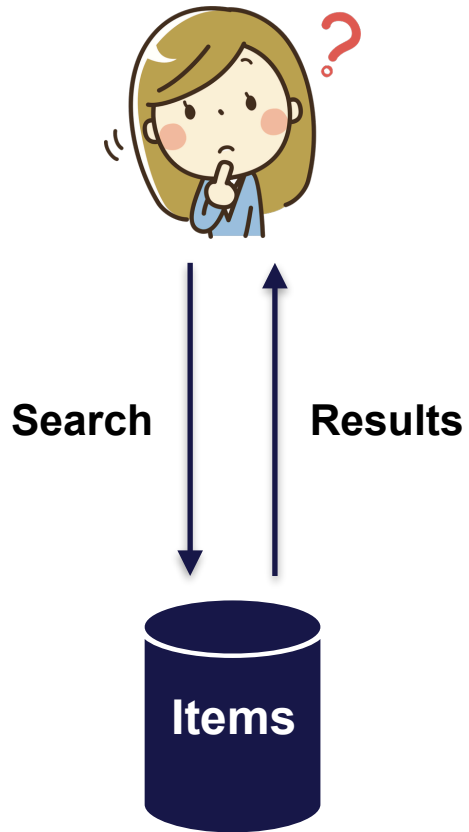


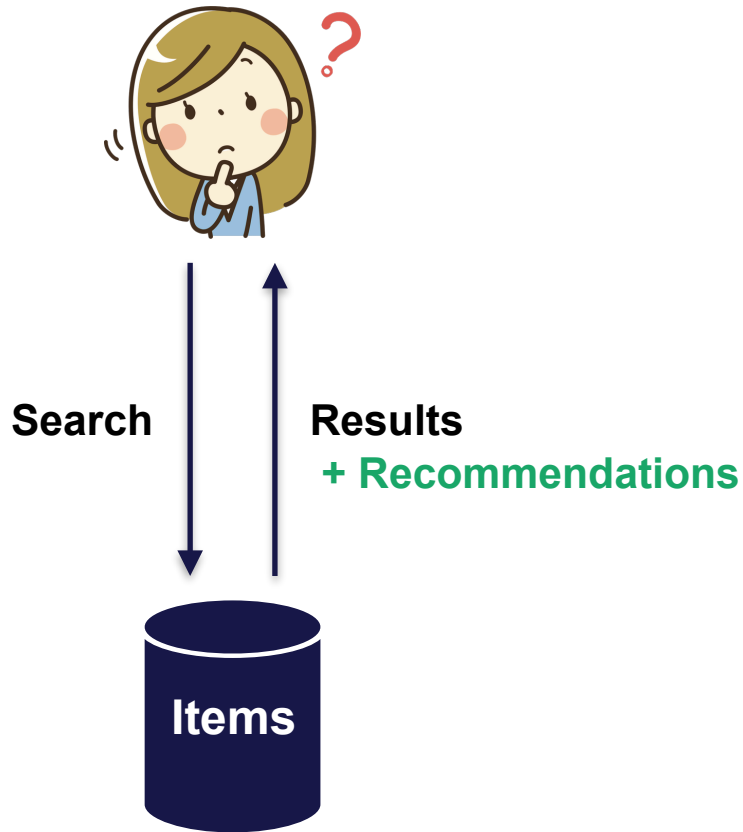
42578 Advanced Business Analytics

Recommender Systems

Recommendations



Recommendations

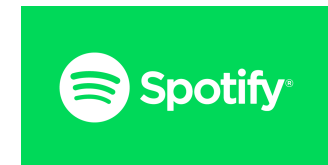


Examples

amazon



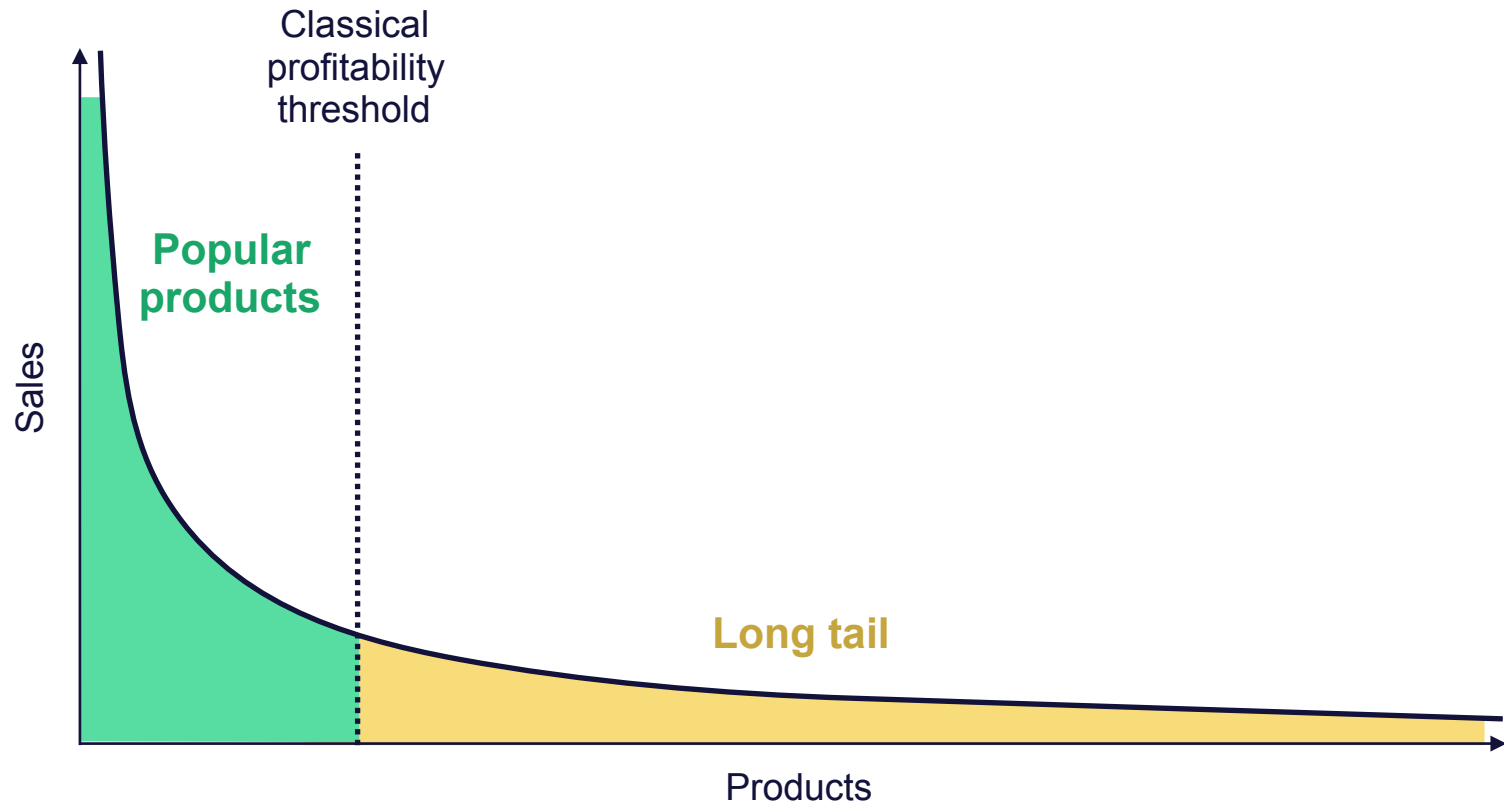
You Tube



From Scarcity to Abundance

- **Shelf space is a scarce commodity for traditional retailers**
 - Also: TV networks, movie theatres,...
- **Web enables near-zero-cost dissemination of information about products**
 - From scarcity to abundance
- **More choice necessitates better filters**
 - Recommendation engines
 - How “Into Thin Air” made “Touching the Void” a bestseller:
<http://www.wired.com/wired/archive/12.10/tail.html>

The Long Tail



Types of Recommendations

- **Editorial and hand-curated**
 - Editors' picks
 - Lists of “essential” items
- **Simple aggregates**
 - Top 10, Most Popular, Recent Uploads
- **Tailored to individual users**
 - Amazon, Netflix, ...

Formal Model

- X = set of **Customers**
- I = set of **Items**
- **Utility function** $u: X \times I \rightarrow R$
 - R = set of ratings
 - R is a totally ordered set
 - e.g., **0-5** stars, real number in **[0,1]** (purchase probability), ...

Rating (Utility) Matrix

	Avatar	LOTR	Matrix	Pirates
Alice	1	?	0.2	?
Bob	?	0.5	?	0.3
Carol	0.2	?	1	?
David	?	?	?	0.4

Key Problems

- **(1) Gathering “known” ratings for the matrix**
 - How to collect the data in the rating matrix
- **(2) Extrapolate unknown ratings from the known ones**
 - Mainly interested in high unknown ratings
 - We are not interested in knowing what you don't like but what you like
- **(3) Evaluating extrapolation methods**
 - How to measure success/performance of recommendation methods

(1) Gathering Ratings

- **Explicit**
 - Ask people to rate items
 - Doesn't work well in practice – people can't be bothered
- **Implicit**
 - Learn ratings from user actions
 - E.g., purchase implies high rating
 - What about low ratings?

