

Recommender Systems: Value, Methods, Measurements

Dietmar Jannach, University of Klagenfurt, Austria

dietmar.jannach@aau.at

Presented at the 12th European Big Data Management & Analytics Summer School
Padua, 2024

Outline (slides: <https://tinyurl.com/eBISS2024>)

- Recommender Systems
 - Why should we have them?
 - Value
 - How do we build them?
 - Methods
 - How do we know they work well?
 - Measurements
- What makes building them difficult?
 - Both in industry and academia

Recommender Systems

- A pervasive part of our daily online user experience
- One of the most widely used applications of machine learning

Recommended For You [Learn more](#)

Grid of recommended content:

- Bo Burnham: Words, Words, Words (TV Special 2010)
- Demetri Martin Live
- Bo Burnham
- Other titles: The Last Week in the Paradise, The Last Week in the Paradise, The Last Week in the Paradise, The Last Week in the Paradise

Navigation: [Prev 6](#) [Next 6](#)

BO BURNHAM

Bo Burnham: Words, Words, Words (TV Special 2010)

TV MA Documentary | Comedy | Music

★★★★★ 8.2/10

The internet (and soon to be movie, TV, radio, etc.) phenomenon, Bo Burnham, brings you his first one-hour stand-up special "Bo Burnham: Words, Words, Words" from the House of Blues in Boston.

[Add to Watchlist](#)

[Next »](#)

Director: Shannon Hartman
Stars: Bo Burnham

Applications

- News
- Books
- Videos
- Music
- Games
- Shopping goods
- Friends
- Groups
- Jobs
- Apps
- Restaurants
- Hotels
- Deals
- Partners
- ...
- Cigars
- Software code
- ...

Part I: Value (and some measurements)

What's their purpose and value?

- Recommenders can **create value** both for **consumers** and the **providers** of the recommendations
 - Academic research (implicitly) mostly focuses on the consumer perspective
- There can be even more **stakeholders**
 - E.g., item providers (manufacturers, property owners, artists), who may benefit from recommendations
 - They may have competing interests

Jannach, D. and Zanker, M.: "Impact and Value of Recommender Systems".

In: Recommender Systems Handbook. Ricci, F., Shapira, B. and Rokach, L. (Eds.), Springer US, 2021

Jannach, D. and Adomavicius, G.: "Recommendations with a Purpose". In: Proceedings of the 10th ACM Conference on Recommender Systems (RecSys 2016). Boston, Massachusetts, USA, 2016, pp. 7-10

Potential value for the consumer

- Examples:
 - Help users find objects that match their long-term preferences (information filtering)
 - Help users explore the item space and improve decision making
 - Make contextual recommendations, e.g.,
 - Show alternatives
 - Show accessories
 - Remind users of what they liked in the past
 - Actively notify consumers of relevant content
 - Establish group consensus

Potential value for the provider

- Examples:
 - Change **user behavior** in desired directions
 - Create additional **demand**
 - Increase (short term) **business success**
 - Enable item “**discoverability**”
 - Increase activity on the site and **user engagement**
 - Provide a valuable **add-on service**
 - **Learn more** about the customers

Multi-stakeholder considerations

- When **goals** are fully **aligned**
 - Better recommendations can lead to more satisfied, returning customers who find what they need
 - This is one implicit assumption of academic research
- When there can be a **goal conflict**
 - Not all recommendable items may have the same business value
 - From a business perspective, it may be better to recommend items with a higher margin (as long as the recommendations are still reasonable)
 - Leads to a **multi-objective** recommendation problem

Measuring the business value

- Typical quotes about value

“35% of Amazon.com’s revenue is generated by its recommendation engine.”

“We think the combined effect of personalization and recommendations save us more than \$1B per year.”

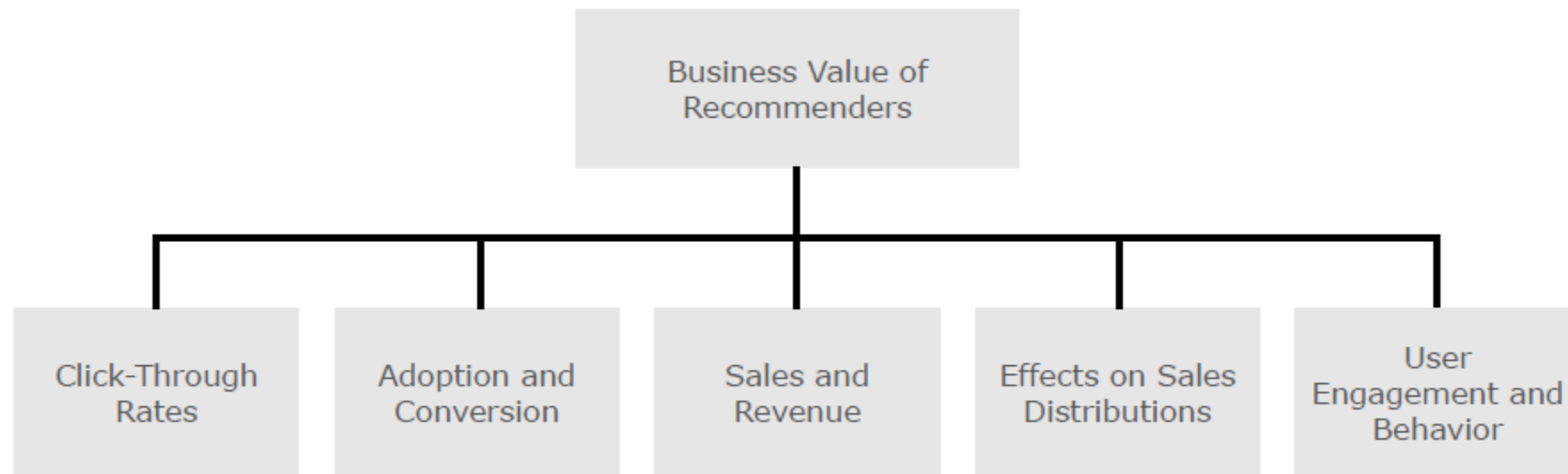
“Netflix says 80 percent of watched content is based on algorithmic recommendations”

Measuring the business value

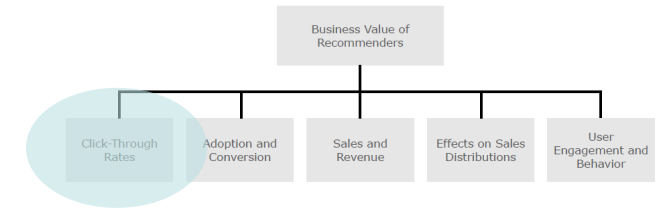
- Measuring the business value can be difficult
 - What does it tell us that 80% of the watched content comes from the recommendations (if everything is a recommendation)?
 - Where do the said savings come from?
- The used measures often largely depend on
 - The business model of the provider
 - The intended effects of the recommendations
 - Assumptions about consumer value

What is measured?

- Considering both the **impact** and **value** perspective

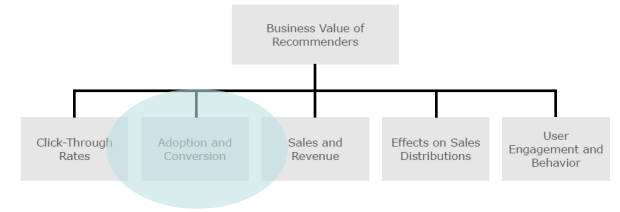


Click-Through Rates



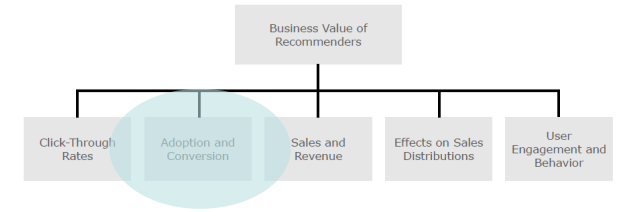
- Measures how many clicks are garnered by recommendations
 - Popular in the news recommendation domain
 - **Google News:** 38% more clicks compared to popularity-based recommendations
 - **Forbes:** 37% improvement through better algorithm compared to time-decayed popularity based method
 - **swissinfo.ch:** Similar improvements when considering only short-term navigation behavior
 - **YouTube:** Almost 200% improvement through co-visitation method (compared to popular recommendations)

Adoption and Conversion Rates



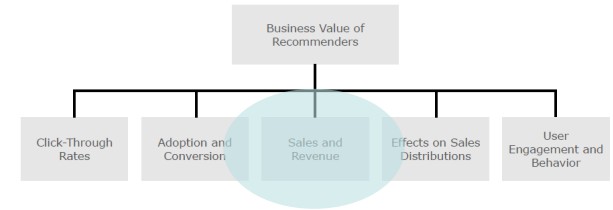
- CTR usually not the ultimate measure
 - Cannot know if users actually liked/purchased what they clicked on (consider also: click baits)
- Therefore
 - Various, domain-specific adoption measures common
- YouTube, Netflix: “Long CTR”/ “Take rate”
 - only count click if certain amount of video was watched

Adoption and Conversion Rates



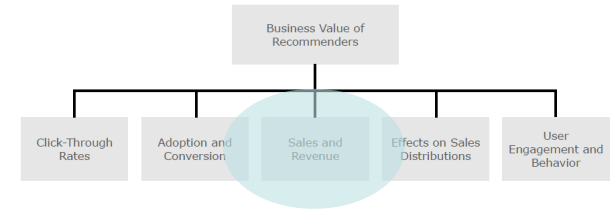
- Alternatives when items cannot be viewed/read:
- eBay:
 - “purchase-through-rate”, “bid-through-rate”
- Other:
 - LinkedIn: Contact with employer made
 - Paper recommendation: “link-through”, “cite-through”
 - E-Commerce marketplace: “click-outs”
 - Online dating: “open communications”, “positive contacts per user”

From CTR to Sales and Revenue



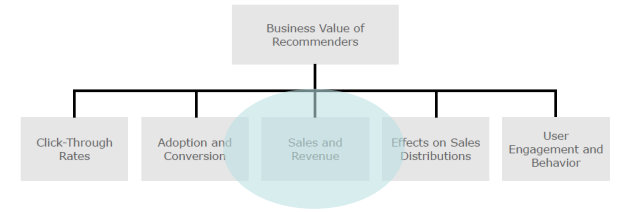
- CTR and adoption measures are good indicators of relevant recommendations
- However:
 - Often unclear how this translates into business value
 - Users might have bought an item anyway (i.e., without recommendation)
 - Substantial increases might be not relevant for business when starting from a very low basis
- Sales and Revenue figures are more direct value measurements

Sales and Revenue



- Only a few published studies, some with limitations
- Video-on-demand study:
 - 15% sales increase after introduction
 - no A/B test, could be novelty effect
- DVD retailer study:
 - 35% lift in sales when using purchased-based recommendation method compared to “no recommendations”
 - Almost no effects when recommendations were based on view statistics
 - Choice of algorithm matters a lot

Sales and Revenue



- e-grocery studies:

- 1.8 % direct increase in sales in one study
- 0.3 % direct effects in another study
- However:

- Up to 26% indirect effects, e.g., where customers were pointed to other categories in the store
- “Inspirational” effect also observed in music recommendation in our own work

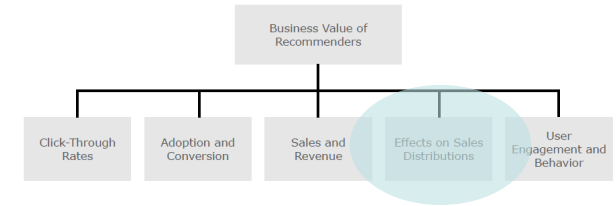
Sales and Revenue

- Book store study:
 - 28 % increase with recommender compared with “no recommender”; could be seasonal effects
 - Drop of 17 % after removing the recommender
- Mobile games (own study)
 - 3.6 % more purchases through best recommender
 - More was possible

The screenshot shows a game store interface with the following sections:

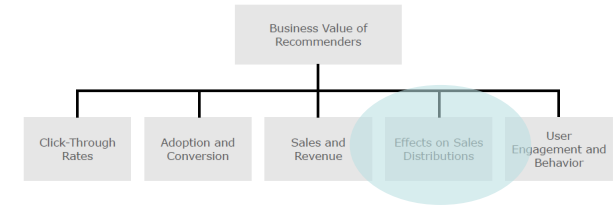
- Spiele** (Games) header with navigation links: [Suche](#) | [Hilfe](#) | [Sexy](#) | [MyGames](#)
- Navigation links: [Meine Empfehlungen](#), [Neu](#), [Top 10](#)
- Section: [Best of December](#)
- Category: [Sexy](#)
- Top Spiele** (Top Games) list:
 - [Gehirnjogging 2](#)
 - [Pizza Manager](#)
 - [Rocket Dream](#)
 - [FreeCell Deluxe For Prizes!](#)
- Meine Empfehlung** (My Recommendation) section:
 - Image of **Jewel Quest!** game box
 - Text: **Jewel Quest 2 For Prizes!**
Räum Gewinne ab!
- Top-Spiele** (Top Games) section:
 - Image of **Bubble Ducky 3in1** game box
 - Text: **Bubble Ducky 3in1**
3 spannende Knobelspiele in 1!
- Trivial Pursuit** section:
 - Image of Trivial Pursuit wheel
 - Text: **Trivial Pursuit**
Die Antwort ist "Spaß"!
- Kategorien** (Categories) list:
 - [A-Z](#)
 - [Premium & 3D](#)
 - [Ab 99 Cent](#)
 - [Action & Shooter](#)

Effects on Sales Distributions



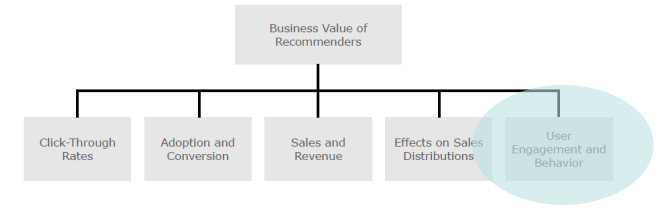
- Goal is maybe not to sell *more* but *different* items
- Influence sales behavior of customers
 - stimulate cross-sales
 - sell off on-stock items
 - promote items with higher margin
 - long-tail recommendations

Effects on Sales Distributions



- Netflix:
 - Measure the “effective catalog size”, i.e., how many items are actually (frequently) viewed
 - Recommenders lead users away from blockbusters
- Online retailer study:
 - Comparison of different algorithms on sales diversity
 - Outcomes
 - Recommenders tend to **decrease** the overall diversity
 - Might increase diversity at individual level though

User Behavior and Engagement



- Assumption:
 - Higher engagement leads to higher re-subscription rates (e.g., at Spotify)
- News domain studies:
 - 2.5 times longer sessions, more sessions when there is a recommender
- Music domain study:
 - Up to 50% more user activity
- LinkedIn:
 - More clicks on job profiles after recommender introduced

Why measuring value is difficult

- Operational aspects
 - Certain measures may actually be misleading
 - in particular CTR
 - Longitudinal effects often unclear
 - even when A/B tests last for weeks
 - A/B testing can be expensive and risky
 - Finding good “offline” proxy metrics is challenging

Why measuring value is difficult

- Strategic/organizational aspects
 - It may be actually unclear what the relevant business KPIs are
 - There may also be competing department interests
 - There may be multiple competing objectives
 - How to balance them, e.g., short-term vs long-term, consumer vs provider benefit
- What to measure then?
 - It all **depends** on the intended purpose and value of the system

The academic perspective

- In academia, we aim to
 - abstract from application specifics, e.g., business models, and
 - develop generalizable methods
- Abstract computational tasks from the literature
 - Find all or some good items
 - Predict the relevance of unseen items
 - Recommend sequence
 - Just browsing

The predominant approach

- Most common task: “Find good items”
- Most common method: “offline experimentation” and accuracy optimization
- Approach
 - Find or create a dataset that contains historical information about which recommendable items were considered “good” for individual users
 - Hide some of the information
 - Predict the hidden information
 - Measure the accuracy of the predictions

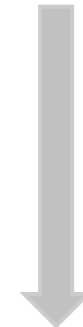
Benefits & Limitations

- Benefits of this approach
 - Well-defined problem
 - Continuous improvement (?)
 - Comparability & reproducibility
- Potential limitations
 - Being accurate is not enough, and higher accuracy not necessarily means better value for the user
 - The value for other stakeholders is not considered
 - Over-simplification of the problem

What to measure: A conceptual framework

- Should help to decide what and how to measure (both in academia and industry)
- Layered structure – strategic to operational
- Considers two viewpoints

Overarching goal of the system, strategic value
Recommendation purpose / Intended utility
System (algorithm) task
Computational metrics



Framework overview

		Consumer's Viewpoint	Provider's Viewpoint
Strategic Perspective	Overarching Goal	“Personal Utility”: Happiness, Satisfaction, Knowledge, ...	“Organizational Utility”: Profit, Revenue, Growth, ...
	Recommendation Purpose	<ul style="list-style-type: none"> • Help users find objects that match the user's long-term preferences • Show alternatives • Help users explore or understand the item space • ... 	<ul style="list-style-type: none"> • Change user behavior in desired directions • Create additional demand • Increase activity on the site • ...
Operational Perspective	System Task	<ul style="list-style-type: none"> • Annotate in context (i.e., estimate preference of a given item) • Find good items • Create diverse set of alternatives • Find suitable accessories • Retrieve novel but relevant items • ... 	
	Computational Metric	Predictive accuracy (e.g., RMSE, MAE), classification accuracy (e.g., precision, recall, AUC), ranking and top-n accuracy (e.g., rank correlation, MRR, NDCG, etc.), item “discoverability” (diversity, novelty, or serendipity measures), recommendation biases (e.g., concentration or popularity biases) and blockbuster effects, survey-based user satisfaction scores, business- and domain-specific measures (e.g., conversion rates or click-through-rates), ...	

		Consumer's Viewpoint	Provider's Viewpoint
Strategic Perspective	Overarching Goal	"Personal Utility": Happiness, Satisfaction, Knowledge, Entertainment, Benefit	"Organizational Utility": Profit, Revenue, Return on Investment, Growth, Customer Retention
	Recommendation Purpose	<ul style="list-style-type: none"> • Help users find objects that match the user's long-term preferences • Show alternatives • Help users explore or understand the item space, ... 	<ul style="list-style-type: none"> • Change user behavior in desired directions • Create additional demand • Help users discover new artists, directors, genres • Increase activity on the site • ...
Operational Perspective	System Task	<ul style="list-style-type: none"> • Annotate in context (i.e., estimate preference of a given item) • Find good items • Create diverse set of alternatives • Find mix of familiar and relevant unknown items • Find suitable accessories • ... 	
	Computational Metric	Predictive accuracy (e.g., RMSE, MAE), classification accuracy (e.g., Precision, Recall , AUC), ranking and top-n accuracy (e.g., rank correlation, MRR, NDCG, etc.), item discoverability (diversity, novelty, or serendipity measures), recommendation biases (e.g., concentration or popularity biases) and blockbuster effects, survey-based user satisfaction scores, business- and domain-specific measures (e.g., conversion rates or click-through-rates), ...	



	Consumer's Viewpoint	Provider's Viewpoint
Strategic Perspective	Overarching Goal "Personal Utility": Happiness, Satisfaction , Knowledge, Entertainment, Benefit	"Organizational Utility": Profit, Revenue, Return on Investment, Growth, Customer Retention
	Recommendation Purpose <ul style="list-style-type: none"> • Help users find objects that match the user's long-term preferences • Show alternatives • Help users explore or understand the item space, ... 	<ul style="list-style-type: none"> • Change user behavior in desired directions • Create additional demand • Help users discover new artists, directors, genres • Increase activity on the site • ...
Operational Perspective	System Task <ul style="list-style-type: none"> • Annotate in context (i.e., estimate preference of a given item) • Find good items • Create diverse set of alternatives • Find mix of familiar and relevant unknown items • Find suitable accessories • ... 	
	Computational Metric <p>Predictive accuracy (e.g., RMSE, MAE), classification accuracy (e.g., Precision, Recall, AUC), ranking and top-n accuracy (e.g., rank correlation, MRR, NDCG, etc.), item discoverability (diversity, novelty, or serendipity measures), recommendation biases (e.g., concentration or popularity biases) and blockbuster effects, survey-based user satisfaction scores, business- and domain-specific measures (e.g., conversion rates or click-through-rates), ...</p>	



	Consumer's Viewpoint	Provider's Viewpoint
Strategic Perspective	Overarching Goal "Personal Utility": Happiness, Satisfaction , Knowledge, Entertainment, Benefit	"Organizational Utility": Profit, Revenue, Return on Investment, Growth, Customer Retention
	Recommendation Purpose <ul style="list-style-type: none"> • Help users find objects that match the user's long-term preferences • Show alternatives • Help users explore or understand the item space, ... 	<ul style="list-style-type: none"> • Change user behavior in desired directions • Create additional demand • Help users discover new artists, directors, genres • Increase activity on the site • ...
Operational Perspective	System Task <ul style="list-style-type: none"> • Annotate in context (i.e., estimate preference of a given item) • Find good items • Create diverse set of alternatives • Find mix of familiar and relevant unknown items • Find suitable accessories • ... 	
	Computational Metric	Predictive accuracy (e.g., RMSE, MAE), classification accuracy (e.g., Precision, Recall , AUC), ranking and top-n accuracy (e.g., rank correlation, MRR, NDCG, etc.), item discoverability (diversity, novelty, or serendipity measures), recommendation biases (e.g., concentration or popularity biases) and blockbuster effects, survey-based user satisfaction scores , business- and domain-specific measures (e.g., conversion rates or click-through-rates), ...



		Consumer's Viewpoint	Provider's Viewpoint
Strategic Perspective	Overarching Goal	"Personal Utility": Happiness, Satisfaction, Knowledge, Entertainment, Benefit	"Organizational Utility": Profit, Revenue, Return on Investment, Growth, Customer Retention
	Recommendation Purpose	<ul style="list-style-type: none"> • Help users find objects that match the user's long-term preferences • Show alternatives • Help users explore or understand the item space, ... 	<ul style="list-style-type: none"> • Change user behavior in desired directions • Create additional demand • Help users discover new artists, directors, genres • Increase activity on the site • ...
Operational Perspective	System Task	<ul style="list-style-type: none"> • Annotate in context (i.e., estimate preference of a given item) • Find good items • Create diverse set of alternatives • Find mix of familiar and relevant unknown items • Find suitable accessories • ... 	
	Computational Metric	Predictive accuracy (e.g., RMSE, MAE), classification accuracy (e.g., Precision, Recall, AUC), ranking and top-n accuracy (e.g., rank correlation, MRR, NDCG, etc.), item discoverability (diversity, novelty, or serendipity measures), recommendation biases (e.g., concentration or popularity biases) and blockbuster effects, survey-based user satisfaction scores, business- and domain-specific measures (e.g., conversion rates or click-through-rates), ...	



		Consumer's Viewpoint	Provider's Viewpoint
Strategic Perspective	Overarching Goal	"Personal Utility": Happiness, Satisfaction, Knowledge, Entertainment, Benefit	"Organizational Utility": Profit, Revenue, Return on Investment, Growth, Customer Retention
	Recommendation Purpose	<ul style="list-style-type: none"> • Help users find objects that match the user's long-term preferences • Show alternatives • Help users explore or understand the item space, ... 	<ul style="list-style-type: none"> • Change user behavior in desired directions • Create additional demand • Help users discover new artists, directors, genres • Increase activity on the site • ...
Operational Perspective	System Task	<ul style="list-style-type: none"> • Annotate in context (i.e., estimate preference of a given item) • Find good items • Create diverse set of alternatives • Find mix of familiar and relevant unknown items • Find suitable accessories • ... 	
	Computational Metric	Predictive accuracy (e.g., RMSE, MAE), classification accuracy (e.g., Precision, Recall, AUC), ranking and top-n accuracy (e.g., rank correlation, MRR, NDCG, etc.), item discoverability (diversity, novelty, or serendipity measures), recommendation biases (e.g., concentration or popularity biases) and blockbuster effects, survey-based user satisfaction scores, business- and domain-specific measures (e.g., conversion rates or click-through-rates), . . .	



