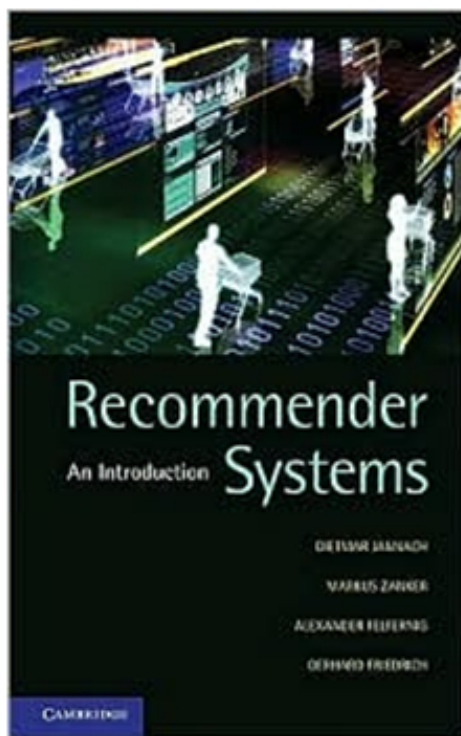

Tutorial: Evaluation of Recommender Systems

ACM Symposium on Applied Computing (SAC 2012)
Riva del Garda, 26 March 2012

Dietmar Jannach
TU Dortmund

Markus Zanker
Alpen-Adria-Universität Klagenfurt



Recommender Systems: An Introduction

by [Dietmar Jannach](#), [Markus Zanker](#), [Alexander Felfernig](#), [Gerhard Friedrich](#)

AVERAGE CUSTOMER RATING:

☆☆☆☆☆ ([Be the first to review](#))



Registrieren, um sehen zu können, was deinen Freunden gefällt.

FORMAT:
Hardcover

NOOKbook (eBook) - not available

[Tell the publisher you want this in NOOKbook format](#)

NEW FROM BN.COM

~~\$65.00~~ List Price

\$52.00 Online Price
(You Save 20%)

Add to Cart

NEW & USED FROM OUR

New starting at **\$56.46** (You Save 13%)
Used starting at **\$51.98** (You Save 20%)

See All Prices

[Table of Contents](#)

Customers who bought this also bought



Agenda

- **What are recommender systems for?**
 - Introduction
- **How do they work ?**
 - Collaborative Filtering
 - Content-based Filtering
 - Knowledge-Based Recommendations
 - Hybridization Strategies
- **How to measure their success?**
 - Evaluation techniques

Introduction



Problem domain

- **Recommendation systems (RS) help to match users with items**
 - Ease information overload
 - Sales assistance (guidance, advisory, persuasion,...)

RS are software agents that elicit the interests and preferences of individual consumers [...] and make recommendations accordingly.

They have the potential to support and improve the quality of the decisions consumers make while searching for and selecting products online.

» [Xiao & Benbasat, MISQ, 2007]

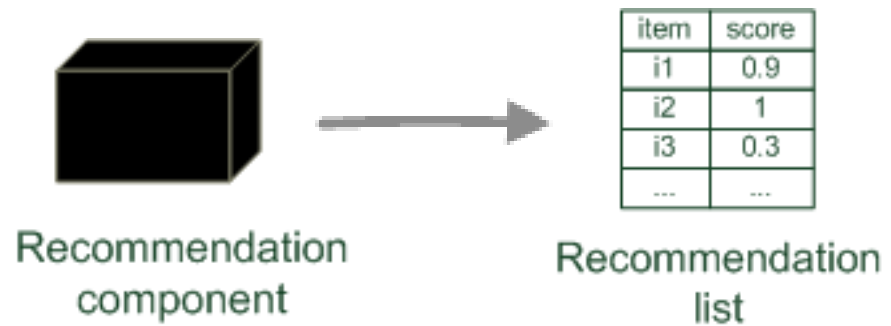


Recommender systems

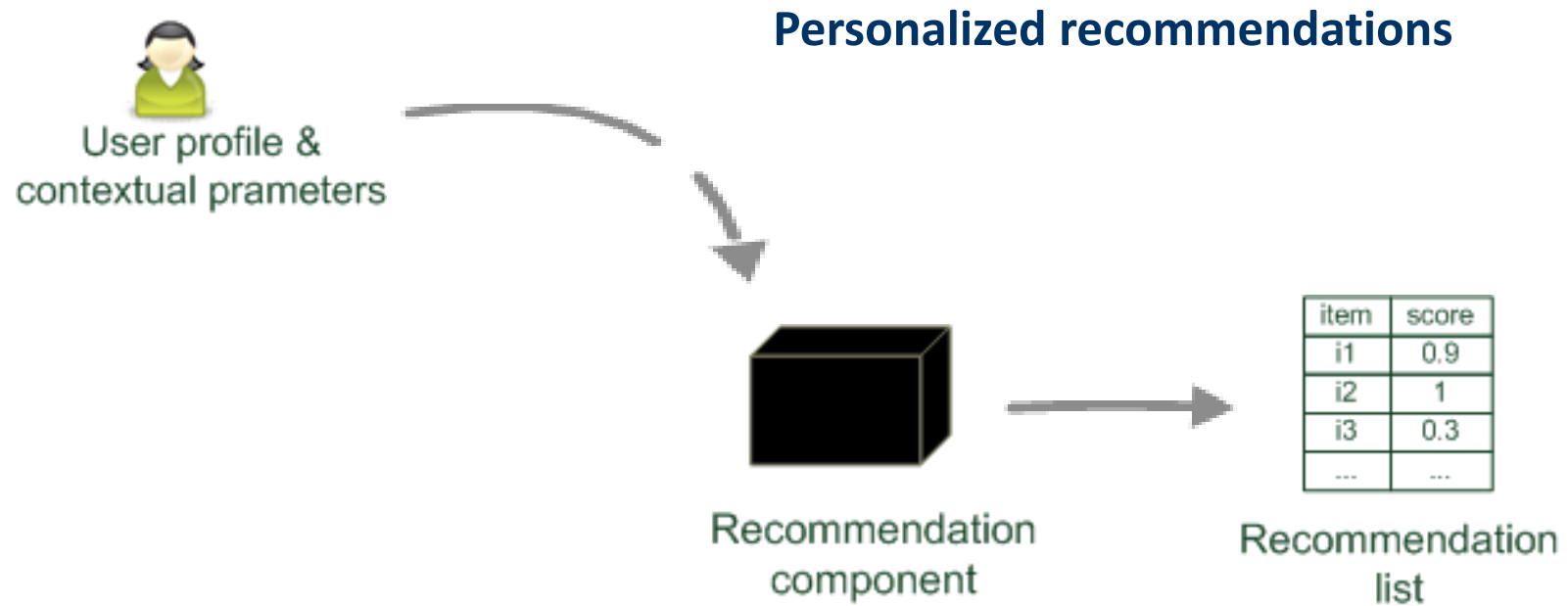
- **RS seen as a function** [AT05]
- **Given:**
 - User model (e.g. ratings, preferences, demographics, situational context)
 - Items (with or without description of item characteristics)
- **Find:**
 - Relevance score. Used for ranking.
- **At the end:**
 - Recommend items that are assumed to be relevant

Paradigms of recommender systems

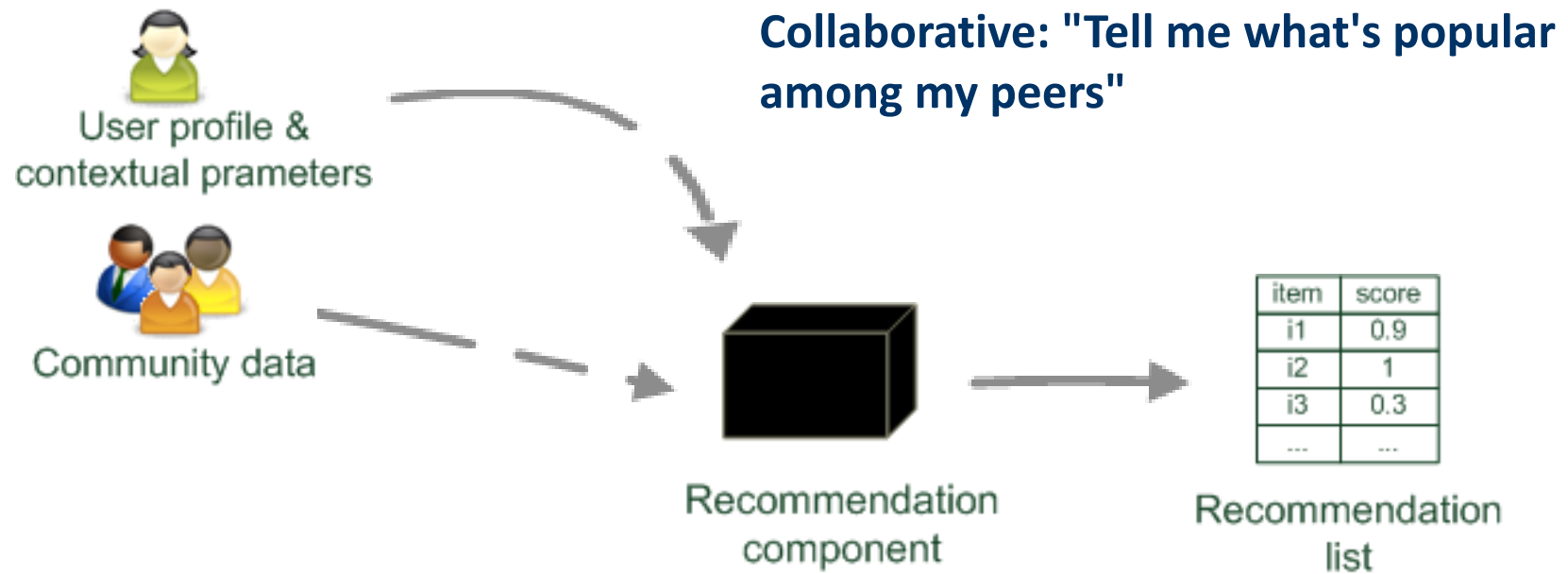
Recommender systems reduce information overload by estimating relevance