(2) Extrapolating Ratings

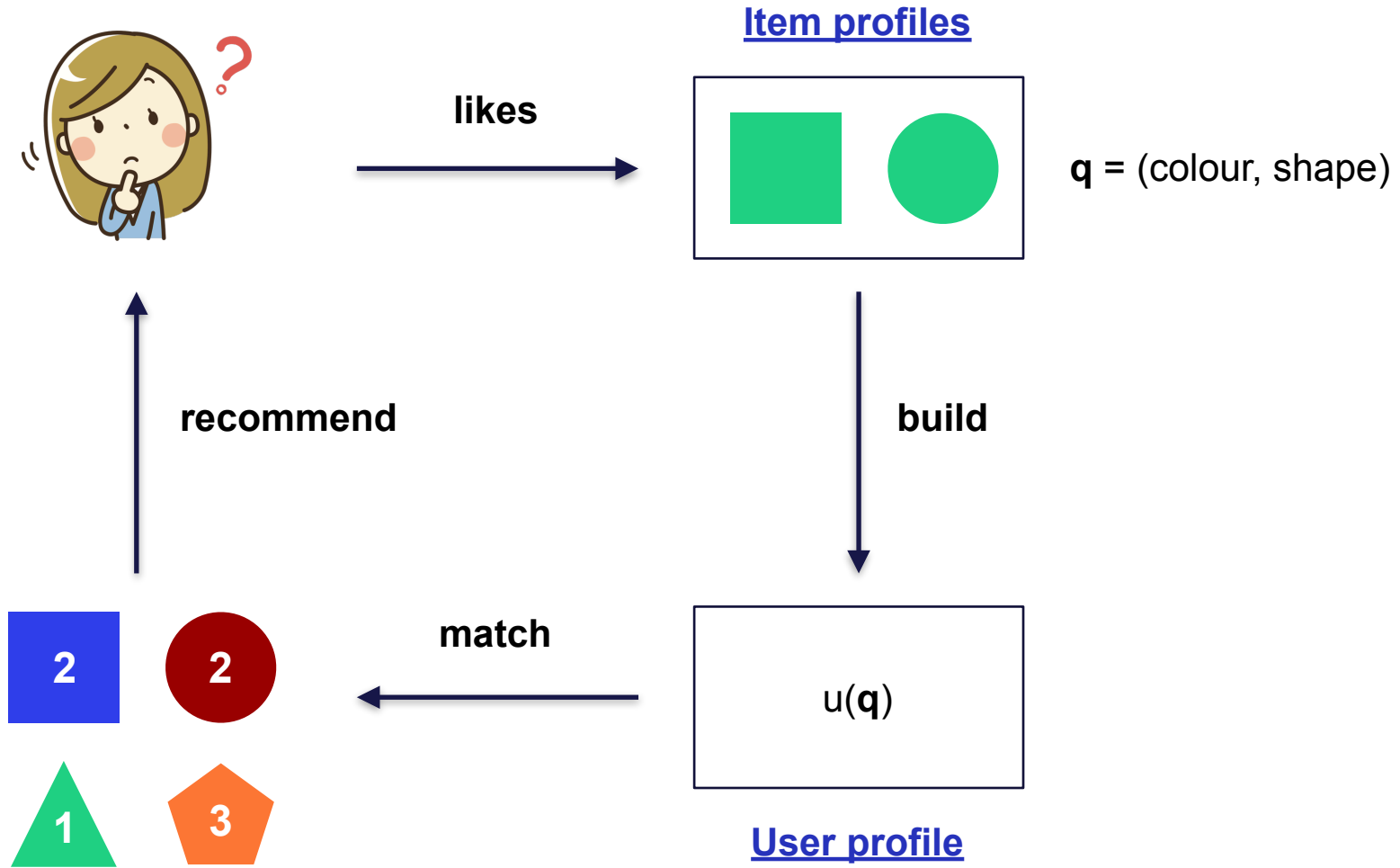
- **Key problem:** Rating matrix \mathbf{R} is **sparse**
 - Most people have not rated most items
 - **Cold start:**
 - New items have no ratings
 - New users have no history
- **Two approaches to recommender systems:**
 - (1) Content-based recommendations
 - (2) Collaborative filtering
 - Memory-based approach: Item-Item, User-User
 - Model-based approach: Matrix factorisation, Latent factors, PCA, Neural nets, ...
 - Hybrid
 - Hybrid

Content-based Recommender Systems

Content-based Recommendations

- **Main idea:** Recommend items to the customer x similar to previous items rated highly by x
- **Example:**
 - **Movie recommendations**
 - Recommend movies with same actor(s), director, genre, ...
 - **Websites, blogs, news**
 - Recommend other sites with “similar” content

Main Idea



Item Profiles

- For each item, create an **item profile**, q_i
- **Item profile is a set (vector) of item features**
 - **Movies:** author, title, actor, director,...
 - **Text:** Set of “important” words in document
 - TF-IDF, topics, ...
 - **Images, music:**
 - Hand-crafted features
 - Unsupervised learning (“embeddings”), e.g., Clustering, Generative modelling

User Profile

- **General case: Utility function $u_x(q)$**
- **Simplest case: Linear model**
 - User profile is a vector of weights \mathbf{p}_x for features in \mathbf{q}
 - linear regression
 - logistic regression (probability of “like”)
- **Nonlinear “black-box” models**
 - Lack of interpretability

Pros: Content-based Approach

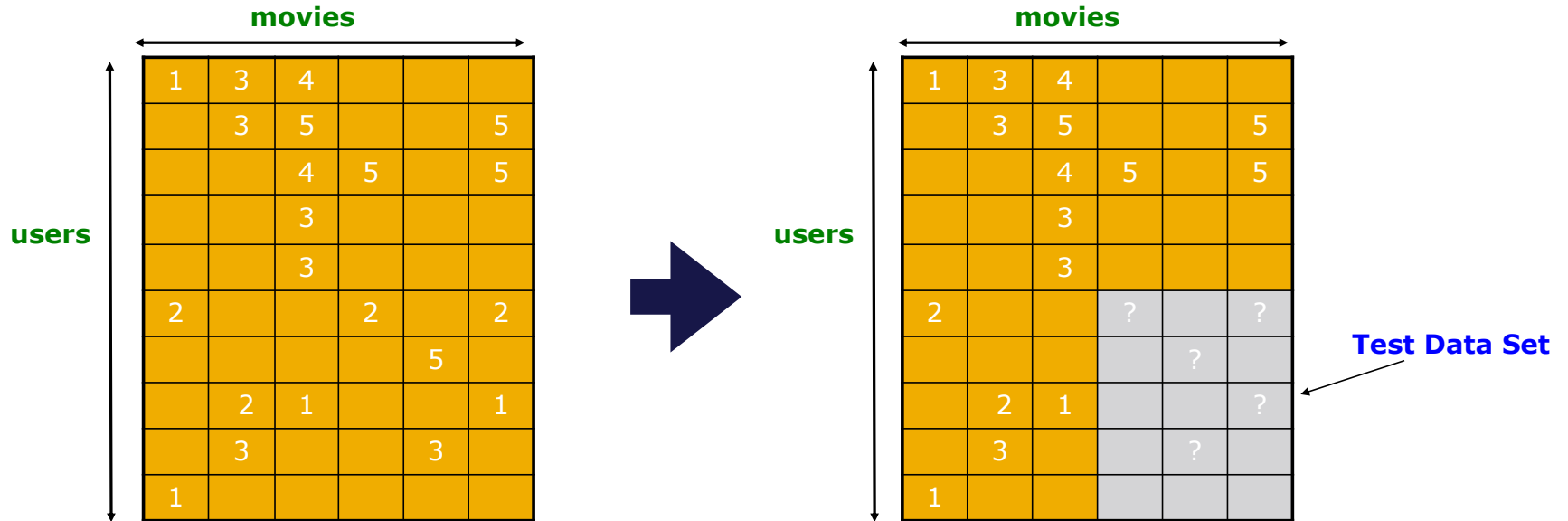
- + Does not heavily depend on data on other users
- + Able to recommend to users with unique tastes
- + Able to recommend new & unpopular items

Cons: Content-based Approach

- **Finding the appropriate features is hard**
 - E.g., images, movies, music
- **Recommendations for new users**
 - Irrelevant recommendations in the beginning
- **Overspecialisation**
 - Never recommends items outside user's content profile
 - People might have multiple interests
 - **Unable to exploit quality judgments of other users**

Evaluation of Recommender Systems

Evaluation of Recommender Systems



- Train / Test split
- Cross-validation

Performance Measure

- **Root-mean-square error (RMSE)**

$$\sqrt{\frac{1}{N} \sum_{xi} (\hat{r}_{xi} - r_{xi})^2}$$
 where \hat{r}_{xi} is predicted and r_{xi} is true rating of x on i