		Consumer's Viewpoint	Provider's Viewpoint
Strategic Perspective	Overarching Goal	"Personal Utility": Happiness, Satisfaction, Knowledge, Entertainment, Benefit	"Organizational Utility": Profit, Revenue, Return on Investment, Growth, Customer Retention
	Recommendation Purpose	<ul style="list-style-type: none"> • Help users find objects that match the user's long-term preferences • Show alternatives • Help users explore or understand the item space, ... 	<ul style="list-style-type: none"> • Change user behavior in desired directions • Create additional demand • Help users discover new artists, directors, genres • Increase activity on the site • ...
Operational Perspective	System Task	<ul style="list-style-type: none"> • Annotate in context (i.e., estimate preference of a given item) • Find good items • Create diverse set of alternatives • Find mix of familiar and relevant unknown items • Find suitable accessories • ... 	
	Computational Metric	Predictive accuracy (e.g., RMSE, MAE), classification accuracy (e.g., Precision, Recall, AUC), ranking and top-n accuracy (e.g., rank correlation, MRR, NDCG, etc.), item discoverability (diversity, novelty, or serendipity measures), recommendation biases (e.g., concentration or popularity biases) and blockbuster effects, survey-based user satisfaction scores, business- and domain-specific measures (e.g., conversion rates or click-through-rates), ...	



		Consumer's Viewpoint	Provider's Viewpoint
Strategic Perspective	Overarching Goal	"Personal Utility": Happiness, Satisfaction, Knowledge, Entertainment, Benefit	"Organizational Utility": Profit, Revenue, Return on Investment, Growth, Customer Retention
	Recommendation Purpose	<ul style="list-style-type: none"> • Help users find objects that match the user's long-term preferences • Show alternatives • Help users explore or understand the item space, ... 	<ul style="list-style-type: none"> • Change user behavior in desired directions • Create additional demand • Help users discover new artists, directors, genres • Increase activity on the site • ...
Operational Perspective	System Task	<ul style="list-style-type: none"> • Annotate in context (i.e., estimate preference of a given item) • Find good items • Create diverse set of alternatives • Find mix of familiar and relevant unknown items • Find suitable accessories • ... 	
	Computational Metric	Predictive accuracy (e.g., RMSE, MAE), classification accuracy (e.g., Precision, Recall, AUC), ranking and top-n accuracy (e.g., rank correlation, MRR, NDCG, etc.), item discoverability (diversity, novelty, or serendipity measures), recommendation biases (e.g., concentration or popularity biases) and blockbuster effects, survey-based user satisfaction scores, business- and domain-specific measures (e.g., conversion rates or click-through-rates), ...	



		Consumer's Viewpoint	Provider's Viewpoint
Strategic Perspective	Overarching Goal	"Personal Utility": Happiness, Satisfaction, Knowledge, Entertainment, Benefit	"Organizational Utility": Profit, Revenue, Return on Investment, Growth, Customer Retention
	Recommendation Purpose	<ul style="list-style-type: none"> • Help users find objects that match the user's long-term preferences • Show alternatives • Help users explore or understand the item space, ... 	<ul style="list-style-type: none"> • Change user behavior in desired directions • Create additional demand • Help users discover new artists, directors, genres • Increase activity on the site • ...
Operational Perspective	System Task	<ul style="list-style-type: none"> • Annotate in context (i.e., estimate preference of a given item) • Find good items • Create diverse set of alternatives • Find mix of familiar and relevant unknown items • Find suitable accessories • ... 	
	Computational Metric	Predictive accuracy (e.g., RMSE, MAE), classification accuracy (e.g., Precision, Recall, AUC), ranking and top-n accuracy (e.g., rank correlation, MRR, NDCG, etc.), item discoverability (diversity, novelty, or serendipity measures), recommendation biases (e.g., concentration or popularity biases) and blockbuster effects, survey-based user satisfaction scores, business- and domain-specific measures (e.g., conversion rates or click-through-rates), ...?	



Summary of first part

- Demonstrated business value of recommenders in many domains
- Size of impact however depends on many factors like baselines, domain specifics etc.
- Measuring impact and value is generally not trivial
 - It all depends on the purpose
 - Choice of the evaluation measure matters a lot
 - CTR can be misleading
- “Metric-Task-Purpose-Fit” to be considered

Part II: Methods

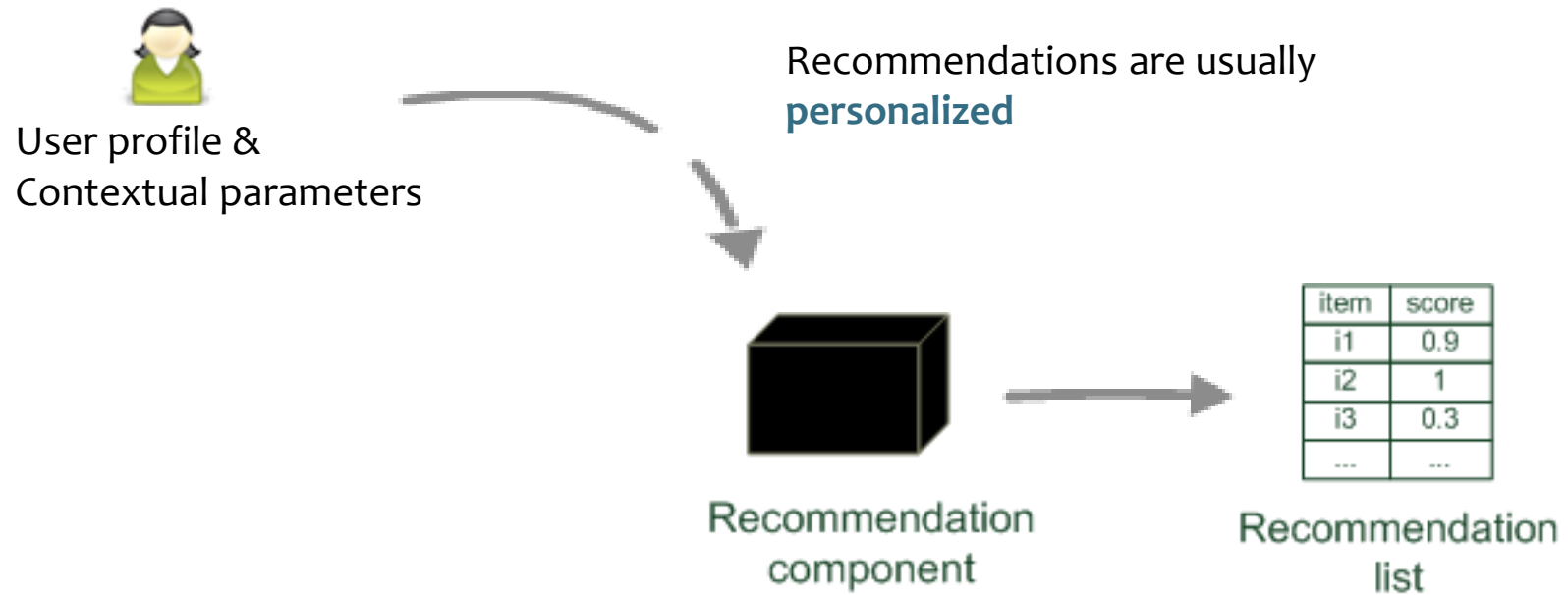
History and roots

- Roots in various fields
 - e.g., Information Retrieval, Machine Learning, Human Computer Interaction
- Their design can furthermore be influenced by insights from more distant fields
 - e.g., Consumer behavior, Psychology, Marketing
- Typical goals:
 - Avoid information overload (filtering)
 - Active promotion of content
- Personalization often as a central concept

A common categorization, based on the used information

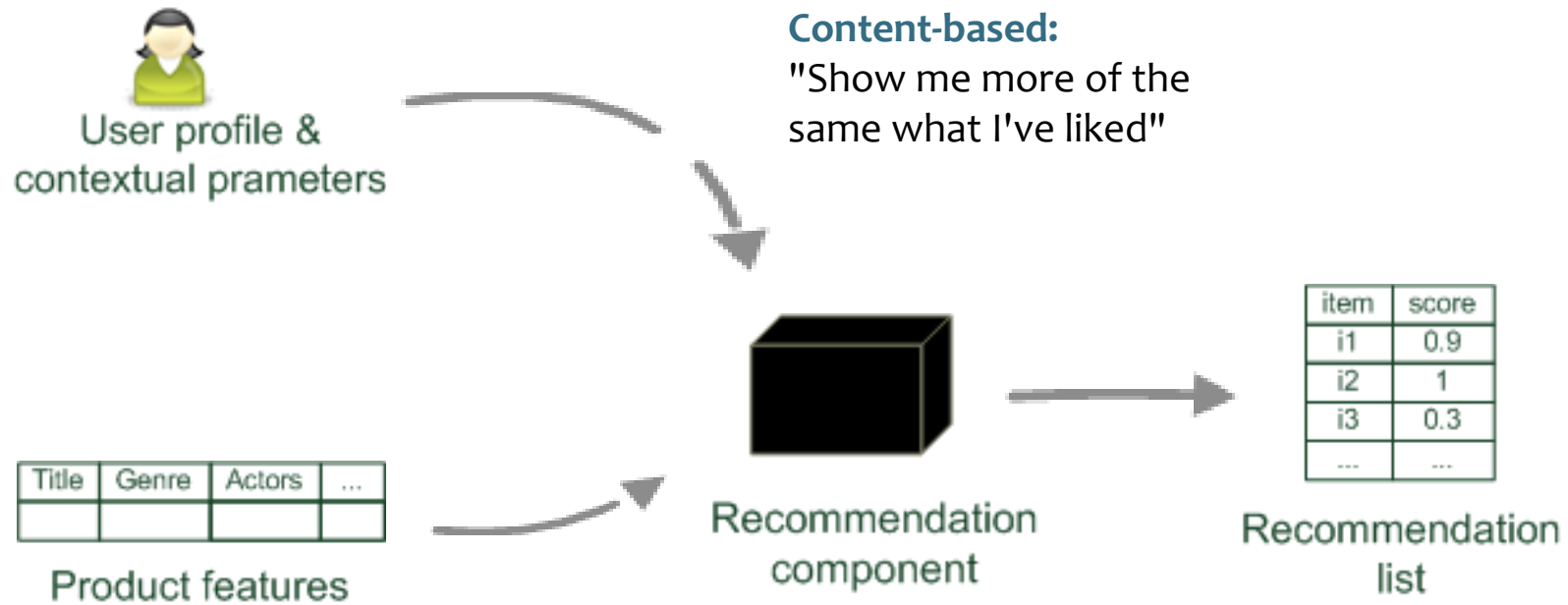
- Content-based Filtering
 - Collaborative Filtering
 - Hybrid Systems
 - Knowledge-based Systems
-
- Interactive and Conversational Recommendation

The Basic Principle (simplified)*



*One-shot scenario, no further interaction, no explanations etc.