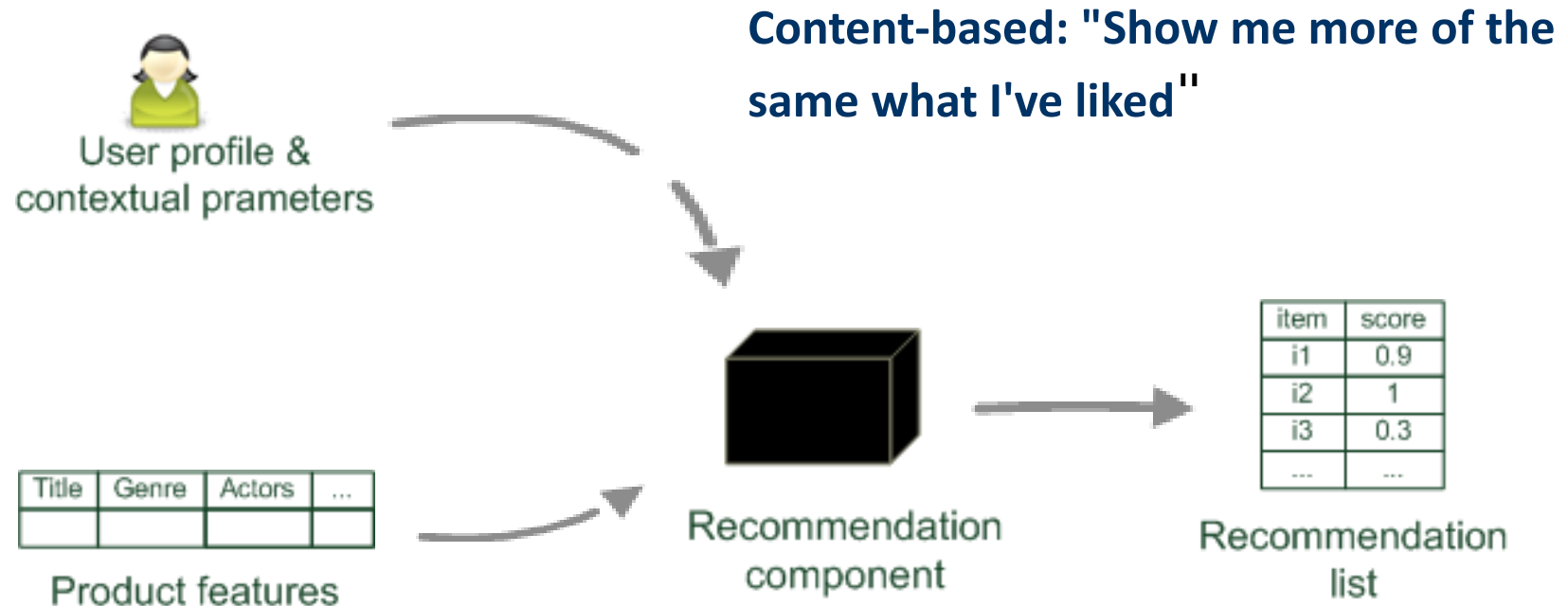
Paradigms of recommender systems



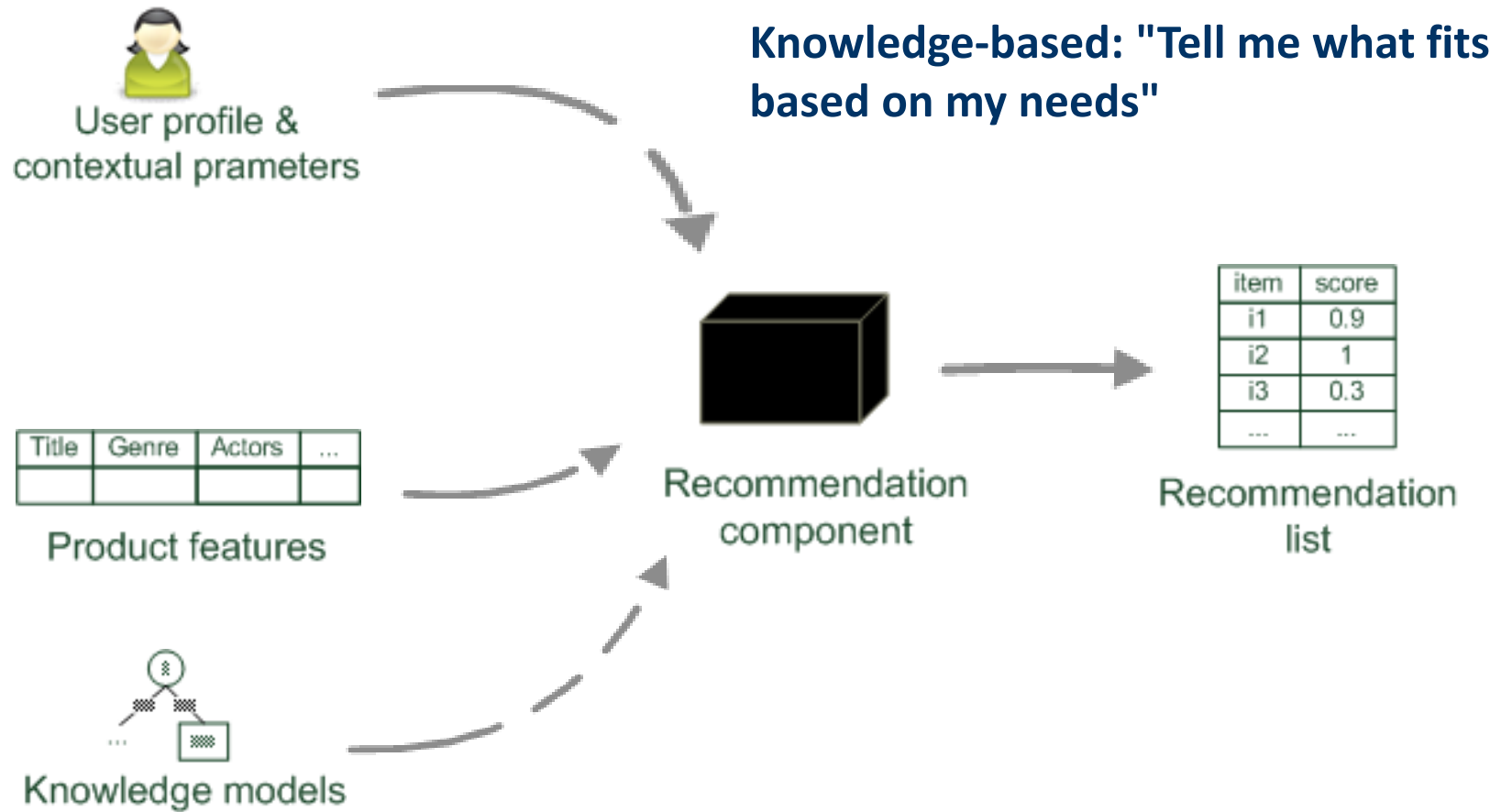
Paradigms of recommender systems



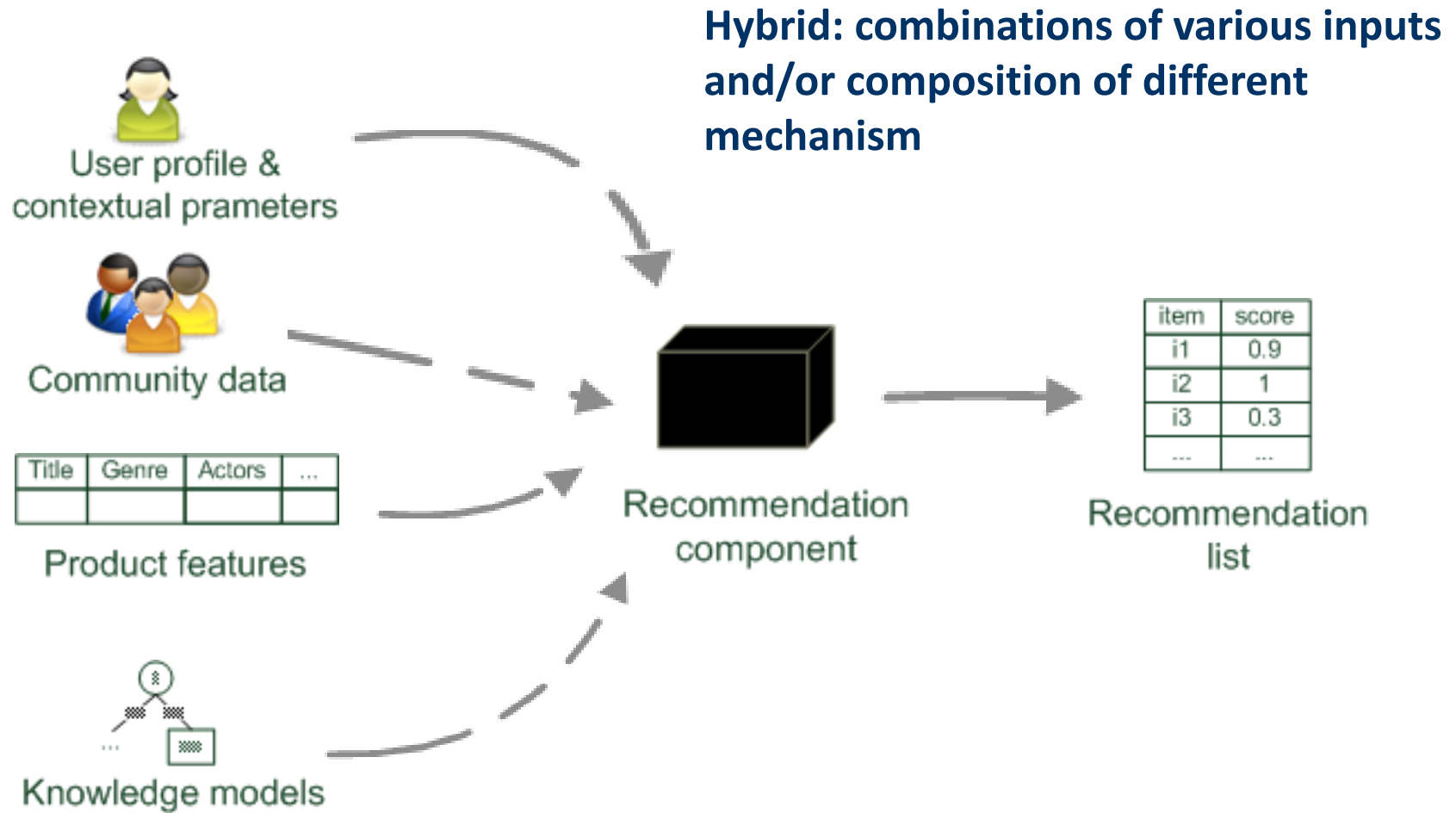
Paradigms of recommender systems



Paradigms of recommender systems



Paradigms of recommender systems



Collaborative Filtering

Collaborative Filtering (CF)

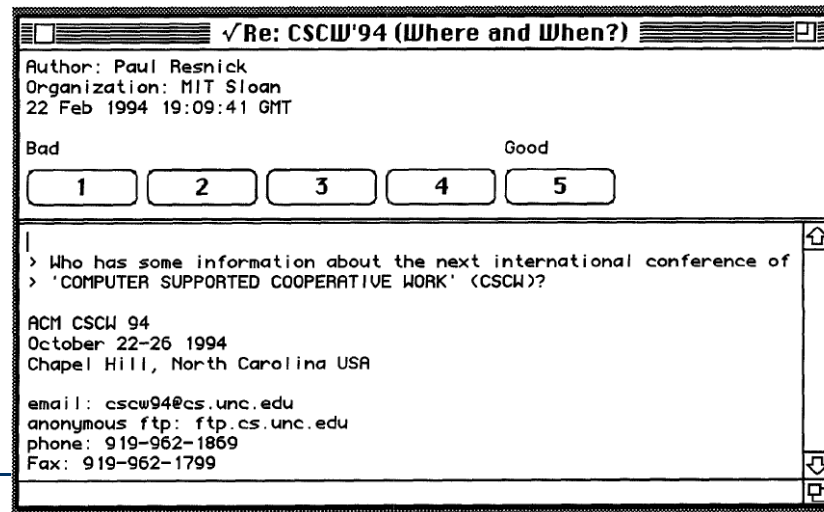
- **The most prominent approach to generate recommendations**
 - used by large, commercial e-commerce sites
 - well-understood, various algorithms and variations exist
 - applicable in many domains (book, movies, DVDs, ..)
- **Approach**
 - use the "wisdom of the crowd" to recommend items
- **Basic assumption and idea**
 - Users give ratings to catalog items (implicitly or explicitly)
 - Customers who had similar tastes in the past, will have similar tastes in the future

1992: *Using collaborative filtering to weave an information tapestry*, D. Goldberg et al., Communications of the ACM

- **Basic idea: "Eager readers read all docs immediately, casual readers wait for the eager readers to annotate"**
- **Experimental mail system at Xerox Parc that records reactions of users when reading a mail**
- **Users are provided with personalized mailing list filters instead of being forced to subscribe**
 - Content-based filters (topics, from/to/subject...)
 - Collaborative filters
- **E.g. Mails to [all] which were replied by [John Doe] and which received positive ratings from [X] and [Y].**

1994: *GroupLens: an open architecture for collaborative filtering of netnews*, P. Resnick et al., ACM CSCW