- **cons:** In practice, we care only to predict high ratings, and RMSE might penalise a method that does well for high ratings and badly for others
- **Precision at top k**
 - % of those in top k predictions (in practice, often $k=10$)
- **Rank (Spearman's) correlation**

Baselines

- **Global mean rating, μ**
 - average of all ratings
- **User mean rating, μ_x**
 - average of all ratings from the user x
- **Item mean rating, μ_i**
 - average of all ratings for the item i
- **Combination**
 - $r_{xi} = \mu + b_x + b_i = \mu + (\mu_x - \mu) + (\mu_i - \mu)$

User's bias
(deviation of the
user's mean from
the global mean)

Item's bias
(deviation of the
item's mean from
the global mean)

Exercises

- “Recommender Systems.ipynb”: Sections 1 - 3.1 (Content-based filtering)

Memory-based Collaborative Filtering



Person A

- likes private jets
- likes beer



Person A

- likes private jets
- likes beer



Person B

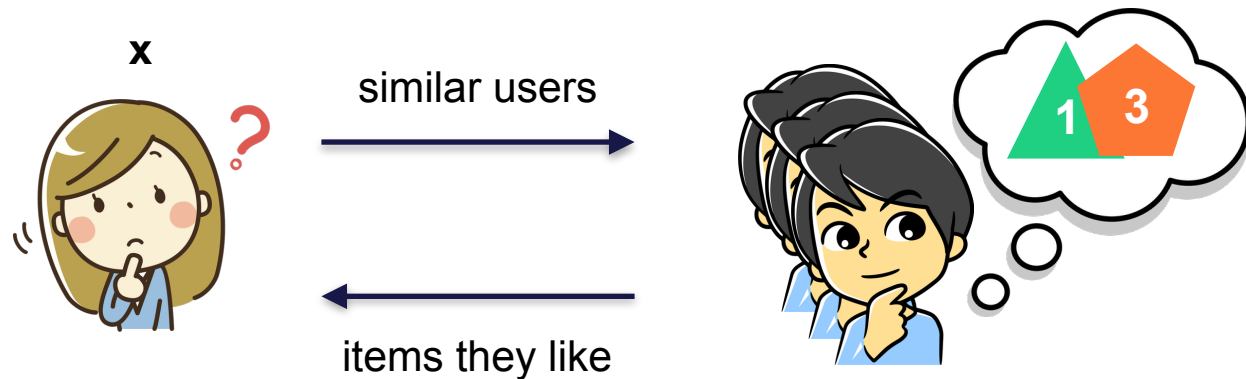
- wants to buy beer
- recommend to buy a private jet?

* Do not judge a person based on just one photograph

<https://www.buzzfeednews.com/article/miriamelder/russians-who-look-like-hollywood-stars>

User-User Collaborative Filtering

- Consider user x
- Find similarity of x to **all other users**
- Estimate x 's ratings using ratings of the other users weighted by their similarity to x
- **K-Nearest-Neighbor approach:** all other users $\rightarrow K$ most similar users



Finding “Similar” Users

- Let \mathbf{r}_x be the vector of user \mathbf{x} 's ratings: $\mathbf{r}_x = (1, ?, ?, 4, 5)$
and \mathbf{r}_y be the vector of user \mathbf{y} 's ratings: $\mathbf{r}_y = (?, 1, ?, 3, 4)$
- **Pearson correlation coefficient**
 - \mathbf{S}_{xy} = items rated by both users \mathbf{x} and \mathbf{y} : (4, 5) and (3, 4)

$$\text{sim}(x, y) = \frac{\sum_{s \in \mathbf{S}_{xy}} (r_{xs} - \bar{r}_x)(r_{ys} - \bar{r}_y)}{\sqrt{\sum_{s \in \mathbf{S}_{xy}} (r_{xs} - \bar{r}_x)^2} \sqrt{\sum_{s \in \mathbf{S}_{xy}} (r_{ys} - \bar{r}_y)^2}}$$

Rating Predictions

- Predict r_{xi} — rating of the item i by the user x
- Let \mathbf{K} be the set of K users most similar to x who have rated item i
- **Prediction:**

$$\triangleright r_{xi} = \frac{1}{K} \sum_{y \in \mathbf{K}} r_{yi}$$

$$\triangleright r_{xi} = \frac{\sum_{y \in \mathbf{K}} \text{sim}(x, y) r_{yi}}{\sum_{y \in \mathbf{K}} \text{sim}(x, y)}$$

Item-Item Collaborative Filtering

- So far: **User-user collaborative filtering**
- **Another view: Item-item**
 - For item i , find other similar items
 - Estimate rating for item i based on ratings for similar items
 - Can use same similarity metrics and prediction functions as in the user-user model

$$r_{xi} = \frac{\sum_{j \in \mathbf{K}} \text{sim}(i, j) r_{xj}}{\sum_{j \in \mathbf{K}} \text{sim}(i, j)}$$

- \mathbf{K} is a set of K most similar items to i rated by x

Item-Item Collaborative Filtering: Example ($K=2$)

users

movies

	1	2	3	4	5	6	7	8	9	10	11	12
1	1		3			5			5		4	
2			5	4			4			2	1	3
3	2	4		1	2		3		4	3	5	
4		2	4		5			4			2	
5			4	3	4	2					2	5
6	1		3		3			2			4	



- unknown rating



- rating between 1 to 5

Item-Item Collaborative Filtering: Example ($K=2$)

users

	1	2	3	4	5	6	7	8	9	10	11	12
movies	1		3		?	5			5		4	
2			5	4			4			2	1	3
3	2	4		1	2		3		4	3	5	
4		2	4		5			4			2	
5			4	3	4	2					2	5
6	1		3		3			2			4	



- estimate rating of movie **1** by user **5**

Item-Item Collaborative Filtering: Example ($K=2$)

		users												sim(1,m)
		1	2	3	4	5	6	7	8	9	10	11	12	
movies	1	1		3		?	5			5		4		1.00
	2			5	4			4			2	1	3	-0.18
	<u>3</u>	2	4		1	2		3		4	3	5		<u>0.41</u>
	4		2	4		5			4			2		-0.10
	5			4	3	4	2					2	5	-0.31
	<u>6</u>	1		3		3			2			4		<u>0.59</u>

Neighbor selection:

Identify movies similar to movie 1, rated by user 5

Here we use Pearson correlation as similarity:

1) Subtract mean rating m_i from each movie i

$$m_1 = (1+3+5+5+4)/5 = 3.6$$

row 1: [-2.6, 0, -0.6, 0, 0, 1.4, 0, 0, 1.4, 0, 0.4, 0]

2) Compute cosine similarities between rows

Item-Item Collaborative Filtering: Example ($K=2$)

		users												
		1	2	3	4	5	6	7	8	9	10	11	12	$\text{sim}(1,m)$
movies	1	1		3		2.6	5			5		4		1.00
	2			5	4			4			2	1	3	-0.18
	<u>3</u>	2	4		1	2		3		4	3	5		<u>0.41</u>
	4		2	4		5			4			2		-0.10
	5			4	3	4	2					2	5	-0.31
	<u>6</u>	1		3		3			2			4		<u>0.59</u>

Similarity weights:

$$s_{1,3}=0.41, s_{1,6}=0.59$$

Predict by taking weighted average:

$$r_{1.5} = (0.41 \cdot 2 + 0.59 \cdot 3) / (0.41 + 0.59) = 2.6$$

Collaborative Filtering: Common Practice

$$r_{xi} = \frac{\sum_{y \in \mathbf{K}} \text{sim}(x, y) r_{yi}}{\sum_{y \in \mathbf{K}} \text{sim}(x, y)}$$



$$r_{xi} = \boxed{b_{xi}} + \boxed{\frac{\sum_{y \in \mathbf{K}} \text{sim}(x, y) (r_{yi} - b_{yi})}{\sum_{y \in \mathbf{K}} \text{sim}(x, y)}}$$

Baseline estimate for r_{xi}

Deviation from the baseline

$$b_{xi} = \mu + b_x + b_i$$

Item-Item vs. User-User

- **In theory**, should perform approximately the same
- **In practice**, it has been observed that **item-item often works better** than user-user
- **Why?** Items are simpler, users have multiple tastes

Pros/Cons of Collaborative Filtering

+ Works for any kind of item

- No feature selection needed

- Cold Start:

- Need enough users in the system to find a match

- Sparsity:

- The user/ratings matrix is sparse
- Hard to find users that have rated the same items

- First rater:

- Cannot recommend an item that has not been previously rated
- New items, Esoteric items

- Popularity bias:

- Cannot recommend items to someone with unique taste
- Tends to recommend popular items

- Expensive step is finding **K most similar customers/items (linear complexity)**

- Solution: caching, indexing, nearest-neighbour search, etc.