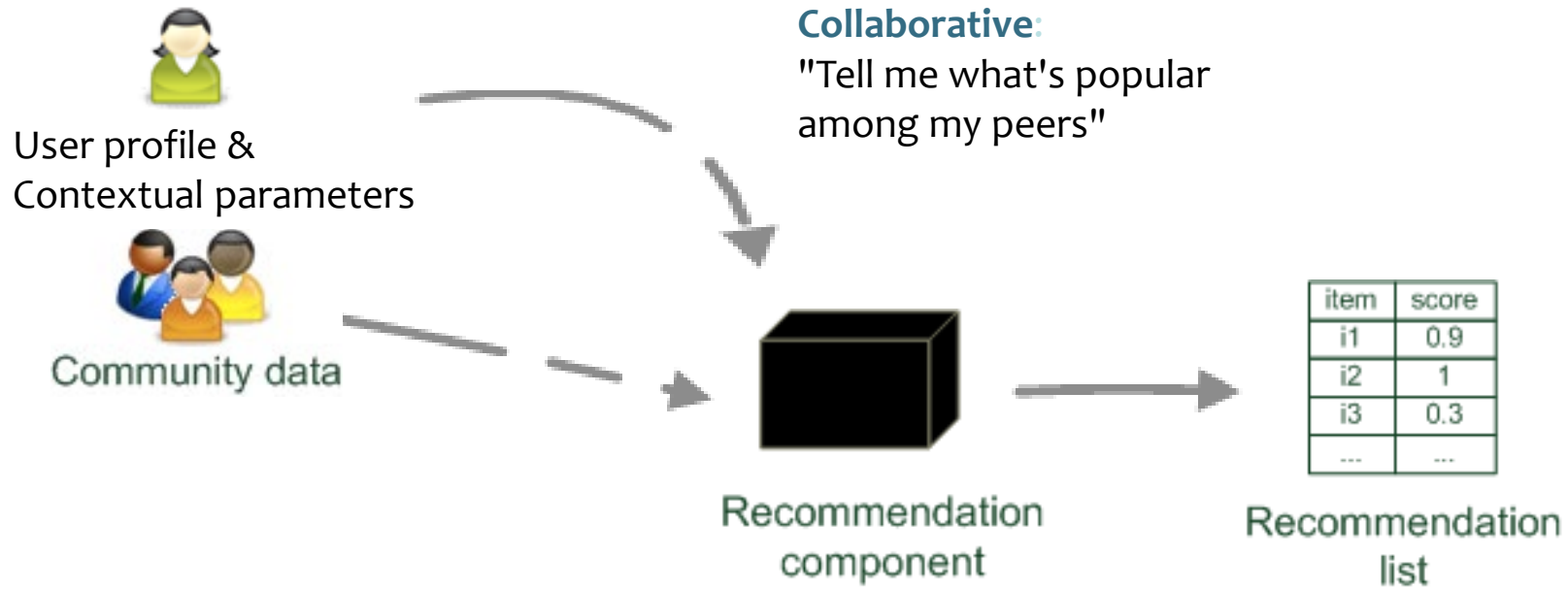
Content-based Filtering



Content-based Filtering

- Basic approach
 - Automated construction of a user profile based on characteristics of items the user previously liked/consumed/etc.
 - Various types of item information considered in the literature
 - Recommending: Matching of user profile with item profiles
 - By design leads to “more of the same” recommendations and limited discovery
- Long history since the 1960s
 - Similar to personalized information retrieval, the term “content-based” stems from early use cases such as news recommendation.

Collaborative Filtering

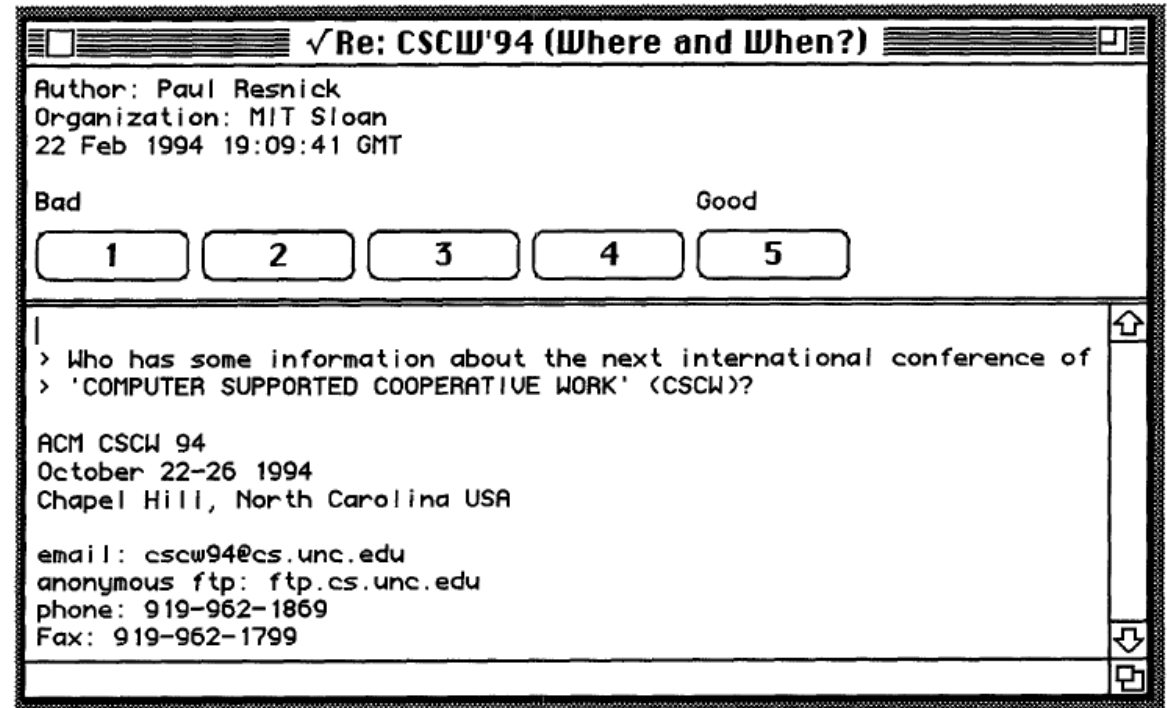


Collaborative Filtering

- Basic Approach
 - Leverage the opinion or observed behavior of other users when recommending
 - Assumption (simplified): When users agreed in their opinion in the past, they will probably agree in the future
- History since early 1990s
 - Term introduced in 1992 with experimental “Tapestry” system
 - “GroupLens” system and related systems proposed in 1994
 - Fully automated prediction of user preferences

Collaborative Filtering in GroupLens

- The GroupLens system
 - User-item ratings as the only input
 - Recommendations based on nearest-neighbor approach
 - Original paper proposed a **system**, and not only an algorithm



Matrix Completion

- Recommendation considered as matrix completion (“matrix filling”) problem

	Item1	Item2	Item3	Item4	Item5
Alice	5	?	4	4	?
User1	3	?	2	3	?
User2	?	3	4	?	?
User3	?	3	1	?	4
User4	1	?	5	2	1

- Items are recommended based on predicted ratings

GroupLens: User-based K-nearest-neighbors (kNN)

- Given an "active user" (Alice) and an item I not yet seen by Alice
 - The *goal is to estimate Alice's rating for this item*, e.g., by
 - find a set of users (peers) who liked the same items as Alice in the past **and** who have rated item I
 - use, e.g., the average of their ratings to predict if Alice will like item I
 - do this for all items Alice has not seen and recommend the best-rated
 - Assessing the predictions without users (offline):
 - Hide & predict some ratings, compute the average prediction error, e.g., based on the Root Mean Square Error

$$RMSE = \sqrt{\frac{\sum_{(u,i) \in K} (\hat{r}_{ui} - r_{ui})^2}{|K|}}$$

Collaborative Filtering (CF): A Success Story

- 1998:
 - Dimensionality reduction for CF, clustering
 - Collaborative/Content-based Hybrids
- 1999: It works in e-commerce!
 - First reports on successful applications in practice (e-commerce, music, video)
- 2000: Item-to-item collaborative filtering
- 2003: Amazon.com
 - Report on the successful use of recommendations at Amazon.com using item-to-item filtering

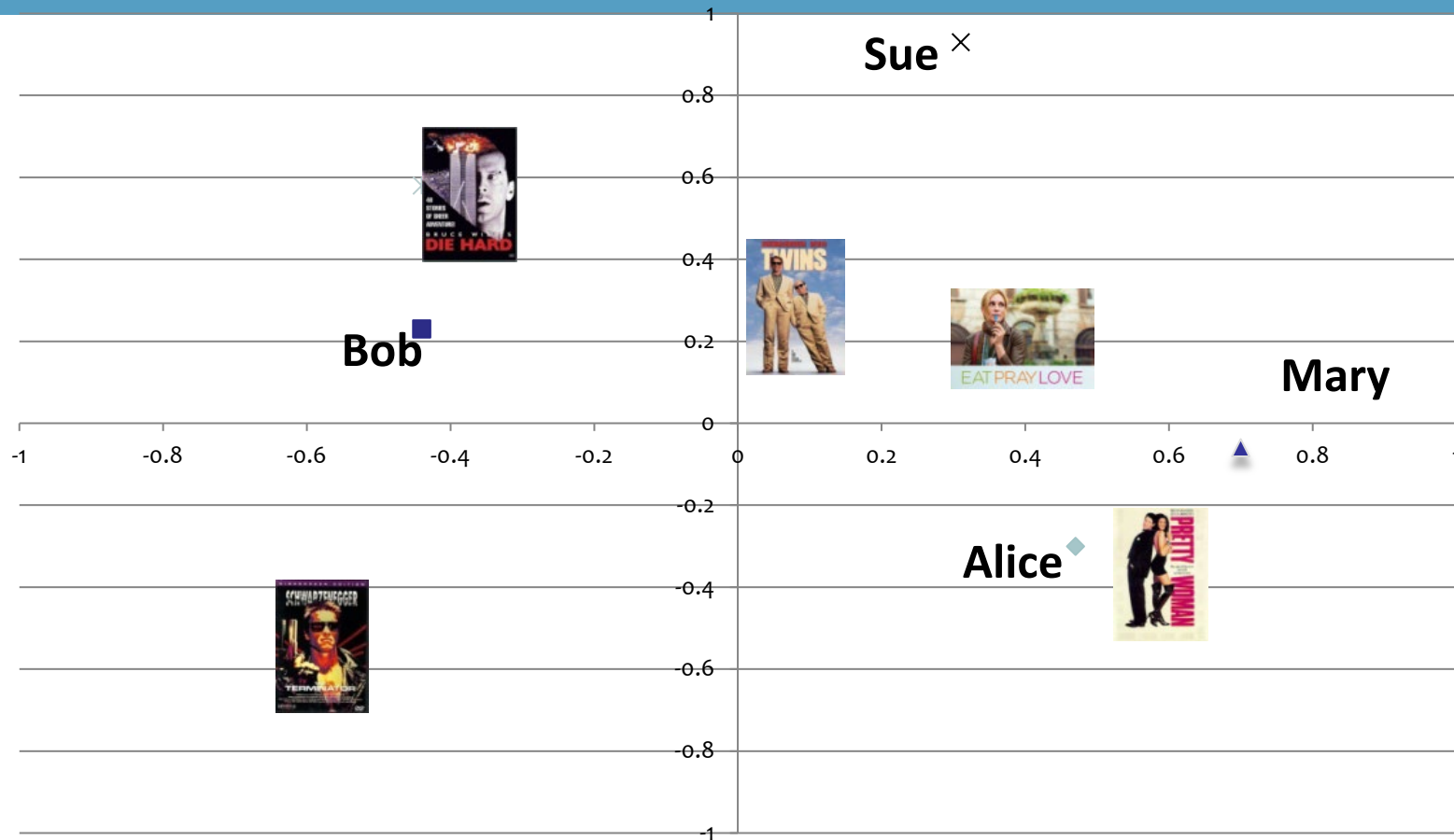
The Netflix Prize (2006-2009)

- Netflix announced a 1 million dollar prize in 2006
 - For beating their system by 10% in terms of the prediction error
 - Provided at that time huge dataset
- Effects
 - Further boosted research on matrix completion
- Contest ended in 2009, some winning ingredients:
 - Matrix factorization, ensemble methods

Matrix Factorization

- **2000:** Early experiments with Singular Value Decomposition
 - Use SVD for dimensionality reduction
 - Capture the most important factors/aspects in the data
 - Should also help to reduce noise
- **2006 and later:** MF variants using, e.g., gradient descent optimization

Projection into lower-dimensional spaces



- Nowadays called “embeddings”

Post-Netflix-Prize Developments

- Rating prediction increasingly considered **irrelevant** in practice
 - Item **relevance prediction** still important
- Various ranking-based methods (“**learning-to-rank**”) proposed around 2009
- More focus on situations where only implicit feedback is available
- Probably hundreds of CF algorithms per year
- During the last few years, **deep learning** techniques dominate the landscape

