and values in the user vectors represent user preferences for a given feature.*

*Ratings are constructed by taking dot products of user & item vectors in the latent feature space.*

*This approach is domain independent, and requires no explicit user or item profiles to be created.*



*This approach is domain independent, and requires no explicit user or item profiles to be created.*

*It combines predictive accuracy, scalability, and enough flexibility for practical modeling (we'll see what this means in a moment).*

*This approach is domain independent, and requires no explicit user or item profiles to be created.*

*It combines predictive accuracy, scalability, and enough flexibility for practical modeling (we'll see what this means in a moment).*

*Since the conclusion of the Netflix prize, these latent factor methods for collaborative filtering have been regarded as the state-of-the-art in recsys technology.*

*But they do have some drawbacks:*

- lots of (high-dimensional) ratings data needed*
- data is typically very sparse (in the Netflix prize dataset, ~99% of possible ratings were missing)*
- susceptible to fraud (eg shilling attacks)*
- **cold start problem:** need lots of data on new user or item before recommendations can be made*

*The cold start problem arises because we've been relying only on ratings data, or on **explicit feedback** from users.*

*The cold start problem arises because we've been relying only on ratings data, or on **explicit feedback** from users.*

*Until a user rates several items, we don't know anything about her preferences!*

*The cold start problem arises because we've been relying only on ratings data, or on **explicit feedback** from users.*

*Until a user rates several items, we don't know anything about her preferences!*

*We can get around this by enhancing our recommendations using **implicit feedback**, which may include things like item browsing behavior, search patterns, purchase history, etc.*

*While explicit feedback (ratings, likes, purchases) leads to high quality ratings, the data is sparse and cold starts are problematic.*

*While explicit feedback (ratings, likes, purchases) leads to high quality ratings, the data is sparse and cold starts are problematic.*

*Meanwhile implicit feedback (browsing behavior, etc) leads to less accurate ratings, but the data is much more dense (and less invasive to collect).*



*While explicit feedback (ratings, likes, purchases) leads to high quality ratings, the data is sparse and cold starts are problematic.*

*Meanwhile implicit feedback (browsing behavior, etc) leads to less accurate ratings, but the data is much more dense (and less invasive to collect).*

*Implicit feedback can help to infer user preferences when explicit feedback is not available, therefore easing the cold start problem.*

**Hybrid filtering methods** *provide another way to get around the cold start problem by combining filtering methods (eg, by using content-based info to “boost” a collaborative model).*

**Hybrid filtering methods** *provide another way to get around the cold start problem by combining filtering methods (eg, by using content-based info to “boost” a collaborative model).*

*This content-based info can be item-based as above, or even user-based (eg, demographic info).*

**Hybrid filtering methods** *provide another way to get around the cold start problem by combining filtering methods (eg, by using content-based info to “boost” a collaborative model).*

*This content-based info can be item-based as above, or even user-based (eg, demographic info).*

*Hybrid methods can also make the data sparsity issue easier to deal with, by broadening the set of features under consideration.*

# **III. A SIMPLE MATRIX FACTORIZATION MODEL**

*Matrix factorization decomposes the ratings matrix and maps users and items into a low-dimensional vector space spanned by a basis of latent factors.*

*Matrix factorization decomposes the ratings matrix and maps users and items into a low-dimensional vector space spanned by a basis of latent factors.*

*Predicted ratings are given by inner products in this space, so for user  $u$  and item  $i$  we can write:*

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

*Factoring the ratings matrix via SVD leads to difficulty, since the matrix is typically sparse and therefore our information about the data is incomplete.*



*Factoring the ratings matrix via SVD leads to difficulty, since the matrix is typically sparse and therefore our information about the data is incomplete.*

*Interpolating missing values is an expensive process and can lead to inaccurate predictions, so we need another way to perform this factorization.*

*One possibility is to learn the feature vectors using the observed ratings only. Since this dramatically reduces the size of the ratings matrix, we have to be careful to avoid overfitting.*

*One possibility is to learn the feature vectors using the observed ratings only. Since this dramatically reduces the size of the ratings matrix, we have to be careful to avoid overfitting.*

*We can learn these feature vectors by minimizing the loss function:*

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

*where  $\kappa$  denotes the set of known ratings, and  $\lambda$  is a hyperparameter.*

*One possibility is to learn the feature vectors using the observed ratings only. Since this dramatically reduces the size of the ratings have to be careful to avoid overfitting.*

**NOTE**

The loss function has two unknowns ( $q$ ,  $p$ ) and so is not convex!