Hybrid Methods

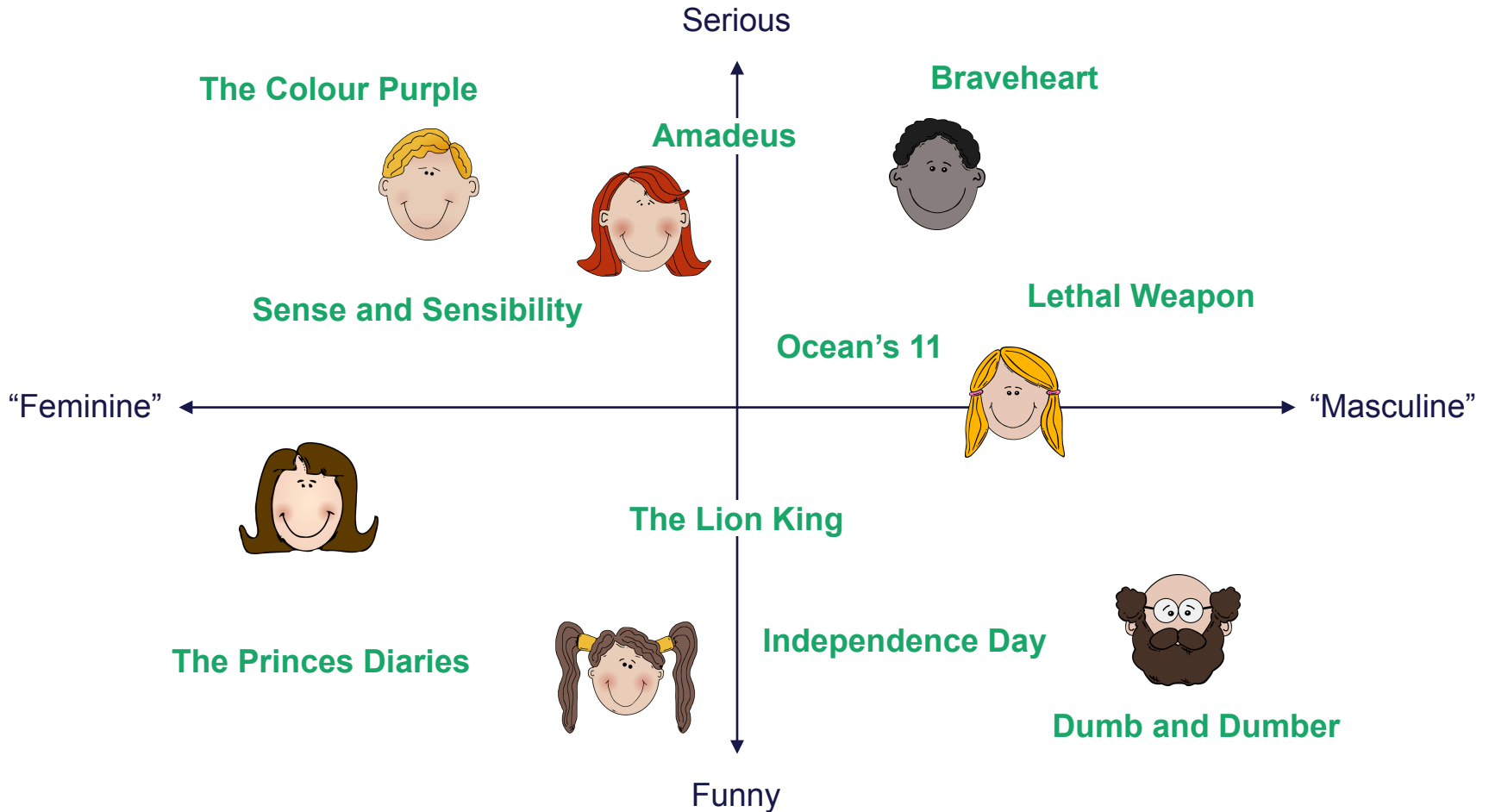
- **Implement two or more different recommenders and combine predictions**
 - Perhaps using a linear model
- **Add content-based methods to collaborative filtering**
 - Item profiles for new item problem
 - Demographics to deal with new user problem

Exercises

- Part 2 - Memory-based recommendations

Model-based Collaborative Filtering

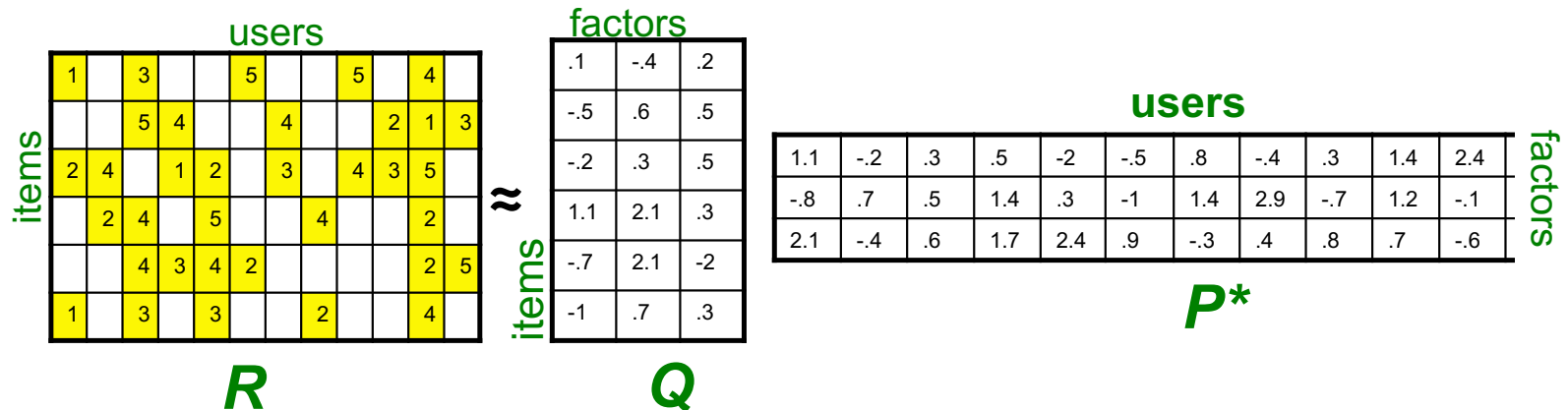
Latent Factor Models



Latent Factor Models

- “SVD” on the utility matrix: $R \approx Q \cdot P^*$

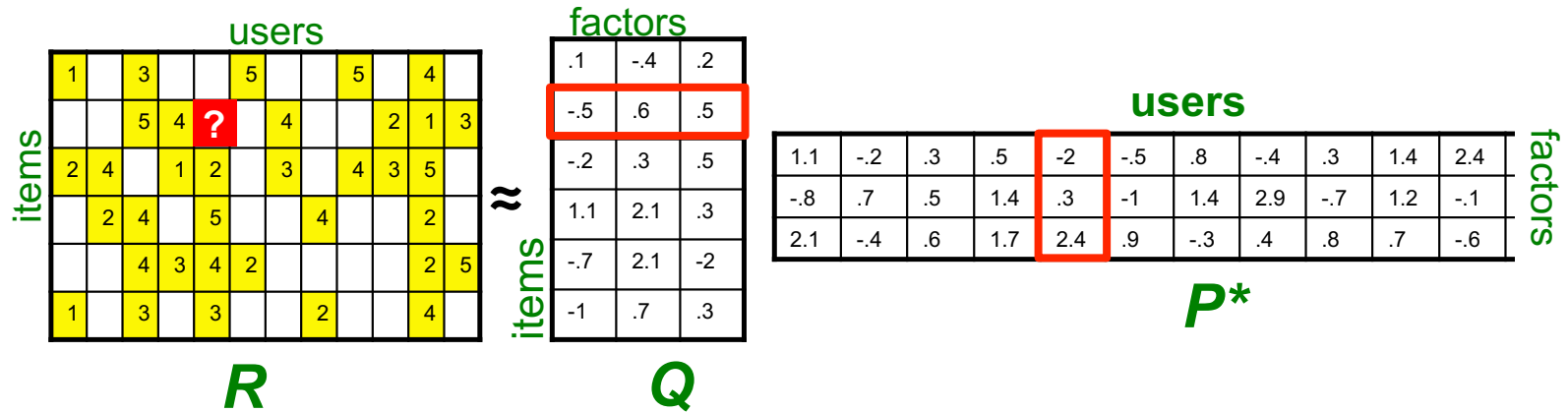
SVD: $D = U \Sigma T^*$



- For now let's assume we can approximate the rating matrix R as a product of “thin” $Q \cdot P^*$
 - R has missing entries but let's ignore that for now!
 - Basically, we will want the reconstruction error to be small on known ratings and we don't care about the values on the missing ones