- Tapestry system does not aggregate ratings and requires knowing each other
- Basic idea: "People who agreed in their subjective evaluations in the past are likely to agree again in the future"
- Builds on newsgroup browsers with rating functionality



User-based nearest-neighbor collaborative filtering (1)

- **The basic technique:**

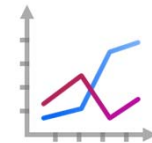
- Given an "active user" (Alice) and an item I not yet seen by Alice
- The *goal is to estimate Alice's rating for this item, e.g., by*
 - find a set of users (peers) who liked the same items as Alice in the past **and** who have rated item I
 - use, e.g. the average of their ratings to predict, if Alice will like item I
 - do this for all items Alice has not seen and recommend the best-rated

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

User-based nearest-neighbor collaborative filtering (2)

- **Some first questions**

- How do we measure similarity?
- How many neighbors should we consider?
- How do we generate a prediction from the neighbors' ratings?



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

Measuring user similarity

- A popular similarity measure in user-based CF: Pearson correlation

$$sim(a, b) = \frac{\sum_{p \in P} (r_{a,p} - \bar{r}_a)(r_{b,p} - \bar{r}_b)}{\sqrt{\sum_{p \in P} (r_{a,p} - \bar{r}_a)^2} \sqrt{\sum_{p \in P} (r_{b,p} - \bar{r}_b)^2}}$$

a, b : users

$r_{a,p}$: rating of user a for item p

P : set of items, rated both by a and b

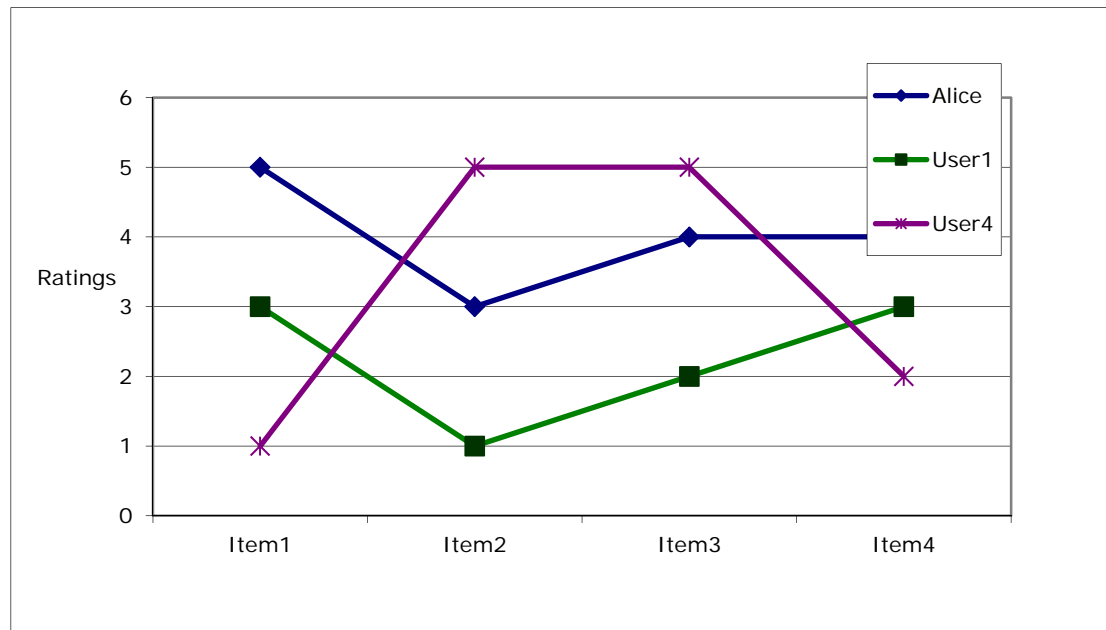
Possible similarity values between -1 and 1; \bar{r}_a, \bar{r}_b = user's average ratings