Recommendation as a Supervised ML Problem

- While the techniques change over time, the underlying machine learning (ML) problem formulation remained mostly constant, where the goal is to (Adomavicius & Tuzhilin, 2005)
 - learn a function from noisy data, which
 - given a user u and an item i ,
 - predicts the relevance of *item i for user u .*
- Learning-to-rank models may work a bit differently, but ultimately also predict relevance scores for a given set of items

Matrix Completion - Limitations



- Amazon's contextual recommendations are a guiding scenario in the literature
 - But there are no ratings
 - There apparently is not even personalization

Sequence-aware Recommenders

- An alternative problem abstraction
 - Aims to address different various real-world application problems
 - Input is not a rating matrix, but a **sequential log** of recorded user interactions
 - Item views, purchases, listening events
 - Most common problem is to predict items that are relevant in the user's **ongoing session**
 - Often, users are anonymous and the user's intent must be guessed from a small set of interactions (“**session-based** recommendation”)

Session-based Recommendation

- Guessing the intention can be difficult



The image shows a product listing for a Minnow Sports Aluminum Baseball Bat. On the left, there is a vertical strip of six small thumbnail images. The main image shows a silver aluminum baseball bat with a black handle. The bat has 'MINNOW' written on it. To the right of the bat, the text reads 'MINNOW SPORTS Baseball Bat' with a fish logo above and a baseball logo below. Below the bat, it says '32" ▶ 24 oz'. Underneath that, it says 'Roll over image to zoom in'. To the right of the bat, the text reads 'Minnow Sports', 'Minnow Sports Aluminum Baseball Bat For Baseball & Teeball', '★★★★☆ 8 customer reviews', 'Price: ~~\$29.99~~', 'Sale: \$19.99', 'You Save: \$10.00 (33%)', 'In Stock.', 'This item does not ship to Germany. Please check other sellers who may ship internationally. Learn more', 'Sold by BBro Store and Fulfilled by Amazon. Gift-wrap available.', 'Item Display Length: 32.0 inches', and a list of bullet points: '• Made from lightweight high grade Aluminum alloy for faster swing speed', '• Ultra-thin 32" handle with All Sports grip for increased stability and accuracy', '• Stylish design featuring full rolled-over end for ultimate performance', '• Ideal for all levels of baseball players from practice to matches', '• 32 inches in length & 24 ounces'.



Session-based Recommendation

- Also in online music recommendation
- Our user searched and listened to “Last Christmas” by Wham!
- Should we, ...
 - Play more songs by Wham!?
 - More pop Christmas songs?
 - More popular songs from the 1980s?
 - Play more songs with controversial user feedback?



Session-aware Recommendation

- In some domains, **past sessions of the current user** are also known,
 - potential for personalization
 - possibility to remind users of objects
- We call this problem “**session-aware**” recommendation
- One main problem is to effectively combine long-term and short-term preference models

Long-term and short-term models

- Being able to predict which kinds of things a certain user **generally** likes, is important
- Here's what the customer looked at or purchased during the last weeks



- Now, he or she return to the shop and browse these items

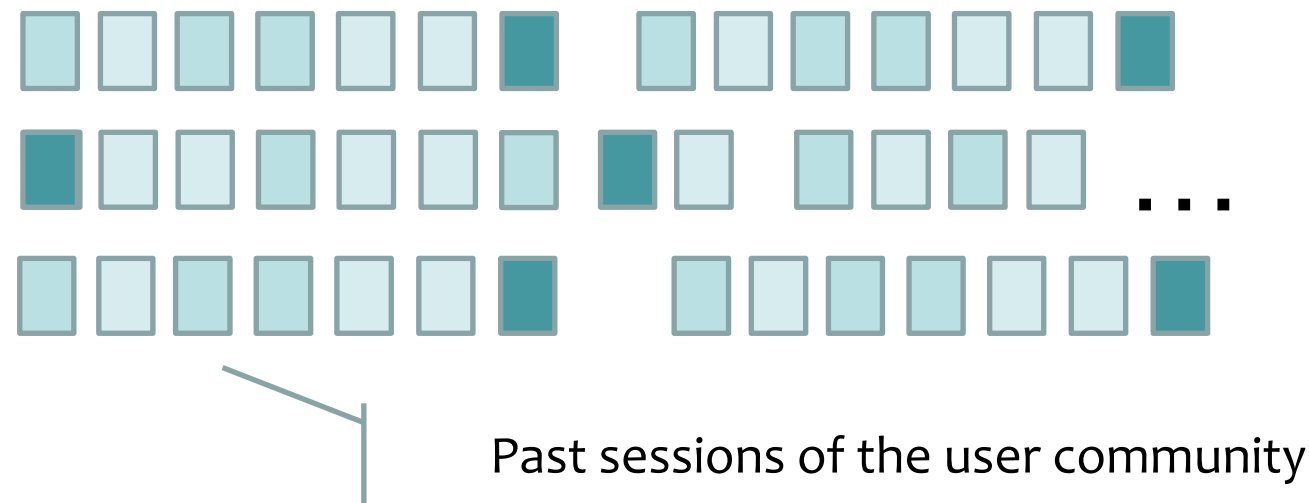
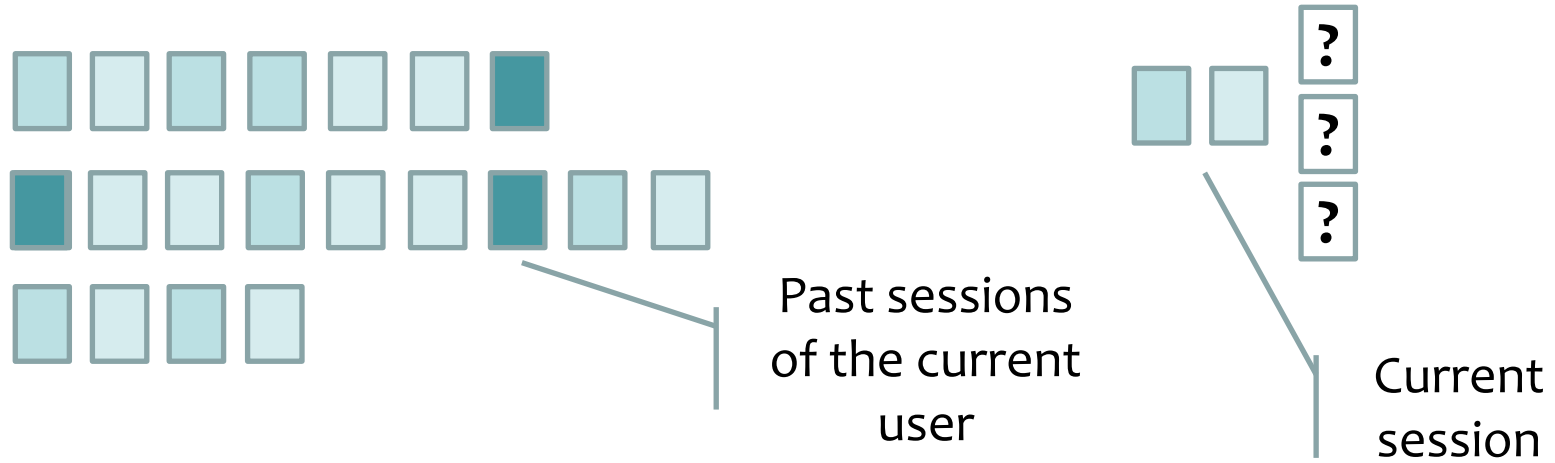


What to recommend?

- Some plausible options
 - Only shoes or only watches?
 - Mostly Nike shoes?
 - Maybe also some T-shirts?
- Considerations and observations
 - Using the matrix completion formulation, the system will mostly recommend T-shirts and trousers
 - Research indicates that both models are relevant, but that the short-term model is much more important



A Problem Abstraction



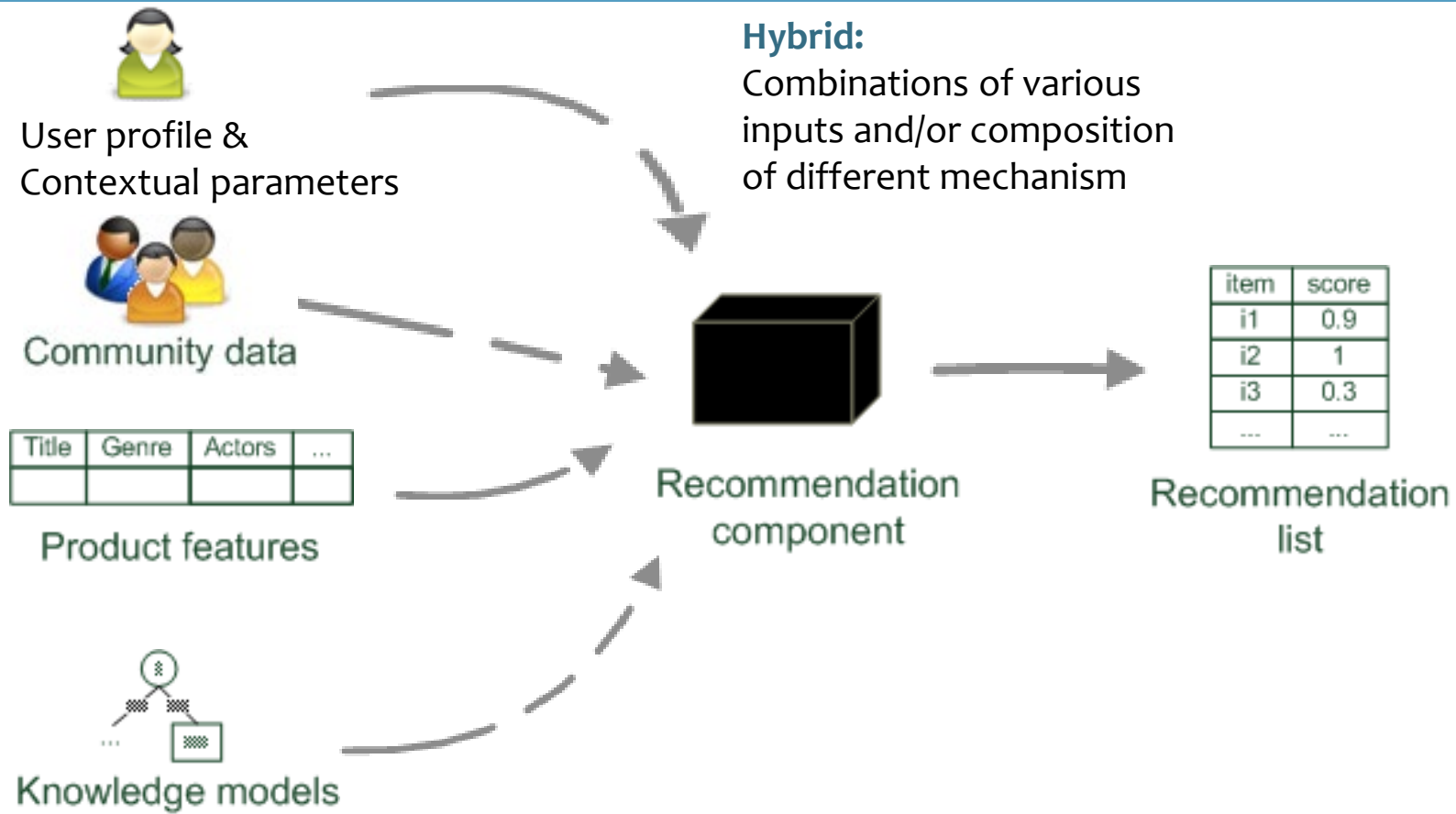
Technical Approaches

- Basic techniques
 - Item co-occurrences: “Customers who bought ... also bought”
 - Markov Chains and Sequential Rules
- Nearest neighbors
 - Find past sessions that are similar to the current (ongoing) one, predict items from neighbor sessions
- Sequence learning / modeling
 - Embeddings, Recurrent Neural Networks, Attention/Transformer

Applications and History

- Early applications for next-page prediction in web browsing
- Next-track music recommendations and automated radio stations, video playlists
- Next-POI recommendation in travel and tourism applications
- E-commerce applications, increasingly since 2015
 - In particular many neural methods proposed recently
 - Publicly available datasets