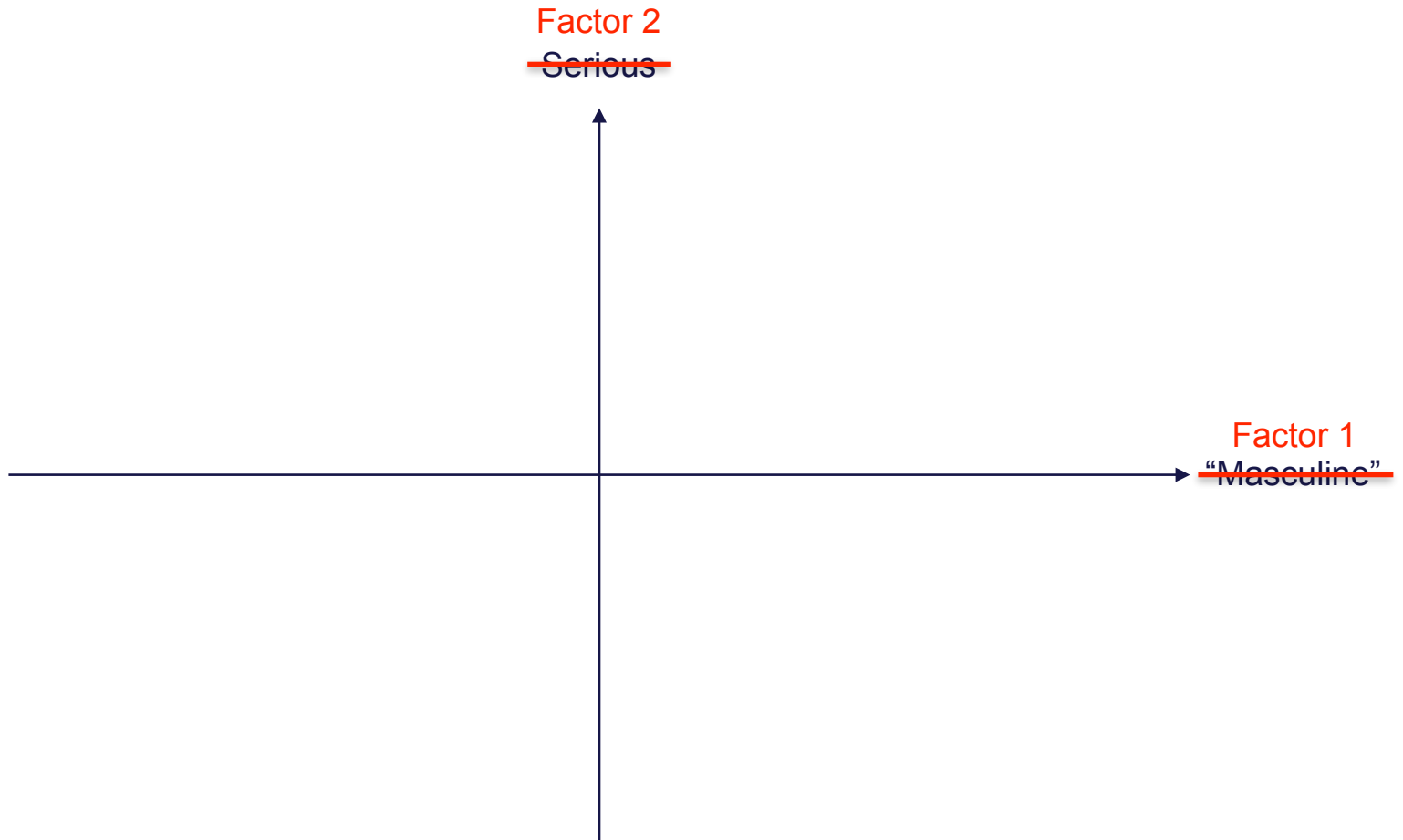
Ratings as Products of Factors

- How to estimate the missing rating of user x for item i ?

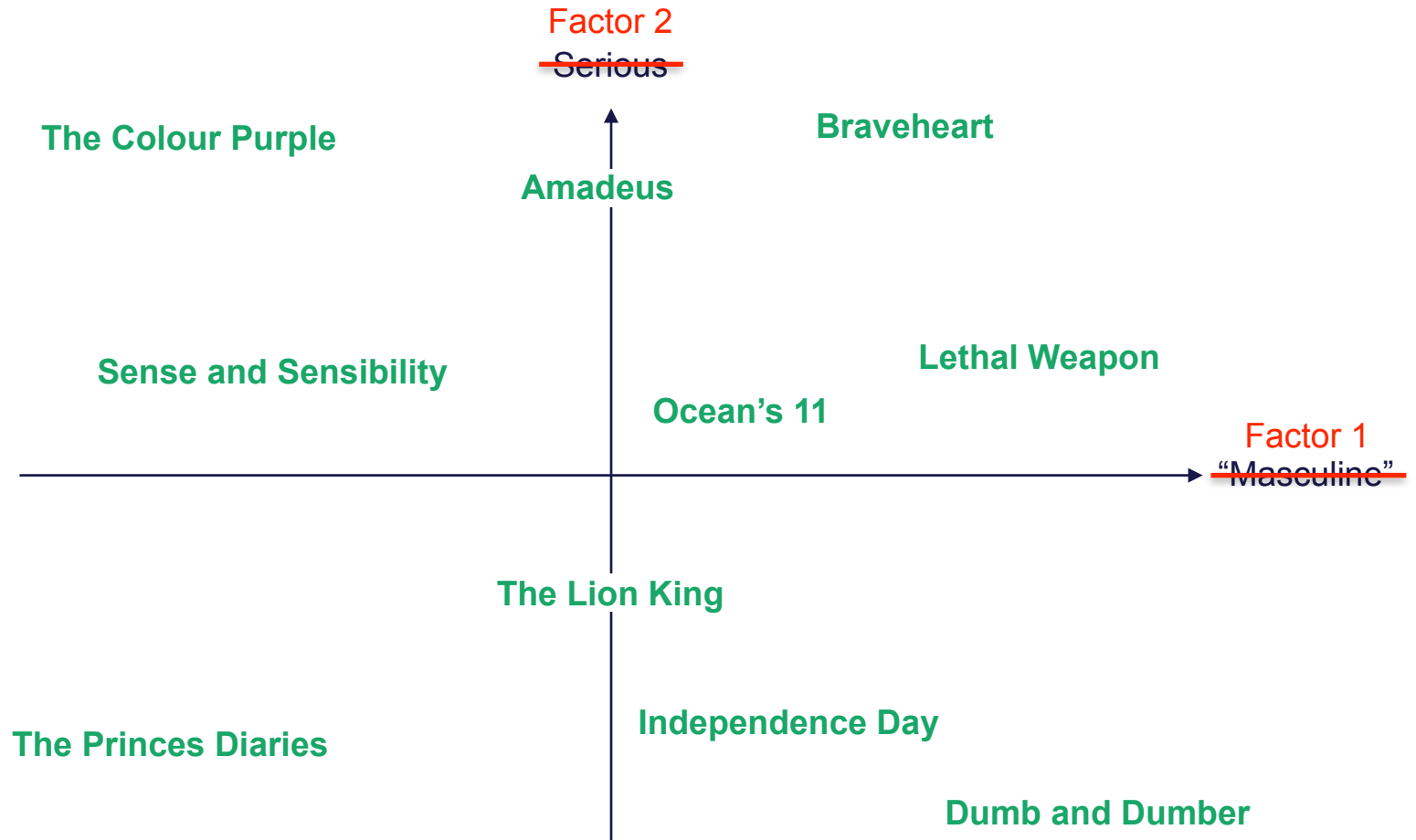


$$\hat{r}_{xi} = \mathbf{q}_i \mathbf{p}_x^* = \sum_j q_{ij} p_{xj} = -0.5 * -2 + 0.6 * 0.3 + 0.5 * 2.4 = 2.4$$