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

sim = 0,85
 sim = 0,70
 sim = -0,79

Pearson correlation

- Takes differences in rating behavior into account



- Works well in usual domains, compared with alternative measures
 - such as cosine similarity

Making predictions

- A common prediction function:

$$pred(a, p) = \bar{r}_a + \frac{\sum_{b \in N} sim(a, b) * (r_{b,p} - \bar{r}_b)}{\sum_{b \in N} sim(a, b)}$$



- Calculate, whether the neighbors' ratings for the unseen item *i* are higher or lower than their average
- Combine the rating differences – use the similarity with as a weight
- Add/subtract the neighbors' bias from the active user's average and use this as a prediction

Improving the metrics / prediction function

- **Not all neighbor ratings might be equally "valuable"**
 - Agreement on commonly liked items is not so informative as agreement on controversial items
 - **Possible solution:** Give more weight to items that have a higher variance
- **Value of number of co-rated items**
 - Use "significance weighting", by e.g., linearly reducing the weight when the number of co-rated items is low
- **Case amplification**
 - Intuition: Give more weight to "very similar" neighbors, i.e., where the similarity value is close to 1.
- **Neighborhood selection**
 - Use similarity threshold or fixed number of neighbors

2001: *Item-based collaborative filtering recommendation algorithms*, B. Sarwar et al., WWW 2001

- **Basic idea:**
 - Use the similarity between items (and not users) to make predictions
- **Example:**
 - Look for items that are similar to Item5
 - Take Alice's ratings for these items to predict the rating for Item5

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

More on ratings

- **Pure CF-based systems only rely on the rating matrix**
- **Explicit ratings**
 - Most commonly used (1 to 5, 1 to 7 Likert response scales)
 - Challenge
 - Users not always willing to rate many items; sparse rating matrices
 - How to stimulate users to rate more items?
- **Implicit ratings**
 - clicks, page views, time spent on some page, demo downloads ...
 - Can be used in addition to explicit ones; question of correctness of interpretation

Data sparsity problems

- **Cold start problem**
 - How to recommend new items? What to recommend to new users?
- **Straightforward approaches**
 - Ask/force users to rate a set of items
 - Use another method (e.g., content-based, demographic or simply non-personalized) in the initial phase
- **Alternatives**
 - Use better algorithms (beyond nearest-neighbor approaches)
 - Example:
 - In nearest-neighbor approaches, the set of sufficiently similar neighbors might be too small to make good predictions
 - Assume "transitivity" of neighborhoods

Memory-based and model-based approaches

- **User-based CF is said to be "memory-based"**
 - the rating matrix is directly used to find neighbors / make predictions
 - does not scale for most real-world scenarios
 - large e-commerce sites have tens of millions of customers and millions of items

- **Model-based approaches**
 - based on an offline pre-processing or "model-learning" phase
 - at run-time, only the learned model is used to make predictions
 - models are updated / re-trained periodically
 - large variety of techniques used
 - model-building and updating can be computationally expensive

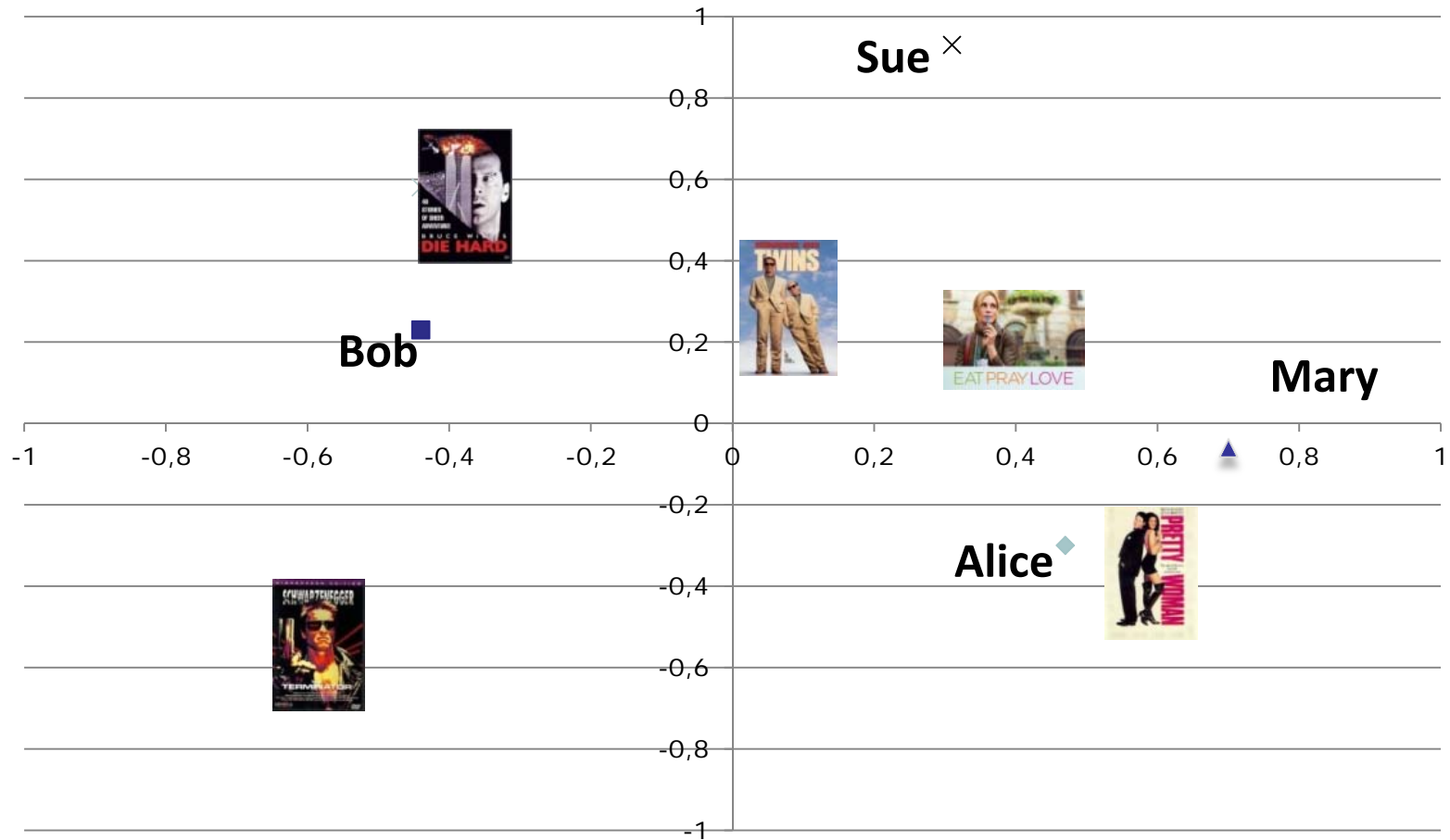
Model-based approaches

- **Plethora of different techniques proposed in the last years, e.g.,**
 - Matrix factorization techniques, statistics
 - singular value decomposition, principal component analysis
 - Association rule mining
 - compare: shopping basket analysis
 - Probabilistic models
 - clustering models, Bayesian networks, probabilistic Latent Semantic Analysis
 - Various other machine learning approaches
- **Costs of pre-processing**
 - Usually not discussed
 - Incremental updates possible?

2000: *Application of Dimensionality Reduction in Recommender System*, B. Sarwar et al., WebKDD Workshop

- **Basic idea: Trade more complex offline model building for faster online prediction generation**
- **Singular Value Decomposition for dimensionality reduction of rating matrices**
 - Captures important factors/aspects and their weights in the data
 - factors can be genre, actors but also non-understandable ones
 - Assumption that k dimensions capture the signals and filter out noise ($K = 20$ to 100)
- **Constant time to make recommendations**
- **Approach also popular in IR (Latent Semantic Indexing), data compression,...**

A picture says ...





2008: *Factorization meets the neighborhood: a multifaceted collaborative filtering model*, Y. Koren, ACM SIGKDD

- **Merges neighborhood models with latent factor models**
- **Latent factor models**
 - good to capture weak signals in the overall data
- **Neighborhood models**
 - good at detecting strong relationships between close items
- **Combination in one prediction single function**
 - Local search method such as stochastic gradient descent to determine parameters
 - Add penalty for high values to avoid over-fitting

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

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

Collaborative Filtering Issues

- **Pros:** 
 - well-understood, works well in some domains, no knowledge engineering required
- **Cons:** 
 - requires user community, sparsity problems, no integration of other knowledge sources, no explanation of results

Content-based recommendation

Content-based recommendation

- **While CF – methods do not require any information about the items,**
 - it might be reasonable to exploit such information; and
 - recommend fantasy novels to people who liked fantasy novels in the past
- **What do we need:**
 - some information about the available items such as the genre ("content")
 - some sort of *user profile* describing what the user likes (the preferences)
- **The task:**
 - learn user preferences
 - locate/recommend items that are "similar" to the user preferences

What is the "content"?