Hybrid Recommendation Approaches



Hybridization Designs

- Various forms proposed in the literature
 - Monolithic exploiting different features
 - E.g., Combining different signals in one system
 - Parallel use of several systems
 - E.g., switching, based on recorded user interactions
 - Pipelined invocation of different systems
 - E.g., use one recommender for filtering items and another one for ranking what remains

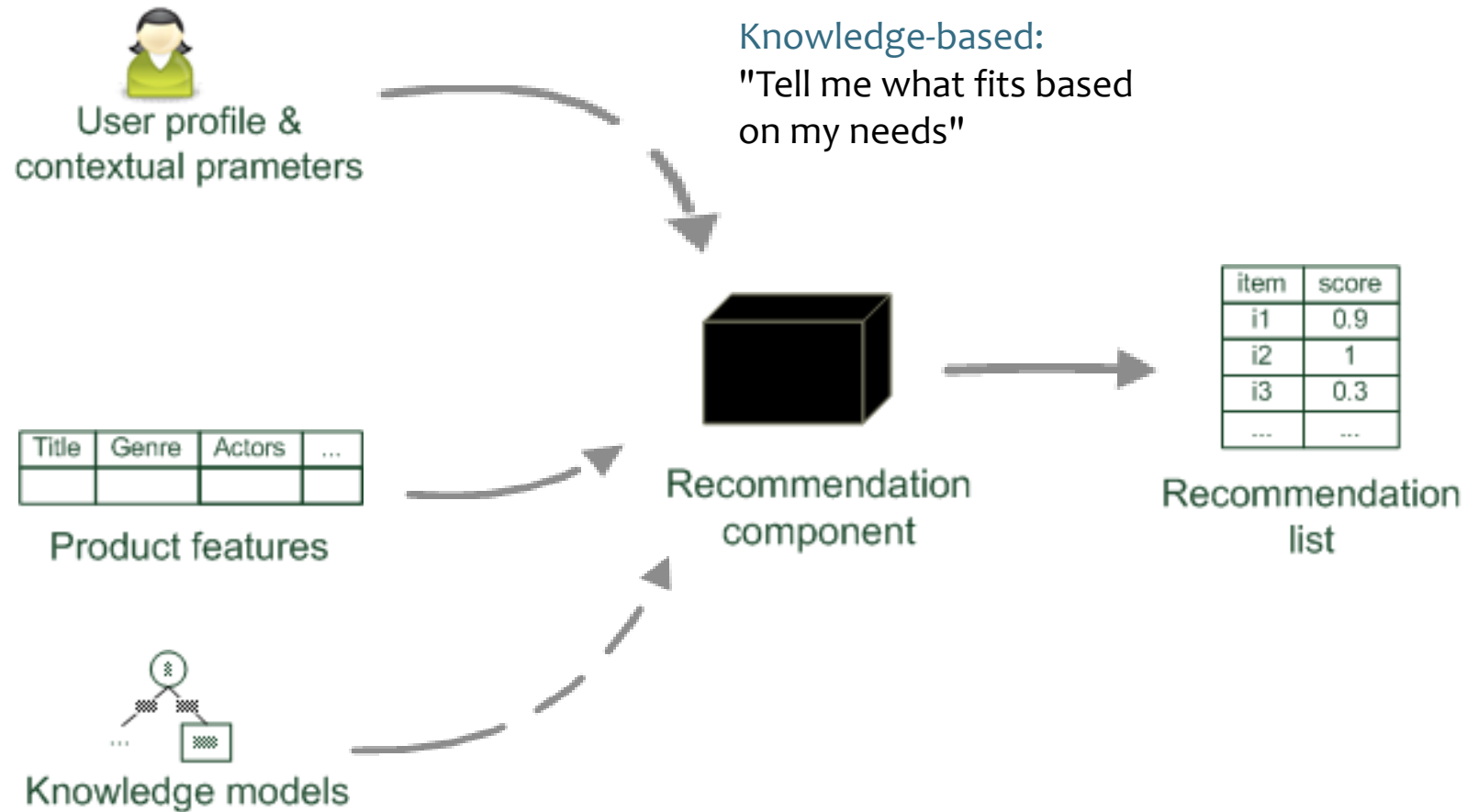
Robin Burke. 2002. Hybrid Recommender Systems: Survey and Experiments. *User Modeling and User-Adapted Interaction* 12, 4 (November 2002), 331–370.

Jannach, D., Zanker, M., Felfernig, A. and Friedrich, G.: "Recommender Systems - An Introduction". Cambridge University Press, 2010.

Collaborative Filtering with Side Information

- Pure content-based techniques are rarely used for recommendation
 - They are limited to finding similar items
 - Content encodings (e.g., TFIDF, embeddings) tell us little about the general quality of the items
 - Recommendations can be obscure or too niche
- Very common, however:
 - Leverage information about items or users in combination with collaborative filtering approaches

Knowledge-based Systems



Knowledge-based Systems

- Explicitly encode recommendation knowledge
- Usually no learning, but knowledge engineering
- Used for certain application domains, e.g.,
 - One-time investments and decisions
 - Domains where technical constraints have to be considered
 - **Interactive/conversational** recommendations, chat bots

Is this even a recommender?

The screenshot shows a web browser window with the URL `http://www.configworks-gmbh.online.de - VIBE - the virtual adviser for the Warmbad-Villach spa reso...`. The page header includes the VIBE logo and the text "VIBE VIRTUAL ADVISER". Navigation links for "HOME", "CALL BACK SERVICE", and "RECOMMENDATIONS" are visible.

The main content area features a woman in a red dress on the left, with a speech bubble that reads: "Think about what you'd really like and I'll see what I can come up with for you." To her right, a question asks: "Mr Jannach, how do you feel right now? What would you like to improve if it were possible?". Below this question is a list of six options, each with a checkbox:

- I feel quite tired and would like to recharge my batteries
- I would like to improve my fitness.
- I would like to lose some weight and be slimmer.
- I often feel tense and sometimes have problems with my back.
- I would like to do something about my appearance and my image.
- I feel perfectly healthy and would simply like to relax for a few days.


At the bottom of the form, there are three buttons: "Direct to result", "Back", and "Next". The status bar at the bottom left shows the word "Fertig" and a green checkmark icon on the right.

Is this even a recommender?

http://www.configworks-gmbh.online.de - VIBE - the virtual adviser for the Warmbad-Villach spa reso... HOME CALL BACK SERVICE RECOMMENDATIONS

VIBE
VIRTUAL ADVISER

Did you know that...



Wonderful, we've now got to your final selection. Here's my recommendation for you ...

Feel well week

Length of stay:	per week (7 nights) per person
Meals:	Half board
Accommodation:	The Warmbaderhof
Dates:	At any season
Rate in single room:	from € 1595
Rate in double room:	from € 1595

[Details](#)
[Why?](#)

I can also recommend the following packages:

- You can book a personal massage or a whole massage programme for your stay at any time.

Golf & Spa

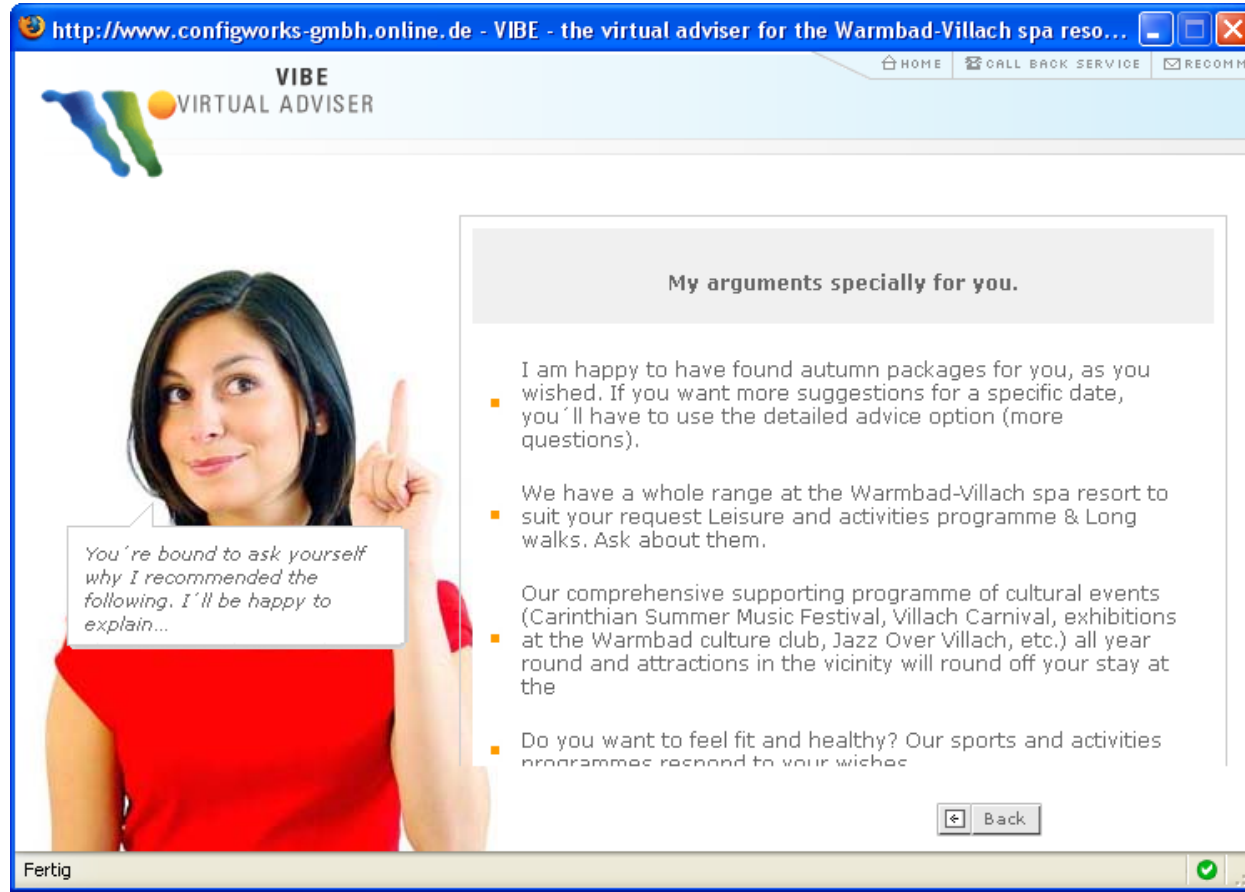
Length of stay:	per week (7 nights) per person
Meals:	Half board
Accommodation:	The Warmbaderhof
Dates:	01.04.2008-31.10.2008

[Details](#)
[Why?](#)

Back Restart Print Online-request

Fertig

Is this even a recommender?



http://www.configworks-gmbh.online.de - VIBE - the virtual adviser for the Warmbad-Villach spa reso...

VIBE
VIRTUAL ADVISER

HOME CALL BACK SERVICE RECOMMENDATIONS

You're bound to ask yourself why I recommended the following. I'll be happy to explain...

My arguments specially for you.

- I am happy to have found autumn packages for you, as you wished. If you want more suggestions for a specific date, you'll have to use the detailed advice option (more questions).
- We have a whole range at the Warmbad-Villach spa resort to suit your request Leisure and activities programme & Long walks. Ask about them.
- Our comprehensive supporting programme of cultural events (Carinthian Summer Music Festival, Villach Carnival, exhibitions at the Warmbad culture club, Jazz Over Villach, etc.) all year round and attractions in the vicinity will round off your stay at the
- Do you want to feel fit and healthy? Our sports and activities programmes respond to your wishes

Back

Fertig

A common categorization, based on the used information

- Content-based Filtering
 - Collaborative Filtering
 - Hybrid Systems
 - Knowledge-based Systems
-

- **Interactive and Conversational Recommendation**

From Algorithms to User Experience

- Most academic research focuses on algorithmic aspects
 - e.g., learning to predict / “post-dict” hidden ratings
- But a recommender *system* is more than the algorithm, see later lectures
- The UI can have a huge impact on adoption
 - Garcin et al., for example, report a more than 100% increase in the CTR when changing the position of the recommendations

Konstan, J.A. & Riedl, J.. “Recommender systems: from algorithms to user experience”
User Model User-Adap Inter (2012) 22: 101.