This can be minimized using a method called *alternating least squares*.

*We can learn these feature vectors by minimizing the loss f*

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

*where  $\kappa$  denotes the set of known ratings, and  $\lambda$  is a hyperparameter.*

*It turns out that much of the variation in observed ratings is due to user or item biases (eg, some users are very critical, or some items are universally popular).*

*It turns out that much of the variation in observed ratings is due to user or item biases (eg, some users are very critical, or some items are universally popular).*

*We can capture these biases in our model by generalizing  $\hat{r}_{ui}$ .*

$$\hat{r}_{ui} = \mu + b_i + b_u + q_i^T r_u$$

*Here  $\mu$  is a global average rating,  $b_i$  is the item bias,  $b_u$  is the user bias, and  $q_i^T r_u$  is the user-item interaction.*

*With this generalization, our minimization problem becomes:*

$$\min_{q,p,b} \sum_{(u,i) \in \kappa} (r_{ui} - \mu - b_u - b_i - q_i^T p_u)^2 + \lambda(\|q_i\|^2 + \|p_u\|^2 + b_u^2 + b_i^2)$$

*With this generalization, our minimization problem becomes:*

$$\min_{q,p,b} \sum_{(u,i) \in \kappa} (r_{ui} - \mu - b_u - b_i - q_i^T p_u)^2 + \lambda(\|q_i\|^2 + \|p_u\|^2 + b_u^2 + b_i^2)$$

*Further modifications can be made to this model (incorporating implicit feedback, capturing temporal effects, attaching confidence scores to predictions), and you can look up the details in the references.*



# **IV. THE NETFLIX PRIZE**

*The Netflix prize was a competition to see if anyone could make a 10% improvement to Netflix's recommendation system (accuracy measured by RMSE).*

*The Netflix prize was a competition to see if anyone could make a 10% improvement to Netflix's recommendation system (accuracy measured by RMSE).*

*The grand prize was \$1mm dollars, with annual \$50k progress prizes to the leader at the end of each year if the 10% threshold had not yet been met. Approx 50k teams participated from >180 countries.*

*The Netflix prize was a competition to see if anyone could make a 10% improvement to Netflix's recommendation system (accuracy measured by RMSE).*

*The grand prize was \$1mm dollars, with annual \$50k progress prizes to the leader at the end of each year if the 10% threshold had not yet been met. Approx 50k teams participated from >180 countries.*

*The ratings matrix contained >100mm numerical entries (1-5 stars) from ~500k users across ~17k movies. The data was split into train/quiz/test sets to prevent overfitting on the test data by answer submission (this was a clever idea!)*

*The competition began in 2006, and the grand prize was eventually awarded in 2009. The winning entry was a stacked ensemble of 100's of models (including neighborhood & matrix factorization models) that were blended using boosted decision trees.*

*The competition began in 2006, and the grand prize was eventually awarded in 2009. The winning entry was a stacked ensemble of 100's of models (including neighborhood & matrix factorization models) that were blended using boosted decision trees.*

*Ultimately, the competition ended in a photo finish. The winning strategy came down to last-minute team mergers & creative blending schemes to shave 3<sup>rd</sup> & 4<sup>th</sup> decimals off RMSE (concerns that would not be important in practice).*

*The competition began in 2006, and the grand prize was eventually awarded in 2009. The winning entry was a stacked ensemble of 100's of models (including neighborhood & matrix factorization models) that were blended using boosted decision trees.*

*Ultimately, the competition ended in a photo finish. The winning strategy came down to last-minute team mergers & creative blending schemes to shave 3<sup>rd</sup> & 4<sup>th</sup> decimals off RMSE (concerns that would not be important in practice).*

*The competition did much to spur interest and research advances in recsys technology, and the prize money was donated to charity.*

*Though they adopted some of the modeling techniques that emerged from the competition, Netflix never actually implemented the prizewinning solution.*

*Why do you think that's true?*