- **The genre is actually not part of the content of a book**
- **Most CB-recommendation methods originate from Information Retrieval (IR) field:**
 - goal is to find and rank interesting text documents (news articles, web pages)
 - the item descriptions are usually automatically extracted (important words)
- **Fuzzy border between content-based and "knowledge-based" RS**
- **Here:**
 - classical IR-based methods based on keywords
 - no expert recommendation knowledge involved
 - User profile (preferences) are rather learned than explicitly elicited

Content representation and item similarities

Title	Genre	Author	Type	Price	Keywords
The Night of the Gun	Memoir	David Carr	Paperback	29.90	Press and journalism, drug addiction, personal memoirs, New York
The Lace Reader	Fiction, Mystery	Brunonia Barry	Hardcover	49.90	American contemporary fiction, detective, historical
Into the Fire	Romance, Suspense	Suzanne Brockmann	Hardcover	45.90	American fiction, Murder, Neo-nazism
...					

Title	Genre	Author	Type	Price	Keywords
...	Fiction, Suspense	Brunonia Barry, Ken Follet, ..	Paperback	25.65	detective, murder, New York

- **Simple approach**

- Compute the similarity of an unseen item with the user profile based on the keyword overlap (e.g. using the Dice coefficient)
- $$\text{sim}(b_i, b_j) = \frac{2 * |\text{keywords}(b_i) \cap \text{keywords}(b_j)|}{|\text{keywords}(b_i)| + |\text{keywords}(b_j)|}$$
- Or combine multiple metrics in a weighted approach

Term-Frequency - Inverse Document Frequency (TF-IDF)

- **Simple keyword representation has its problems**
 - in particular when automatically extracted as
 - not every word has similar importance
 - longer documents have a higher chance to have an overlap with the user profile

- **Standard measure: TF-IDF**
 - Encodes text documents in multi-dimensional Euclidian space
 - weighted term vector
 - TF: Measures, how often a term appears (density in a document)
 - assuming that important terms appear more often
 - normalization has to be done in order to take document length into account
 - IDF: Aims to reduce the weight of terms that appear in all documents

TF-IDF

- **Compute the overall importance of keywords**

- Given a keyword i and a document j

$$TF-IDF(i,j) = TF(i,j) * IDF(i)$$

- **Term frequency (TF)**

- Let $freq(i,j)$ number of occurrences of keyword i in document j
- Let $maxOthers(i,j)$ denote the highest number of occurrences of another keyword of j

- $TF(i,j) = \frac{freq(i,j)}{maxOthers(i,j)}$

- **Inverse Document Frequency (IDF)**

- N : number of all recommendable documents
- $n(i)$: number of documents in which keyword i appears

- $IDF(i) = \log \frac{N}{n(i)}$

Example TF-IDF representation

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	5.25	3.18	0	0	0	0.35
Brutus	1.21	6.1	0	1	0	0
Caesar	8.59	2.54	0	1.51	0.25	0
Calpurnia	0	1.54	0	0	0	0
Cleopatra	2.85	0	0	0	0	0
mercy	1.51	0	1.9	0.12	5.25	0.88
worser	1.37	0	0.11	4.15	0.25	1.95

Figure taken from <http://informationretrieval.org>

More on the vector space model

- **Vectors are usually long and sparse**
 - **Improvements**
 - remove stop words ("a", "the", ..)
 - use stemming
 - size cut-offs (only use top n most representative words, e.g. around 100)
 - use additional knowledge, use more elaborate methods for feature selection
 - detection of phrases as terms (such as United Nations)
 - **Limitations**
 - semantic meaning remains unknown
 - example: usage of a word in a negative context
 - "there is nothing on the menu that a vegetarian would like.."
 - **Usual similarity metric to compare vectors: Cosine similarity (angle)**
-

Recommending items

- **Simple method: nearest neighbors**
 - Given a set of documents D already rated by the user (like/dislike)
 - Find the n nearest neighbors of a not-yet-seen item i in D
 - Take these ratings to predict a rating/vote for i
 - (Variations: neighborhood size, lower/upper similarity thresholds..)
 - Good to model short-term interests / follow-up stories
 - Used in combination with method to model long-term preferences

- **Other methods**
 - Rocchio's feedback
 - Probabilistic methods

Probabilistic methods

- **Recommendation as classical text classification problem**
 - long history of using probabilistic methods
- **Simple approach:**
 - 2 classes: hot/cold
 - simple Boolean document representation
 - calculate probability that document is hot/cold based on Bayes theorem

Doc-ID	recommender	intelligent	learning	school	Label
1	1	1	1	0	1
2	0	0	1	1	0
3	1	1	0	0	1
4	1	0	1	1	1
5	0	0	0	1	0
6	1	1	0	0	?

$$\begin{aligned}
 P(X|Label=1) &= P(\text{recommender}=1|Label=1) \times \\
 &\quad P(\text{intelligent}=1|Label=1) \times \\
 &\quad P(\text{learning}=0|Label=1) \times P(\text{school}=0|Label=1) \\
 &= \frac{3}{3} \times \frac{2}{3} \times \frac{1}{3} \times \frac{2}{3} \\
 &\approx 0.149
 \end{aligned}$$

Improvements

- **Side note: Conditional independence of events does in fact not hold**
 - "New York", "Hong Kong"
 - Still, good accuracy can be achieved
 - **Boolean representation simplistic**
 - positional independence assumed
 - keyword counts lost
 - **More elaborate probabilistic methods**
 - e.g., estimate probability of term v occurring in a document of class C by relative frequency of v in all documents of the class
 - **Other linear classification algorithms (machine learning) can be used**
 - Support Vector Machines, ..
 - **Use other information retrieval methods (used by search engines..)**
-

Limitations of content-based recommendation methods

- **Keywords alone may not be sufficient to judge quality/relevance of a document or web page**
 - up-to-dateness, usability, aesthetics, writing style
 - content may also be limited / too short
 - content may not be automatically extractable (multimedia)
- **Ramp-up phase required**
 - Some training data is still required
 - Web 2.0: Use other sources to learn the user preferences
- **Overspecialization**
 - Algorithms tend to propose "more of the same"
 - Or: too similar news items

Knowledge-Based Recommender Systems



Knowledge-Based Recommendation

- **Explicit domain knowledge**
 - Sales knowledge elicitation from domain experts
 - System mimics the behavior of experienced sales assistant
 - Best-practice sales interactions
 - Can guarantee “correct” recommendations (determinism) with respect to expert knowledge

- **Conversational interaction strategy**
 - Opposed to one-shot interaction
 - Elicitation of user requirements
 - Transfer of product knowledge (“educating users”)

Knowledge-Based Recommendation

- **Different views on “knowledge”**
 - Similarity functions
 - Determine matching degree between query and item (case-based RS)
 - Utility-based RS
 - E.g. MAUT – Multi-attribute utility theory
 - Logic-based knowledge descriptions (from domain expert)
 - E.g. Hard and soft constraints

- **Hybridization**
 - E.g. merging explicit knowledge with community data
 - Can ensure some policies based on e.g. availability, user context or profit margin

Constraint-based recommendation (Filtering)

Knowledge Base:

	LHS	RHS
C1	TRUE	Brand = Brand pref.
C2	Motives = <i>Landscape</i>	Low. foc. Length ≤ 28
C3	TRUE	Price \leq Max. cost

Current user:

	User model (SRS)	
R1	Motives	<i>Landscape</i>
R2	Brand preference	<i>Canon</i>
R3	Max. cost	<i>350 EUR</i>

Product catalogue:

Powershot XY	
Brand	<i>Canon</i>
Lower focal length	<i>35</i>
Upper focal length	<i>140</i>
Price	<i>420 EUR</i>
Lumix	
Brand	<i>Panasonic</i>
Lower focal length	<i>28</i>
Upper focal length	<i>112</i>
Price	<i>319 EUR</i>