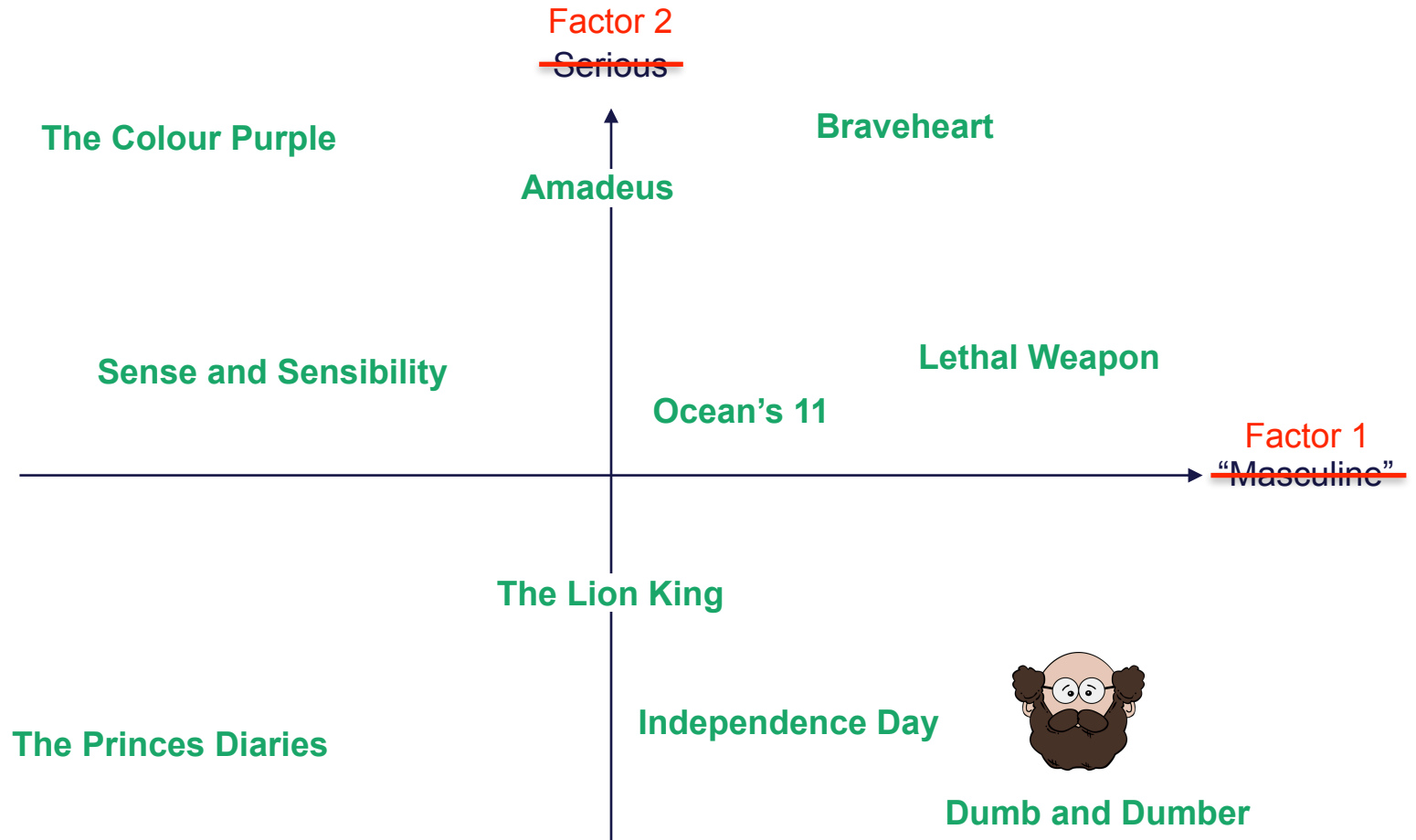
Latent Factor Models



Latent Factor Models

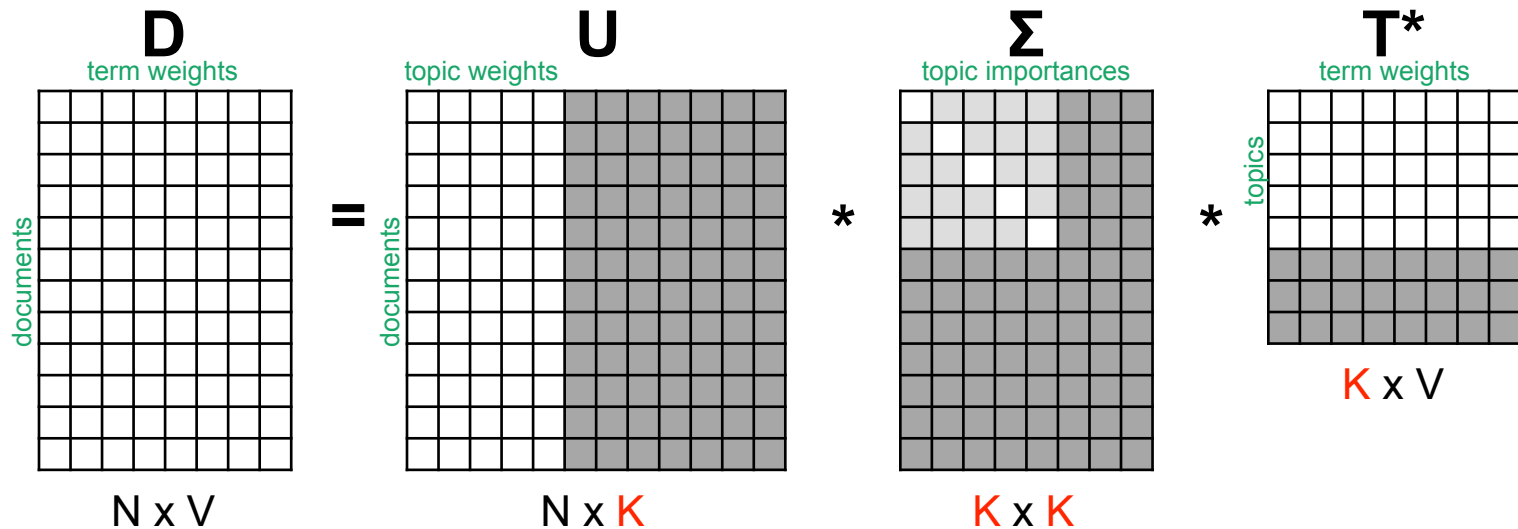


Latent Factor Models



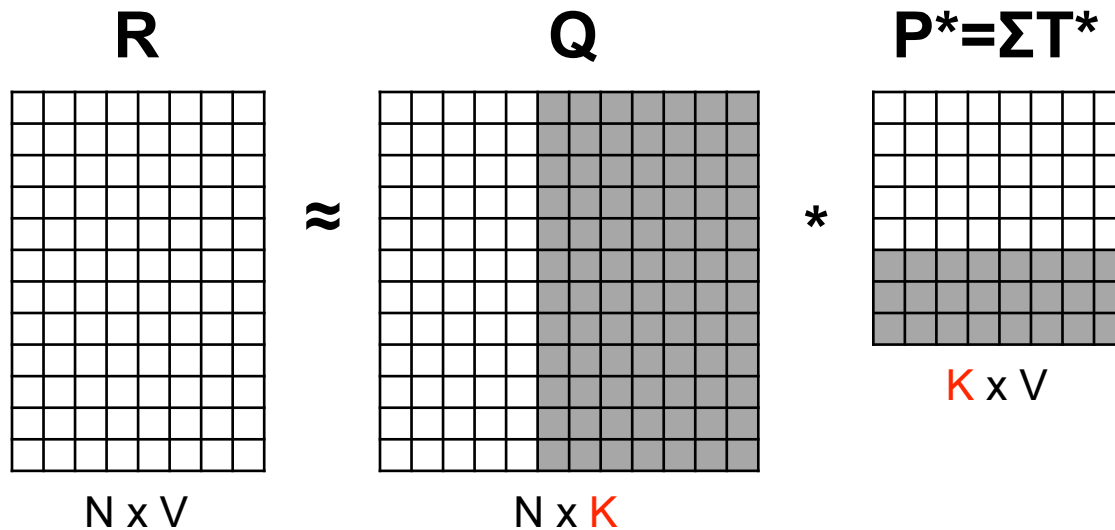
Recap: SVD

- It is a Singular Value Decomposition (generalisation of eigenvalue decomposition) of the document-term matrix $\mathbf{D} = \mathbf{U}\mathbf{\Sigma}\mathbf{T}^*$



- \mathbf{D} is a document-term matrix
- \mathbf{U} is a document-topic map ("topic distribution")
- $\mathbf{\Sigma}$ is an ordered diagonal matrix of singular values ("topic importance")
- \mathbf{T} is a term-topic map ("term distribution")

Low-Rank Matrix Factorisation



- R is a utility matrix
- Q is a item-factor map
- P^* is a user-factor map

SVD vs Low-Rank Matrix Factorisation

- **SVD gives minimum reconstruction error (Sum of Squared Errors):**

$$\min_{U, \Sigma, T} \sum_{i,j} \left(d_{ij} - [U \Sigma T^*]_{ij} \right)^2$$

which is monotonically related to RMSE

- **Complication:** The sum in SVD error term is over all entries (no-rating in interpreted as zero-rating). **But our R has missing entries!**
- **Use optimisation methods to find elements of P and Q (e.g., Gradient Descent)**

$$\min_{P, Q} \sum_{i,j \in R} \left(r_{ij} - [PQ^*]_{ij} \right)^2$$

- **Note:**
 - We don't require cols of **P**, **Q** to be orthogonal/unit length
 - **Non-convex :(**

Finding P and Q

Number of latent factors

- Loss function

$$J(P, Q) = \sum_{\text{training data } i, j \in R} \sum_{k=1}^K \left(r_{ij} - \sum_{k=1}^K q_{ik} p_{kj} \right)^2$$