Garcin, F., Faltings, B., Donatsch, O., Alazzawi, A., Bruttin, C., and Huber, A. 2014. Offline and online evaluation of news recommender systems at swissinfo.ch. In *Proceedings of the 8th ACM Conference on Recommender systems (RecSys '14)*.

Interactive Recommender Systems

- But: A common assumption in many research works: Which user interaction?
 - The system monitors what I do
 - And then shows me stuff
 - Which I can click on

Customers Who Bought This Item Also Bought



Star Wars Trilogy Episodes I-III (Blu-ray + DVD)
Hayden Christiansen
★★★★☆ 2,042
Blu-ray
\$34.96



Star Wars: The Force Awakens (Blu-ray/DVD/Digital HD)
Harrison Ford
★★★★☆ 10,002
Blu-ray
\$24.41



Star Wars: Episode I - The Phantom Menace (Widescreen Edition)
Ewan McGregor
★★★★☆ 3,533
DVD
\$53.24

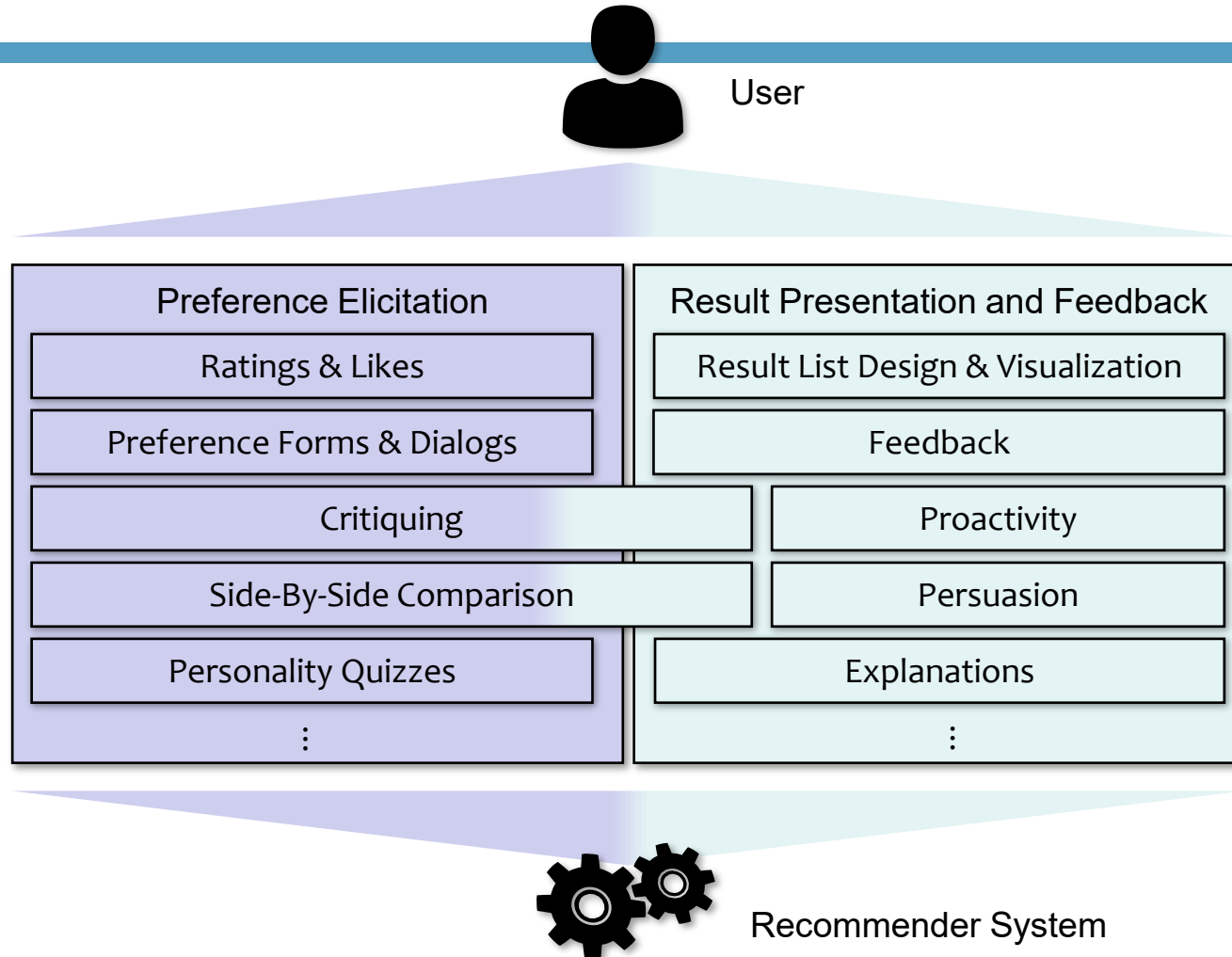


Harry Potter: Complete 8-Film Collection [Blu-ray]
Daniel Radcliffe
★★★★☆ 6,945
Blu-ray
\$65.00

UI research for Recommenders

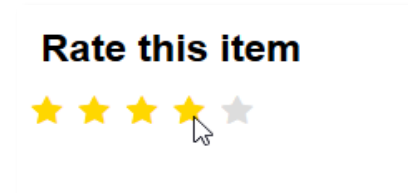
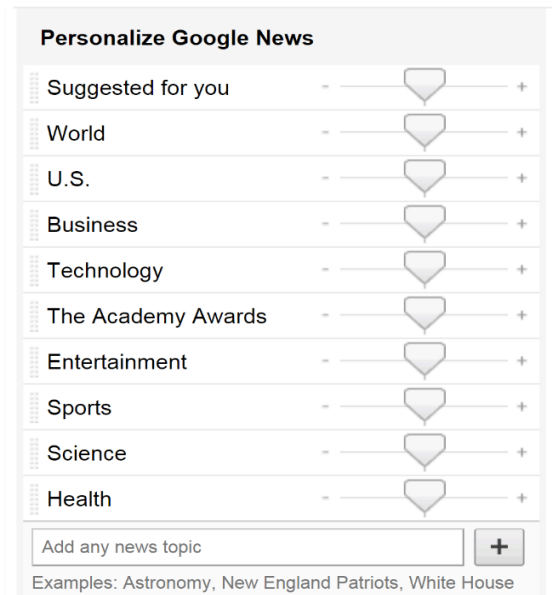
- HCI research is one of the main roots of recommender systems research
- Nonetheless, UI-related aspects seem less explored than algorithmic questions
 - One reason lies in the difficulty of evaluating new proposals
 - Existing research is also largely scattered

Structuring Existing Works

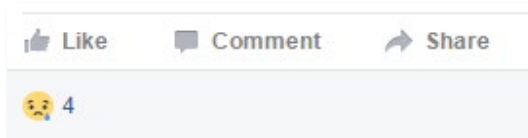


Design Space Examples

- Telling the system **explicitly** what you like
 - Global settings
 - Ratings
 - But how many options? How many categories?



Sources: Facebook.com,
Google.com



Design Space Examples

- What to display as recommendation?
 - The items of course
 - How many? Where on the screen? Multiple lists?
- Should users be able to give feedback?
 - Like/Dislike?
 - Or more?


Tell us why

I've already watched the video

I don't like the video

I'm not interested in this channel: Jimmy Kimmel Live

I'm not interested in recommendations based on:

 **Wild Animals with Dave Salmoni**
by Jimmy Kimmel Live

Cancel Submit

Source: Youtube.com

List Design Considerations

Customers Who Bought This Item Also Bought



The screenshot shows three product recommendations. The first is a Nikon lens, the second is a camera bag, and the third is a Lexar memory card. Each item has a title, a star rating, and a price. A red box highlights the 'Best Seller' badge for the Nikon lens. Annotations on the right point to the 'List label', 'Item description', 'Community rating', and 'Highlighting' (the red box). A bracket at the bottom indicates the 'Number of options'.

Item	Description	Rating	Price
Nikon AF-S FX NIKKOR 50mm f/1.8G Lens with Auto Focus for Nikon DSLR Cameras	Nikon AF-S FX NIKKOR 50mm f/1.8G Lens with Auto Focus for Nikon DSLR Cameras	★★★★★ 1,505	\$216.95 ✓Prime
Loweprro Adventura 140 Camera Shoulder Bag for DSLR or Camcorder	Loweprro Adventura 140 Camera Shoulder Bag for DSLR or Camcorder	★★★★☆ 178	\$26.99 ✓Prime
Lexar Professional 633x 64GB SDXC UHS-I Card w/Image Rescue 5 Software...	Lexar Professional 633x 64GB SDXC UHS-I Card w/Image Rescue 5 Software...	★★★★★ 592	\$22.50 ✓Prime

Source: amazon.com

Number of options

List label

Item description

Community rating

Highlighting

What else to show?

- What to display in addition to a nice picture?
 - Maybe some explanation, but which one?
 - A predicted rating?



The screenshot shows a movie recommendation card for "The Next Three Days". The card features a movie poster on the left with a play button icon. To the right, the title "The Next Three Days" is displayed in a red box, followed by the year "2010", a "PG-13" rating box, and "133 minutes". Below this, a white text box contains a synopsis: "When his wife is sent to jail on murder charges she fervently denies, a college professor hatches a meticulous plan for the ultimate prison escape." A "More Info" link is provided. The cast and director are listed: "Starring: Russell Crowe, Elizabeth Banks" and "Director: Paul Haggis". A recommendation based on user interests is shown: "Based on your interest in: Iron Man 2, John Q and X-Men Origins: Wolverine". At the bottom, there is a predicted rating section: "Our best guess for Xavier:" followed by five stars, four of which are filled. Two buttons are at the bottom: "Not interested" and "In Instant Queue".

The Next Three Days
2010 **PG-13** 133 minutes

When his wife is sent to jail on murder charges she fervently denies, a college professor hatches a meticulous plan for the ultimate prison escape.

[More Info](#)


Starring: Russell Crowe, Elizabeth Banks
Director: Paul Haggis

Based on your interest in: *Iron Man 2*, *John Q* and *X-Men Origins: Wolverine*


Our best guess for Xavier:
★★★★☆

Explanations and Control

- What to display in addition to a nice picture?
 - Maybe some explanation, but which one?
 - Or our logic to recommend this?



Recommended for you




[Guardians of the Galaxy \[Blu-ray\]](#)
Blu-ray ~ Chris Pratt (8 Jan 2015)
In stock
Price: EUR 9,99
[73 used & new](#) from EUR 8,75

Rate this item
 ☆☆☆☆☆

I own it
 Not interested

Because you purchased...



[Mad Max: Fury Road \[Blu-ray\]](#) (Blu-ray)
DVD ~ Charlize Theron

☆☆☆☆☆
 Don't use for recommendations

Interactive and Conversational Systems

- Increased research on Conversational Recommender Systems in recent years
 - Mostly chatbot-like systems, supporting natural language conversations
 - New horizons with ChatGPT

Jannach, D., Manzoor, A., Cai, W. and Chen, L.: "A Survey on Conversational Recommender Systems". ACM Computing Surveys, Vol. 54(5). 2021, pp. 1-26

Chongming Gao, Wenqiang Lei, Xiangnan He, Maarten de Rijke, Tat-Seng Chua, "Advances and challenges in conversational recommender systems: A survey", AI Open, Volume 2, 2021

Summary of second part

- We reviewed the history of technical approaches to build recommenders
- We found algorithmic works based on collaborative filtering to be dominant
 - Recently, sequence-aware recommenders were more in the focus
- In contrast, many questions regarding the design of a recommender system remain open
- The design space for the user interface, for example, is huge, but the literature is comparably scarce

Part III: Measurements

Evaluation aspects

- Computer Science research in this context is mostly about **building** “better” recommenders
 - i.e., systems or algorithms that serve a particular purpose better than alternative approaches
 - Often not about **understanding** what makes things better
- Typical purposes could be (see Part I)
 - Rank relevant items higher in the list
 - Make sure that the list is not monotonous
 - ...
 - Increase the user’s trust in the system
 - Provide a more convenient user interface

How can we know we are better?