Constraint-based recommendation

- **A knowledge-based RS formulated as constraint satisfaction problem**

$$CSP (X_I \cup X_U, D, SRS \cup KB \cup I)$$

- **Def.**

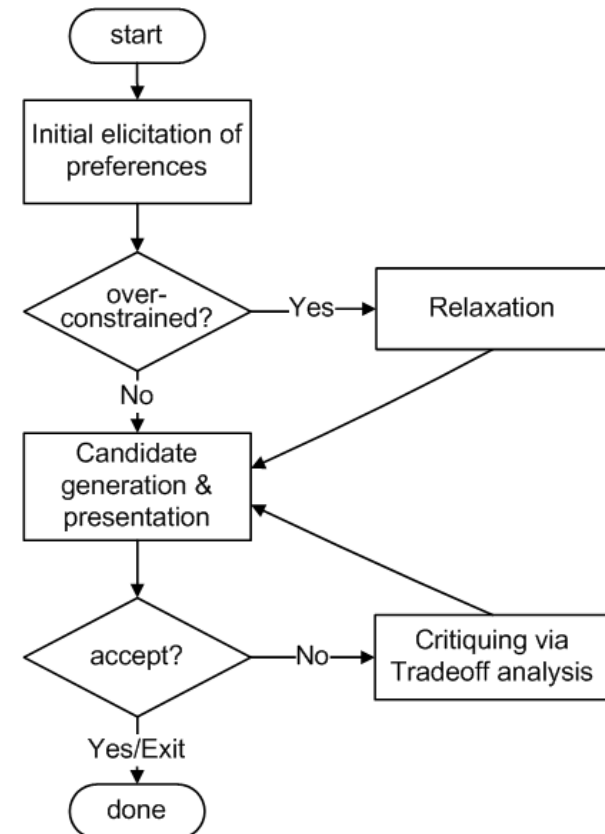
- X_I, X_U : Variables describing items and user model with domain D (e.g. lower focal length, purpose)
- KB : Knowledge base comprising constraints and domain restrictions (e.g. **IF** purpose="on travel" **THEN** lower focal length < 28mm)
- SRS : Specific requirements of a user (e.g. purpose = "on travel")
- I : Product catalog (e.g. $(id=1 \wedge lfl = 28mm) \vee (id=2 \wedge lfl = 35mm) \vee \dots$)

- **Solution: Assignment tuple θ assigning values to all variables X_i**

s.t. $SRS \cup KB \cup I \cup \theta$ **is satisfiable.**

Conversational strategies

- **Process consisting of multiple conversational moves**
 - Resembles natural sales interactions
 - Not all user requirements known beforehand
 - Customers are rarely satisfied with the initial recommendations
- **Different styles of preference elicitation:**
 - Free text query interface
 - Asking technical/generic properties
 - Images / inspiration
 - Proposing and Critiquing



Example: critiquing

*Find your
Favourite restaurant* 

In Vienna you chose:

+43 1 123 123 123 **Biergasthof** 30€-50€
Mariahilferstrasse 123,
1010 Wien Local cuisine

local food, central in the city, weekend brunch, room with a view,
famous for beer, seasonal dishes, group bookings, open all day

For Graz we recommend:

+43 316 45 45 45 **Brauhof** 30€-50€
Brauhoferstrasse 45,
8023 Graz Local cuisine

local food, own beer, weekend lunch, open all day, private function room,
famous for beer, seasonal dishes, group bookings, good transport connection

Less \$\$ *Nicer* *Cuisine* *More Quiet*

Traditional *Creative* *Livelier*

- Similarity-based navigation in item space
- Compound critiques
 - More efficient navigation than with unit critiques

Limitations of knowledge-based recommendation methods

- **Cost of knowledge acquisition**

- From domain experts
- From users
- From web resources

- **Accuracy of preference models**

- Very fine granular preference models require many interaction cycles with the user or sufficient detailed data about the user
- Preferences may depend on each other
- Collaborative filtering models the preference of a user implicitly

- **Instability of preference models**

- E.g. asymmetric dominance effects and decoy items

Hybridization Strategies

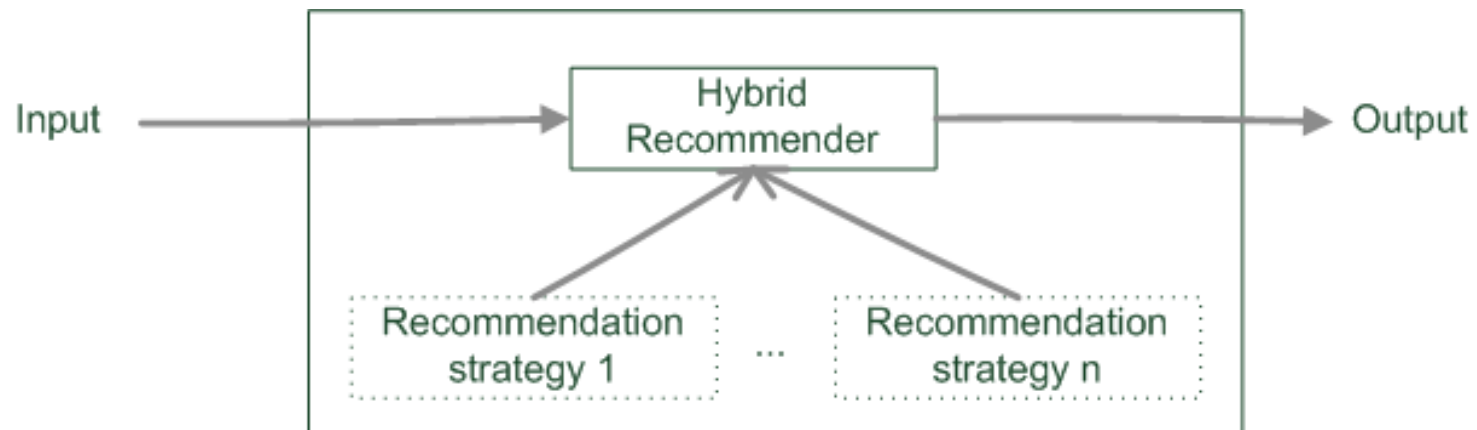


Hybrid recommender systems

- **All three base techniques are naturally incorporated by a good sales assistance (at different stages of the sales act) but have their shortcomings**
- **Idea of crossing two (or more) species/implementations**
 - Avoid some of the shortcomings
 - Reach desirable properties not (or only inconsistently) present in parent individuals

Monolithic hybridization design

- Only a single recommendation component



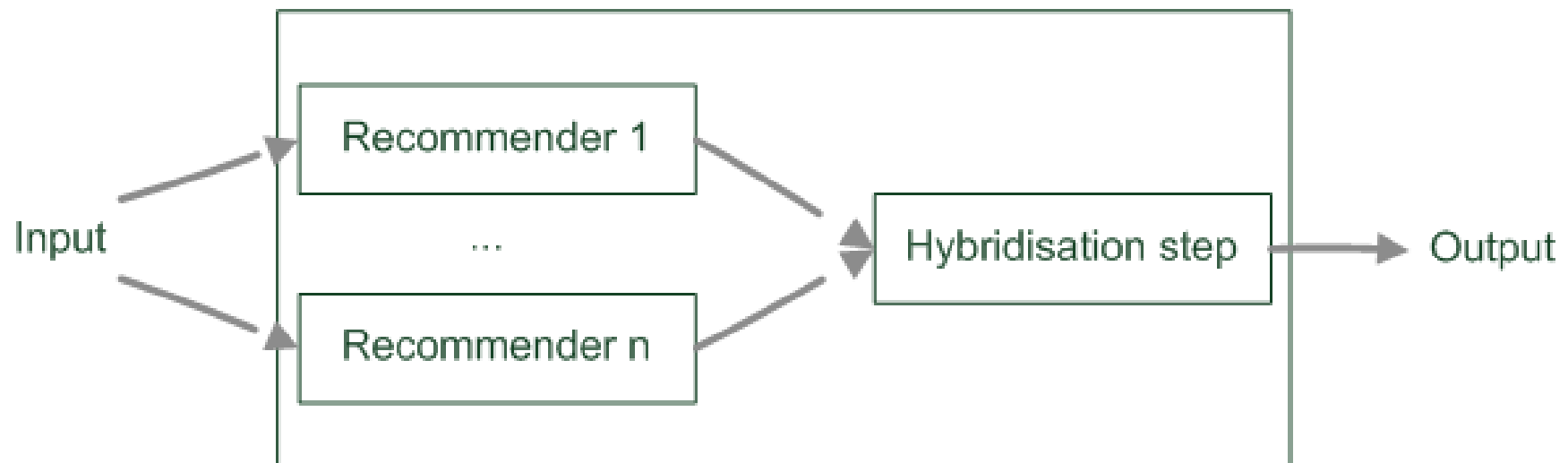
- Hybridization is "virtual" in the sense that
 - Features/knowledge sources of different paradigms are combined

Monolithic hybridization designs: Feature combination

- **"Hybrid" user features:**
 - Social features: Movies liked by user
 - Content features: Comedies liked by user, dramas liked by user
 - Hybrid features: users who like many movies that are comedies, ...