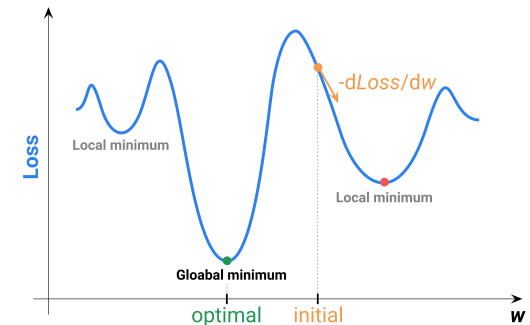
- Gradients for q_{ij} (the same for p_{ij})

$$\frac{\partial J}{\partial q_{ij}} = \sum_{\text{training data } i, j \in R} \sum_{k=1}^K -2 \left(r_{ij} - \sum_{k=1}^K q_{ik} p_{kj} \right) q_{ij}$$

- Updates (the same for p_{ij}) — until convergence

$$q_{ij} \leftarrow q_{ij} - \eta \frac{\partial J}{\partial q_{ij}}$$

Learning rate



Finding P and Q

Number of latent factors

- Loss function

$$J(P, Q) = \sum_{\text{training data } i, j \in R} \sum_{k=1}^K \left(r_{ij} - \sum_{k=1}^K q_{ik} p_{kj} \right)^2$$

- Gradients for q_{ij} (the same for p_{ij})

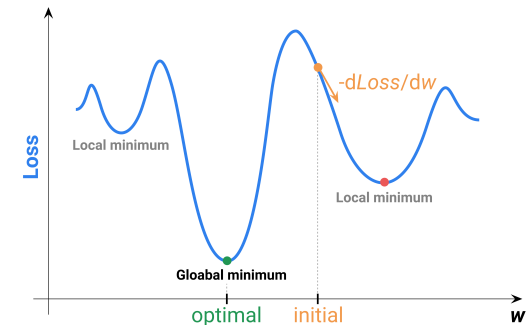
$$\frac{\partial J}{\partial q_{ij}} = \sum_{\text{training data } i, j \in R} \sum_{k=1}^K -2 \left(r_{ij} - \sum_{k=1}^K q_{ik} p_{kj} \right) q_{ij}$$

A few training data points (“batch”) — Stochastic Gradient Descent — Faster & Better in high dimensions

- Updates (the same for p_{ij}) — until convergence

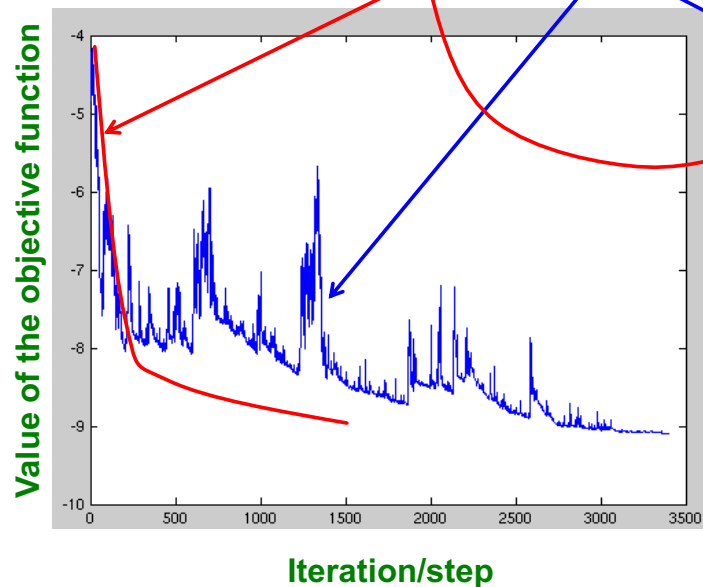
$$q_{ij} \leftarrow q_{ij} - \eta \frac{\partial J}{\partial q_{ij}}$$

Learning rate



Sidenote: GD vs. SGD

■ Convergence of **GD** vs. **SGD**



GD improves the value of the objective function at every step.

SGD improves the value but in a “noisy” way.

GD takes fewer steps to converge but each step takes much longer to compute.

In practice, **SGD** is much faster!

Finding P and Q

Number of latent factors

Hyperparameter responsible for the model's complexity (**overfitting!**)

- **Loss function**

$$J(P, Q) = \sum_{\text{training data } i, j \in R} \sum \left(r_{ij} - \sum_{k=1}^K q_{ik} p_{kj} \right)^2$$

- **Gradients for q_{ij} (the same for p_{ij})**

$$\frac{\partial J}{\partial q_{ij}} = \sum_{\text{training data } i, j \in R} \sum -2 \left(r_{ij} - \sum_{k=1}^K q_{ik} p_{kj} \right) q_{ij}$$

- **Updates (the same for p_{ij}) — until convergence**

$$q_{ij} \leftarrow q_{ij} - \eta \frac{\partial J}{\partial q_{ij}}$$

Finding P and Q

Number of latent factors

Hyperparameter responsible for the model's complexity
(**overfitting!**)

- **Loss function**

$$J(P, Q) = \sum_{\text{training data } i, j \in R} \sum \left(r_{ij} - \sum_{k=1}^K q_{ik} p_{kj} \right)^2$$

- **Gradients for q_{ij} (the same for p_{ij})**

$$\frac{\partial J}{\partial q_{ij}} = \sum_{\text{training data } i, j \in R} \sum -2 \left(r_{ij} - \sum_{k=1}^K q_{ik} p_{kj} \right) q_{ij}$$

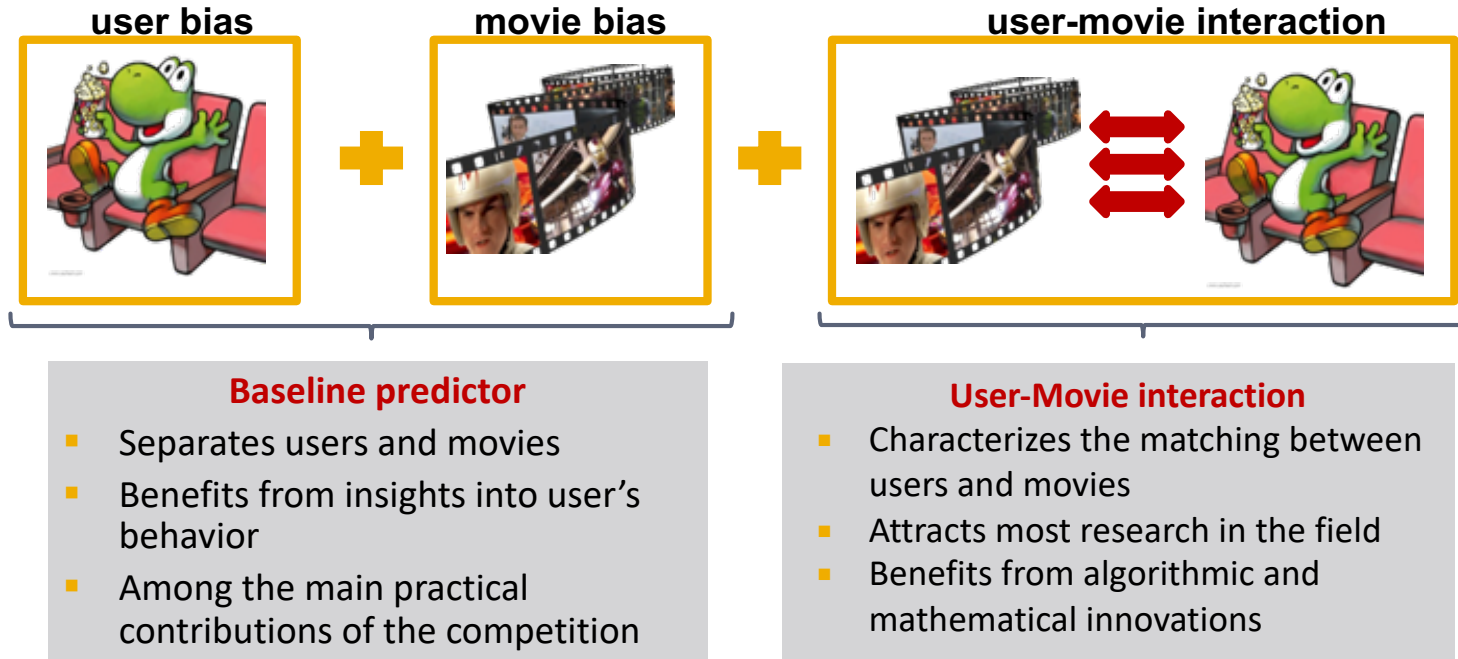
- **Updates (the same for p_{ij}) — until convergence**

$$q_{ij} \leftarrow q_{ij} - \eta \frac{\partial J}{\partial q_{ij}}$$

- regularisation
 $+ \lambda_1 \sum_{ij} q_{ij}^2 + \lambda_2 \sum_{ij} p_{ij}^2$
 (two more HPs to tune)
- direct tuning on a separate validation set / cross-validation on train data

One Last Thing...