- Testing a real application with real users
 - A/B tests (measuring, e.g., sales increase, CTR)
- Laboratory studies
 - Controlled experiments (measuring, e.g., satisfaction with the system), see later lecture
- Offline experiments
 - Simulations using on historical data (measuring, e.g., prediction accuracy, coverage)
- Theoretical analyses
 - For example, regarding scalability

Offline experiments

- Such experiments are, by far, the most common form of empirical research in the CS literature
- Main ingredients:
 - One or two historical dataset containing ratings or implicit feedback
 - A number of existing algorithms to compare the new proposal with
 - A number of established accuracy metrics (RMSE, Precision, Recall) and evaluation procedures to determine the metrics (e.g., cross-validation)

Sounds safe?

- All seems okay, “proving” progress in a reproducible way seems straightforward
 - At least one dataset should be public nowadays, so that others can replicate the results
 - The evaluation protocol and the metrics are well accepted and broadly known
 - The algorithmic proposals are usually laid out in great depth in the papers. Sometimes, even the source code is shared.

Progress can still be limited

- **Reason 1:** “Proving” progress by finding a better model for a very specific experimental setup can be relatively easy
- **Reason 2:** The used metrics are not necessarily helpful to measure improvements as perceived by users in the first place

Potential issues w/ research practice

- Applied ML research often obsessed with accuracy and the hunt for the “best model”
- But, there probably is no best model. The ranking of algorithms can depend on:
 - Given dataset
 - Used pre-processing steps
 - Evaluation measure
 - Choice of baselines
 - **Optimization of baselines**

Worrying observations

- Sometimes, it remains unclear if we truly make progress
 - Armstrong et al. (2009) find that there was not much progress within the previous ten years for a given Information Retrieval Task
 - Lin (2019) and Yang et al. (2019) found that ten years later problems with the choice of baselines still exist for deep learning methods
 - Rendle et al. (2019) run new experiments for classical recommendation tasks and find that recent methods are not necessarily better than previous ones

Worrying observations

- Makridakis (2018) compared various ML methods for time-series prediction, concluding that existing statistics-based methods are often better
- Ludewig et al. (2018-2020) evaluated various session-based recommendation techniques, finding that simple methods are often very competitive
- Ferrari Dacrema et al. (2019/2021) examined recent neural top-n recommendation techniques and found potential issues in terms of the choice and optimization of baselines
- Shehzad and Jannach (2024) found that none of a couple of GNN-based models outperformed basic models

Problems piling up, also in other ML areas

Armstrong, T.G., A. Moffat, W. Webber, and J. Zobel. 2009. “Improvements That Don’t Add Up: Ad-hoc Retrieval Results Since 1998.” In Proceedings of the 18th ACM Conference on Information and Knowledge Management (CIKM ’09), 601–10.

Lin, J. 2019. “The Neural Hype and Comparisons Against Weak Baselines.” ACM SIGIR Forum 52(2): 40–51.

Makridakis, S., E. Spiliotis, and V. Assimakopoulos. 2018. “Statistical and Machine Learning Forecasting Methods: Concerns and Ways Forward.” PloS one 13(3): 1–26.

Rendle, S., W. Krichene, L. Zhang, and J. Anderson. 2020. “Neural Collaborative Filtering vs. Matrix Factorization Revisited.” In Proceedings of the 14th ACM Conference on Recommender Systems (RecSys ’20).

Ferrari Dacrema, M., P. Cremonesi, and D. Jannach. 2019. “Are We Really Making Much Progress? A Worrying Analysis of Recent Neural Recommendation Approaches.” In Proceedings of the 13th ACM Conference on Recommender Systems (RecSys ’19), 101–9.

Ferrari Dacrema, M., S. Boglio, P. Cremonesi, and D. Jannach. 2021. “A troubling analysis of reproducibility and progress in recommender systems research.” ACM Transactions on Information Systems (TOIS), 39(2): 1–49.

Ludewig, M., S. Latifi, N. Mauro, and D. Jannach. “Empirical Analysis of Session-Based Recommendation Algorithms.” User Modeling and User-Adapted Interaction 31(3): 49–181.

Lipton, Z. C., and J. Steinhardt. 2019. “Research for Practice: Troubling Trends in Machine-Learning Scholarship.” Communications of the ACM 62(6): 45–53.

Beyond recommender systems



Thread



Tim Dettmers

@Tim_Dettmers



I got excited about a paper, implement stuff and then see they cheated: (1) Copy baseline results from other paper, (2) do much more hyperparam tuning on their own method, (3) accepted to EMNLP. Results look good, but their method is crap! Why waste people's time like this?

[Tweet übersetzen](#)

7:02 nachm. · 11. Feb. 2023 · **257.200** Mal angezeigt



Please Commit More Blatant Academic Fraud

Posted on May 29, 2021

Explicit academic fraud is, of course, the natural extension of the sort of mundane, day-to-day fraud that most academics in our community commit on a regular basis. **Trying that shiny new algorithm out on a couple dozen seeds, and then only reporting the best few. Running a big hyperparameter sweep on your proposed approach but using the defaults for the baseline. Cherry-picking examples where your model looks good, or cherry-picking whole datasets to test on, where you've confirmed your model's advantage.** Making up new problem settings, new datasets, new objectives in order to claim victory on an empty playing field. Proclaiming that your work is a “promising first step” in your introduction, despite being fully aware that nobody will ever build on it. Submitting a paper to a conference because it's got a decent shot at acceptance and you don't want the time you spent on it go to waste, even though you've since realized that the core ideas aren't quite correct.

Everyone's a Winner!

- Hyperparameters of baselines not well documented in the literature
- We compared eight different models
 - Seven neural ones, one popularity based
- We searched for good hyperparameters for all of them
- We compared them to the non-tuned versions of the others
- If we do not tune all models, we can declare everyone to be “Ours”, and outperform the state-of-the-art

Everyone's a Winner!

Tuned models					
ML-1M		AMZm		Epinions	
<i>Model</i>	<i>nDCG@10</i>	<i>Model</i>	<i>nDCG@10</i>	<i>Model</i>	<i>nDCG@10</i>
Mult-DAE	0,300	NeuMF	0,056	Mult-VAE	0,149
Mult-VAE	0,294	Mult-VAE	0,054	Mult-DAE	0,146
GMF	0,280	GMF	0,051	GMF	0,128
NeuMF	0,277	Mult-DAE	0,048	NeuMF	0,118
ONCF	0,225	<i>MostPop</i>	<i>0,013</i>	ONCF	0,077
<i>MostPop</i>	<i>0,162</i>	ConvMF	0,011	<i>MostPop</i>	<i>0,045</i>
ConvMF	0,160	NGCF	0,008	ConvMF	0,043
NGCF	0,100	ONCF	0,009	NGCF	0,031
Non-tuned models		>	>	>	>
Mult-DAE	0,071	Mult-DAE	0,003	Mult-DAE	0,015
ONCF	0,037	Mult-VAE	0,002	ONCF	0,005
ConvMF	0,022	ConvMF	0,002	NGCF	0,003
NeuMF	0,021	GMF	0,0007	GMF	0,002
GMF	0,016	NGCF	0,0006	Mult-VAE	0,002
NGCF	0,013	ONCF	0,0004	NeuMF	0,0008
Mult-VAE	0,006	NeuMF	0,0004	ConvMF	0,0008

Table 1. Accuracy results (NDCG@10) for tuned and non-tuned models, sorted by NDCG in descending order.

Potential ways forward

- Further increasing reproducibility is advocated
 - Reproducibility should be easy to establish
 - Many researchers use free software tools
 - Sharing images of the experimental environment is easy
 - Code should include everything from algorithm, over data-pre-processing and evaluation
- Choice and optimization of baselines as main problem
 - Often not clear what represents the state-of-the-art
 - Validation against optimized existing methods

Potential ways forward

- Toward more “theory-guided” research
 - Choice of dataset/pre-processing often seems arbitrary
 - Sometimes, researchers claim that their method is suited to make better recommendations
 - Then they use a rating dataset and transform all ratings to ones for evaluating an implicit feedback method
 - What is measured then, however, is how good we are at predicting who will rate what. Which does not necessarily mean better recommendations
 - Choice of evaluation procedures often seems arbitrary and not guided by an application problem
 - Various forms of measures used, cut-off lengths between one and several hundred, cross-validation/leave-one-out ...

Potential ways forward

- In the long run: Need better education and awareness
 - Among students, teachers, researchers, policy makers, ...

Bauer, C., Fröbe, M., Jannach, D., Kruschwitz, U., Rosso, P., Spina, D. and Tintarev, N.: "Overcoming Methodological Challenges in Information Retrieval and Recommender Systems through Awareness and Education". In: Dagstuhl Seminar 23031: Frontiers of Information Access Experimentation for Research and Education. 2023

Offline experiments and computational metrics in general

- Reason 2 from above: The used metrics are not necessarily helpful to measure improvements as perceived by users in the first place
- Generally:
 - Being able to accurately predict the relevance of items for users is and will be a central problem of recommender systems research
 - Increasing the prediction accuracy therefore can be a relevant goal of research

The problems with accuracy

- Accuracy alone is not enough
 - Recommending items that the user might have bought anyway might be of little business value
 - Focusing on accuracy alone can lead to monotone recommendations (e.g., only movies from the Star Wars series) and limited discovery
 - Optimizing for accuracy might lead to recommendations that are considered too “obscure” for users
 - Familiarity with some recommendations might be important to increase the user’s trust in a system

Multi-metric evaluations

- One possible way forward
- Offline experimentation can assess multiple, possibly competing, goals in parallel
 - Accuracy
 - Diversity
 - Novelty
 - Serendipity
 - Long-term effects, e.g., on reinforcement effects
 - Business value for multiple stakeholders
 - Scalability
 - ...

Multi-metric over-simplification

- Using some diversity metric along with accuracy may not be sufficient either
 - We need to validate that the metric matches user perceptions
- Moreover, the same set of recommendations can be good or not, depending on the purpose, context, and application, e.g.,
 - Recommending already popular items can be good for the business or not
 - Recommending things, for example musical songs, that the user already knows can be desirable or not, depending on the user's mood

General problems of offline experiments

- Are offline experiments actually predictive of the perceived value?
 - Gomez-Uribe and Hunt (2015), Netflix, found that offline experiments were **not** found “*to be as highly predictive of A/B test outcomes as we would like.*”
 - In fact, a number of user studies did **not** find that algorithms with higher prediction accuracy led to better quality perceptions by study participants

Accuracy, again

- In some domains, higher prediction accuracy almost directly leads to better systems
 - Language translation tasks
 - Image recognition tasks
- This analogy not necessarily holds for recommender systems
 - A small accuracy increase in a certain offline experiment might not tell us a lot about the quality of the resulting recommendations
- Problem: We measure (only) what we can easily measure
 - The McNamara Fallacy



Possible steps forward

- Toward a more comprehensive approach to recommender systems research
 - Considering the user in the loop
 - Considering the business value for one or more stakeholders
 - Use a richer methodological repertoire

- “From Algorithms to Systems”

User-centric research

- Much richer conceptual models of recommender systems and their impact exist in the field of Information Systems
 - Algorithms are only one of many components
 - Apparently limited knowledge of these works in the computer science community

User-centric research

- Different evaluation frameworks exist, e.g.,
 - Pu et al. (RecSys 2011, UMUAI 2012)
 - Knijnenburg et al. (UMUAI 2012)
- Frameworks describe relevant quality criteria
 - e.g., perceived accuracy, novelty, diversity, context compatibility, interface adequacy, information sufficiency and explainability, usefulness, ease of use
- and evaluation approaches
 - e.g., in terms of questionnaires

Wrap-up

Summary

- Discussed the business value of recommender systems
 - These are systems to have a real-world impact
 - It is great to work on this topic!
- Reviewed briefly how we can build such systems
 - From algorithms to the user experience
- Discussed limitations of current research practices
 - With an outlook on potential ways forward

-
- Thank you for your attention
 - Enjoy the summer school!

- Contact:

dietmar.jannach@aau.at

- Slides:

<https://tinyurl.com/eBISS2024>