Parallelized hybridization design

- **Output of several existing implementations combined**
- **Least invasive design**
- **Weighting or voting scheme applied**
 - Weights can be learned dynamically



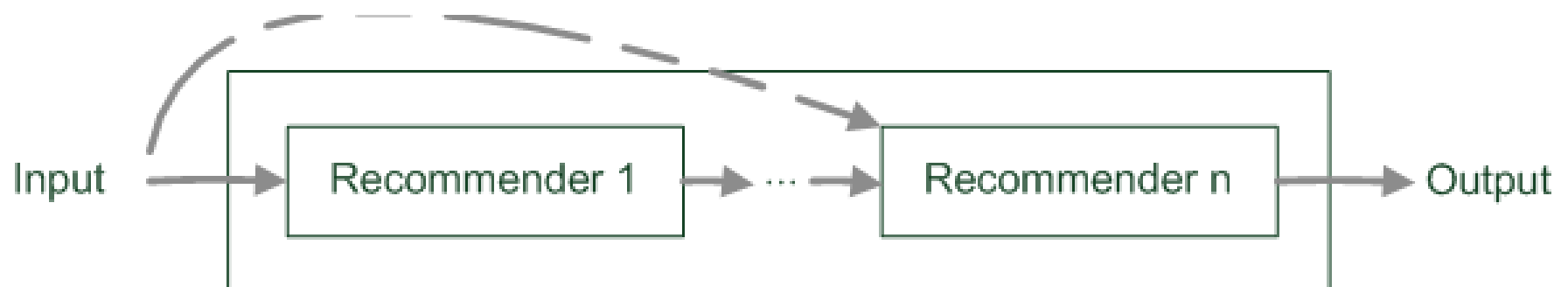
Parallelized hybridization design: Switching

- **Special case of dynamic weights (all weights except one are 0)**
- **Requires an oracle that decides which recommender is used**
- **Example:**
 - Ordering on recommenders and switch based on some quality criteria:
E.g. if too few ratings in the system, use knowledge-based, else collaborative
 - More complex conditions based on contextual parameters, apply classification techniques



Pipelined hybridization designs

- **One recommender system pre-processes some input for the subsequent one**
 - Cascade
 - Meta-level
- **Refinement of recommendation lists (cascade)**
- **Learning of model (e.g. collaborative knowledge-based meta-level)**



Pipelined hybridization designs: Cascade

<i>Recommender 1</i>		
Item1	0.5	1
Item2	0	
Item3	0.3	2
Item4	0.1	3
Item5	0	

<i>Recommender 2</i>		
Item1	0.8	2
Item2	0.9	1
Item3	0.4	3
Item4	0	
Item5	0	

<i>Recommender cascaded (rec1, rec2)</i>		
Item1	0,80	1
Item2	0,00	
Item3	0,40	2
Item4	0,00	
Item5	0,00	

- **Recommendation list is continually reduced**
- **First recommender excludes items**
 - Remove absolute no-go items (e.g. knowledge-based)
- **Second recommender assigns score**
 - Ordering and refinement (e.g. collaborative)

Limitations and success of hybridization strategies

- **Only few works that compare strategies from the meta-perspective**
 - For instance, [Burke02]
 - Most datasets do not allow to compare different recommendation paradigms
 - i.e. ratings, requirements, item features, domain knowledge, critiques rarely available in a single dataset
 - Thus few conclusions that are supported by empirical findings
 - Monolithic: some preprocessing effort traded-in for more knowledge included
 - Parallel: requires careful matching of scores from different predictors
 - Pipelined: works well for two antithetic approaches
- **Netflix competition – “stacking” recommender systems**
 - Weighted design based on >100 predictors – recommendation functions
 - Adaptive switching of weights based on user model, parameters (e.g. number of ratings in one session)

Evaluating Recommender Systems

Agenda

- **What is the current state-of-practice?**

- „How to“ from different perspectives:
 - Empirical research principles
 - Information Retrieval
 - Machine Learning
 - HCI and Decision Support

- Outlook

What is popular?

- **Small quantitative survey in the literature (Jannach et al., 2010)**
 - Evaluation designs ACM TOIS 2004-2010
 - In total 15 articles on RS
 - Nearly 50% movie domain
 - 80% offline experimentation
 - 2 user experiments under lab conditions
 - 1 qualitative research

- **Wide-scale survey (Jannach et al., 2012)**
 - 330 publications from a predefined set of conferences and journals
 - Systematic review of the period Jan. 2006 - Jul. 2011
 - 73 journal publications
 - 20% of total from IS community

Publication outlets 1/2

Conference	Field	#Pub.
ACM Conf. on Human Factors in Comp. Syst. (CHI)	CS	13
ACM Conf. on Recommender Syst. (RecSys)	CS	86
Int. Conf. on Int. User Interfaces (IUI)	CS	17
Int. Conf. on Knowl. Disc. and DM (SIGKDD)	CS	22
Int. Conf. on Res. and Dev. in IR (SIGIR)	CS	33
Int. Conf. on World Wide Web (WWW)	CS	21
Int. Joint Conf. on AI (IJCAI)	CS	13
AAAI Conf. on AI (AAAI)	CS	10
Int. Conf. on Data Mining (ICDM)	CS	5
Americas Conf. on Information Systems (AMCIS)	IS	8
European Conf. on Information Systems (ECIS)	IS	6
Int. Conf. on Information Systems (ICIS)	IS	7
Med. Conf. on Information Systems (MCIS)	IS	5
Pac. Asia Conf. on Information Systems (PACIS)	IS	11

Publication outlets 2/2

Journal	Field	#Pub.
ACM Trans. on Intell. Syst. and Techn. (TOIST)	CS	6
ACM Trans. on the Web (TWeb)	CS	5
AI Communications	CS	12
IEEE Intelligent Systems	CS	14
Int. Jرنل. of Human Computer Studies (IJHCS)	CS	5
World Wide Web (WWW)	CS	3
Dec. Supp. Syst. Jرنل. (DSS)	IS	9
Inf. Syst. Res. (ISR)	IS	3
Int. Jرنل. of Electronic Comm. (IJEC)	IS	7
Jرنل. of Mgt. Information Systems (JMIS)	IS	7
Mgt. Information Systems Quarterly (MISQ)	IS	2

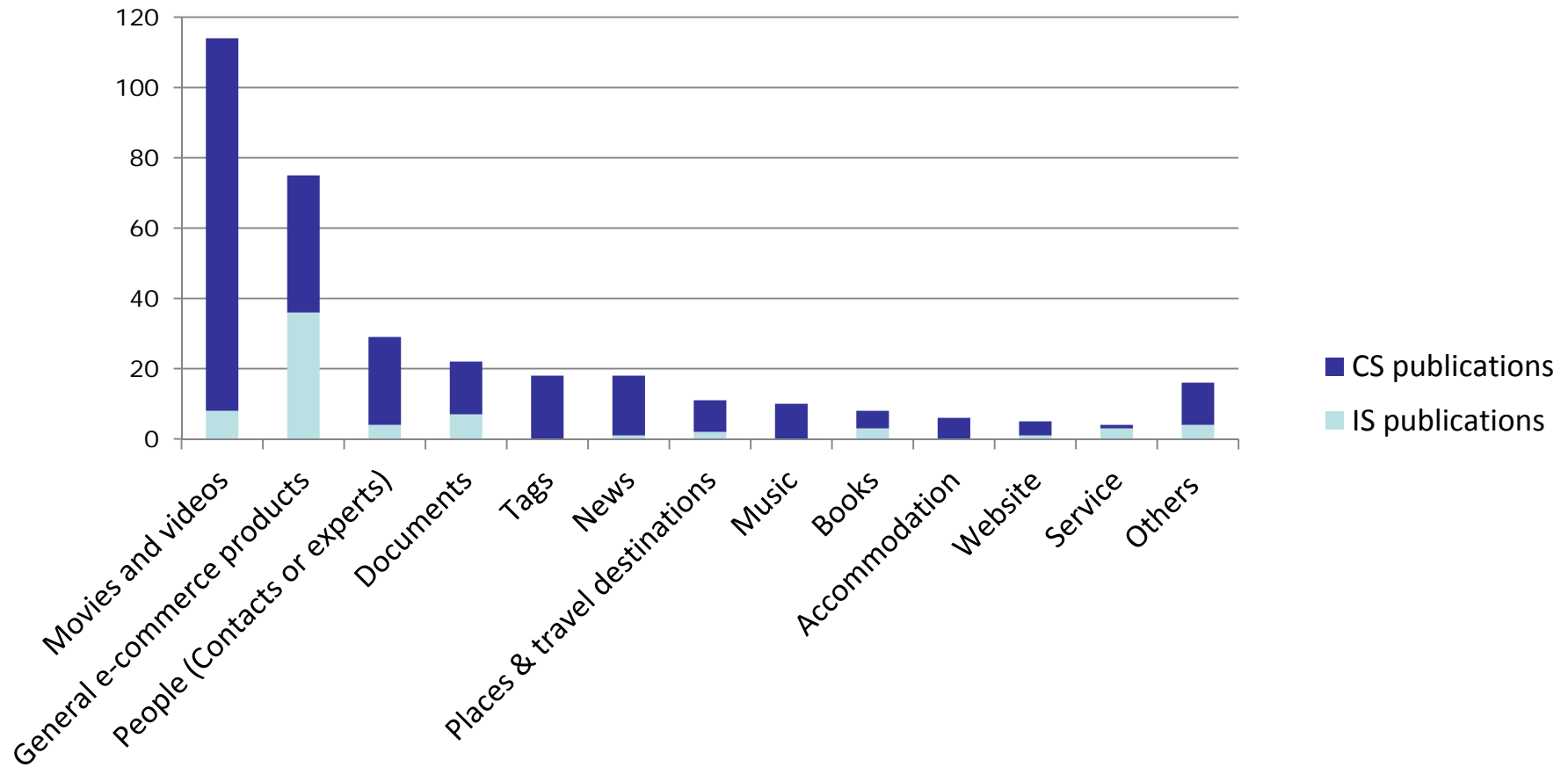
Research contributions

Type of contribution	IS outlets	CS outlets
Technical artifacts (i.e. novel algorithms)	24 (36,9%)	189 (71,3%)
Empirical research	21 (32,3%)	18 (6,8%)
Both	9 (13,8%)	43 (16,2%)
Other	11 (16,9%)	15 (5,7%)