Modelling Biases and Interactions



- μ = overall mean rating
- b_x = bias of user x
- b_i = bias of movie i

Putting It All Together

$$\hat{r}_{xi} = \underbrace{\mu}_{\text{Global mean rating}} + \underbrace{b_x}_{\text{Bias for user } x} + \underbrace{b_i}_{\text{Bias for item } i} + \underbrace{\sum_{k=1}^K q_{ik} p_{kj}}_{\text{User-Item interaction}}$$

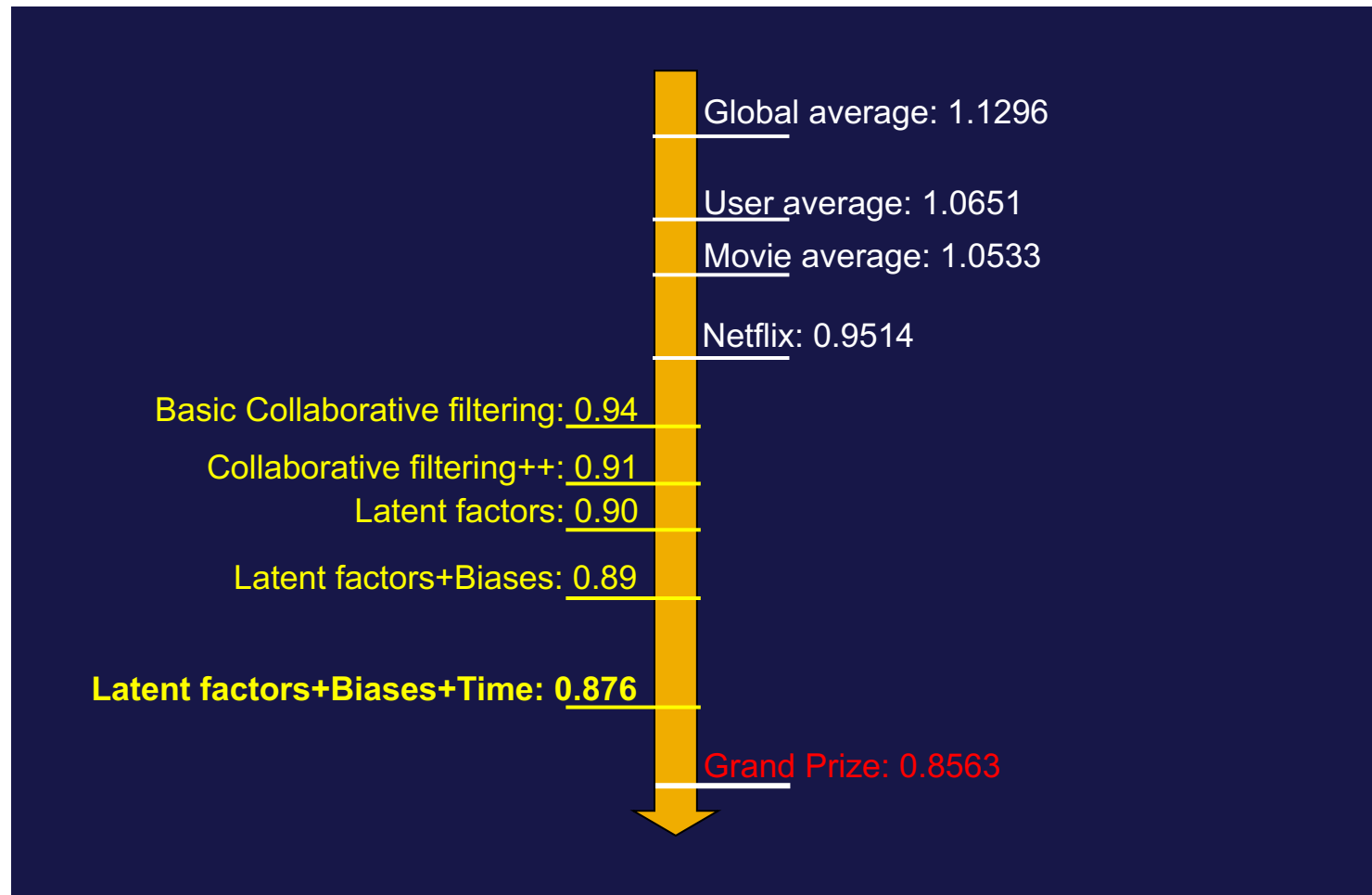
Just another optimisation parameters

$$J(P, Q) = \sum_{\text{training data}} \sum_{i, x \in R} \left(r_{ix} - (\mu + b_x + b_i) - \sum_{k=1}^K q_{ik} p_{kj} \right)^2$$

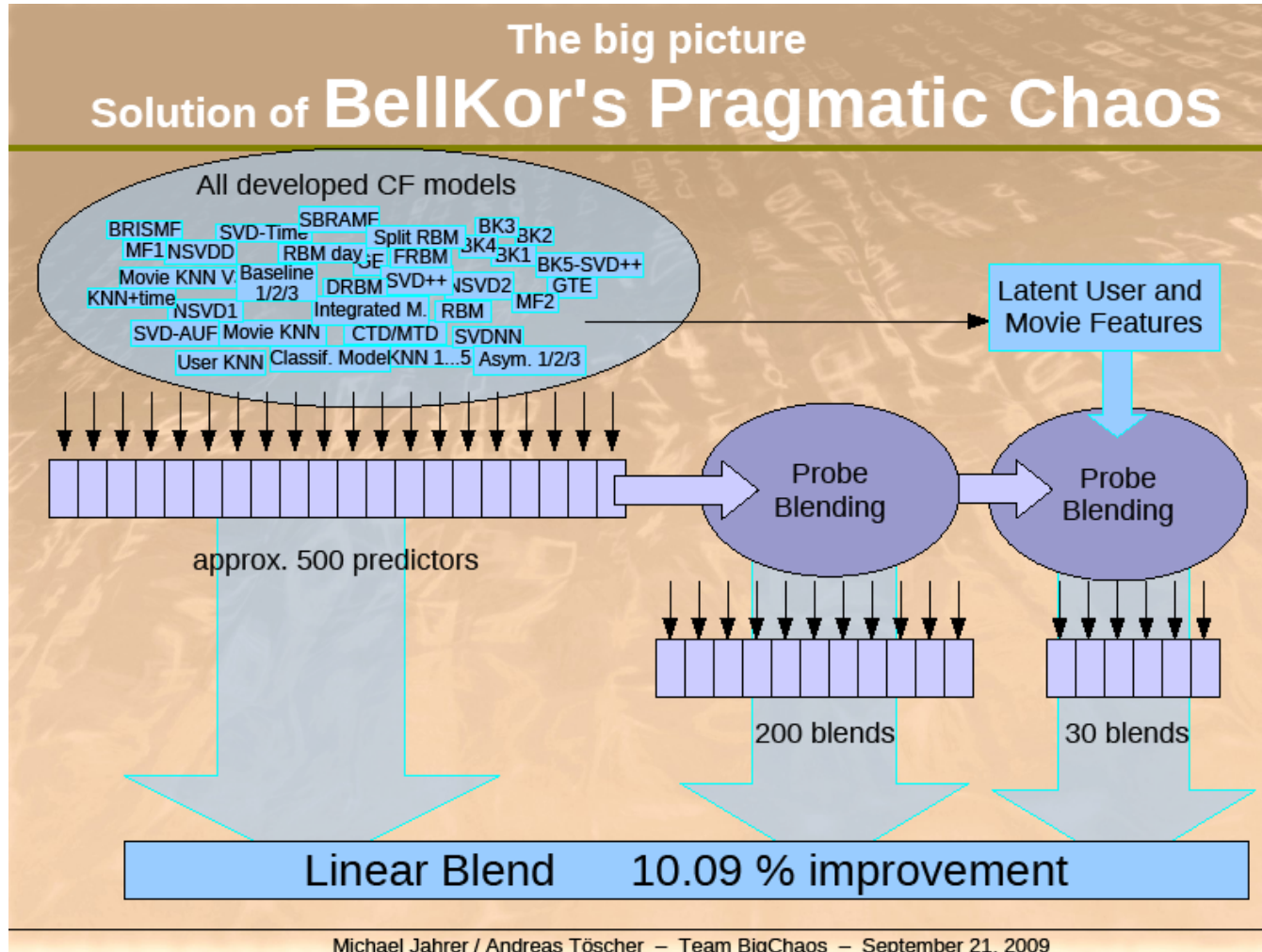
The Netflix Prize (2009)

- **Training data**
 - 100 million ratings, 480,000 users, 17,770 movies
 - 6 years of data: 2000-2005
- **Test data**
 - Last few ratings of each user (2.8 million)
 - Evaluation criterion: Root Mean Square Error (RMSE)
 - **Netflix's system RMSE: 0.9514**
- **Competition**
 - 2,700+ teams
 - **\$1 million prize** for 10% improvement on Netflix

Performance of Various Methods



Winner: A “kitchen sink” approach



Exercises

- “Recommender Systems.ipynb”: Sections 3.2 - 5

Recommender reading

- Chapter 9, “Mining of Massive Datasets” by Jure Leskovec, Anand Rajaraman, and Jeff Ullman
<http://infolab.stanford.edu/~ullman/mmds/ch9.pdf>