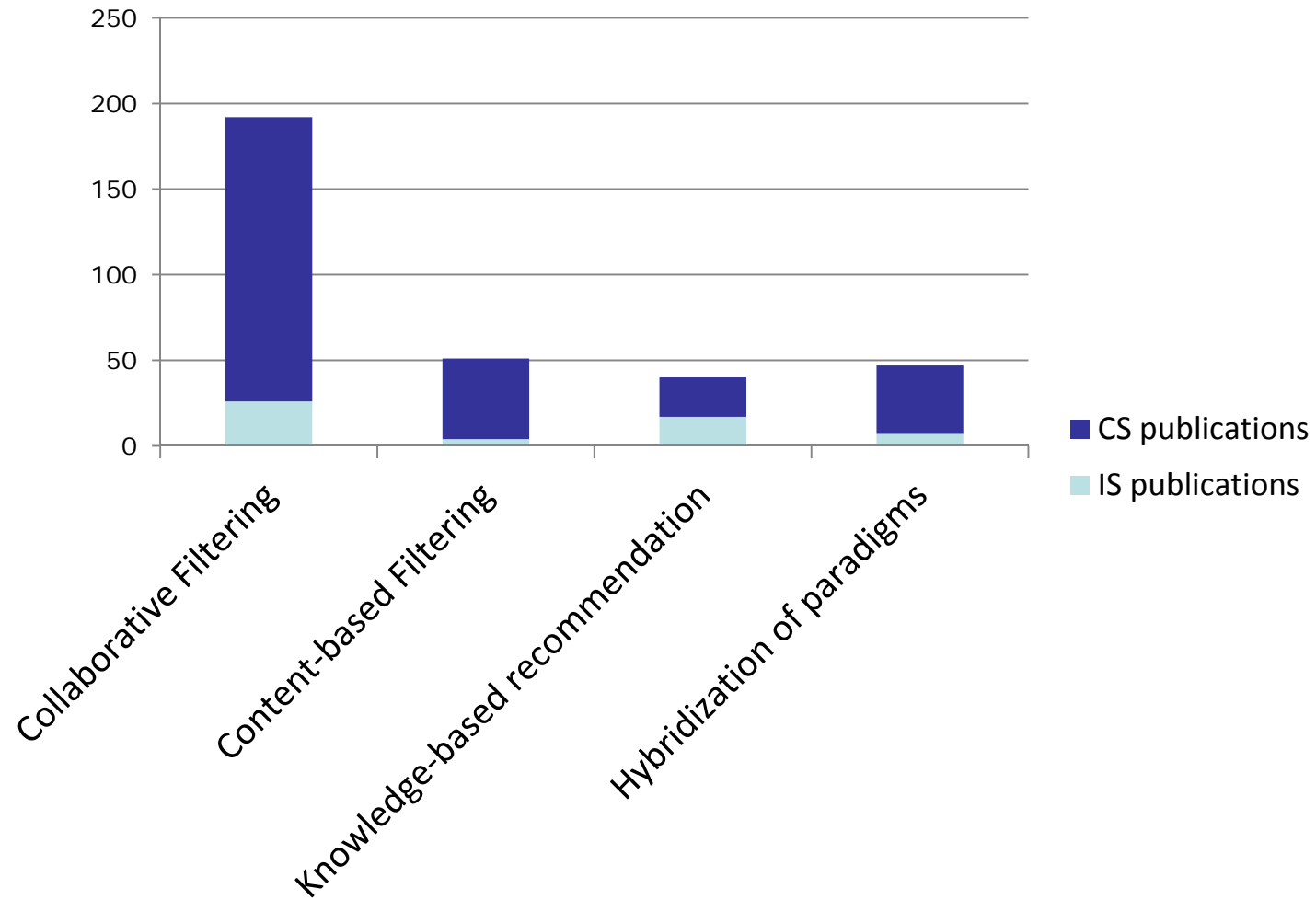
- **Topics:**

- Social and semantic web (25% of CS papers, only 6% of IS papers)
- Scalability, privacy (15% of CS papers)
- Cold-start recommendations (10% of CS papers, 5% of IS papers)
- UI design (CS 5,8%, IS 12,3%)
- Transparency (CS 6,8%, IS 10,8%)

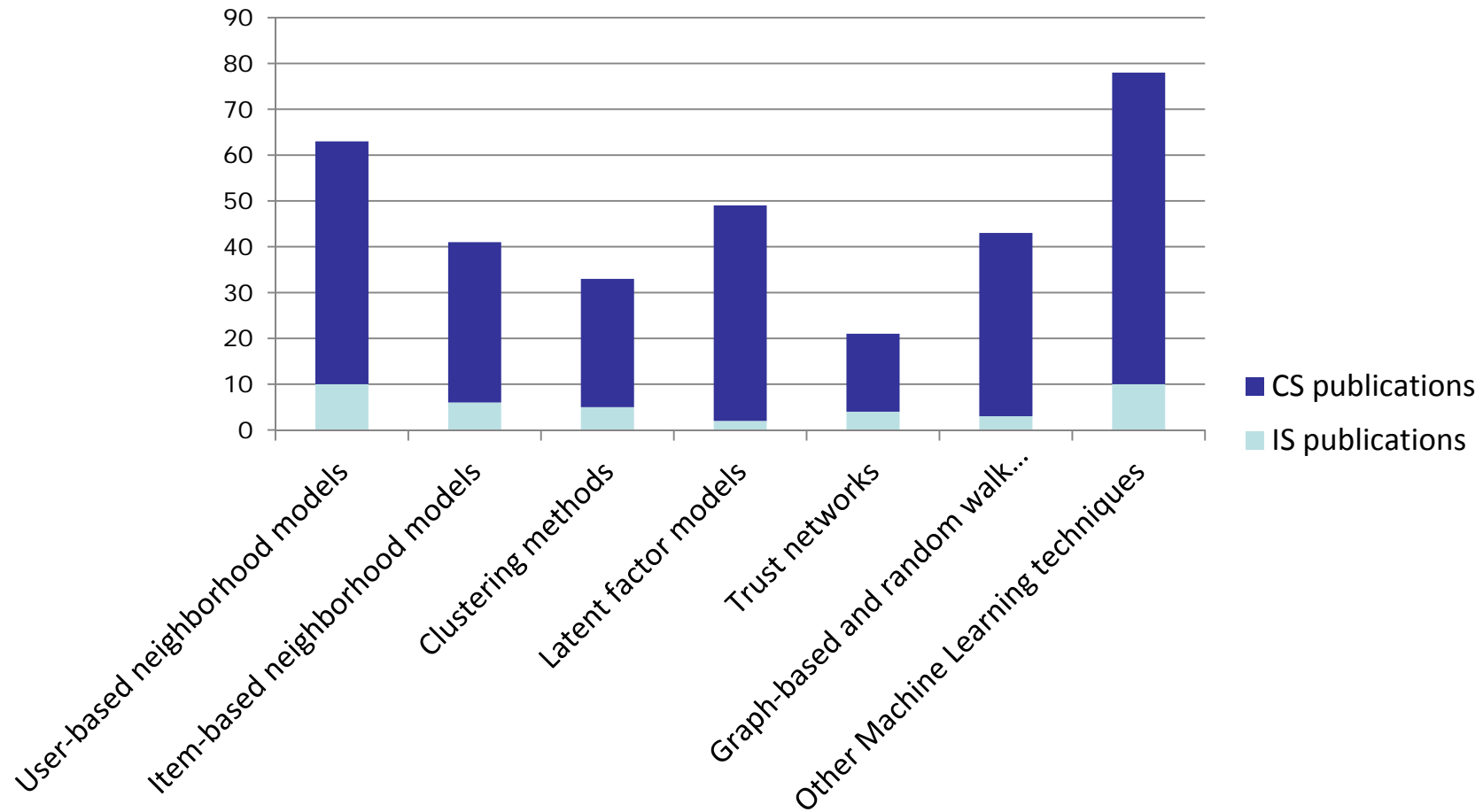
Recommended items



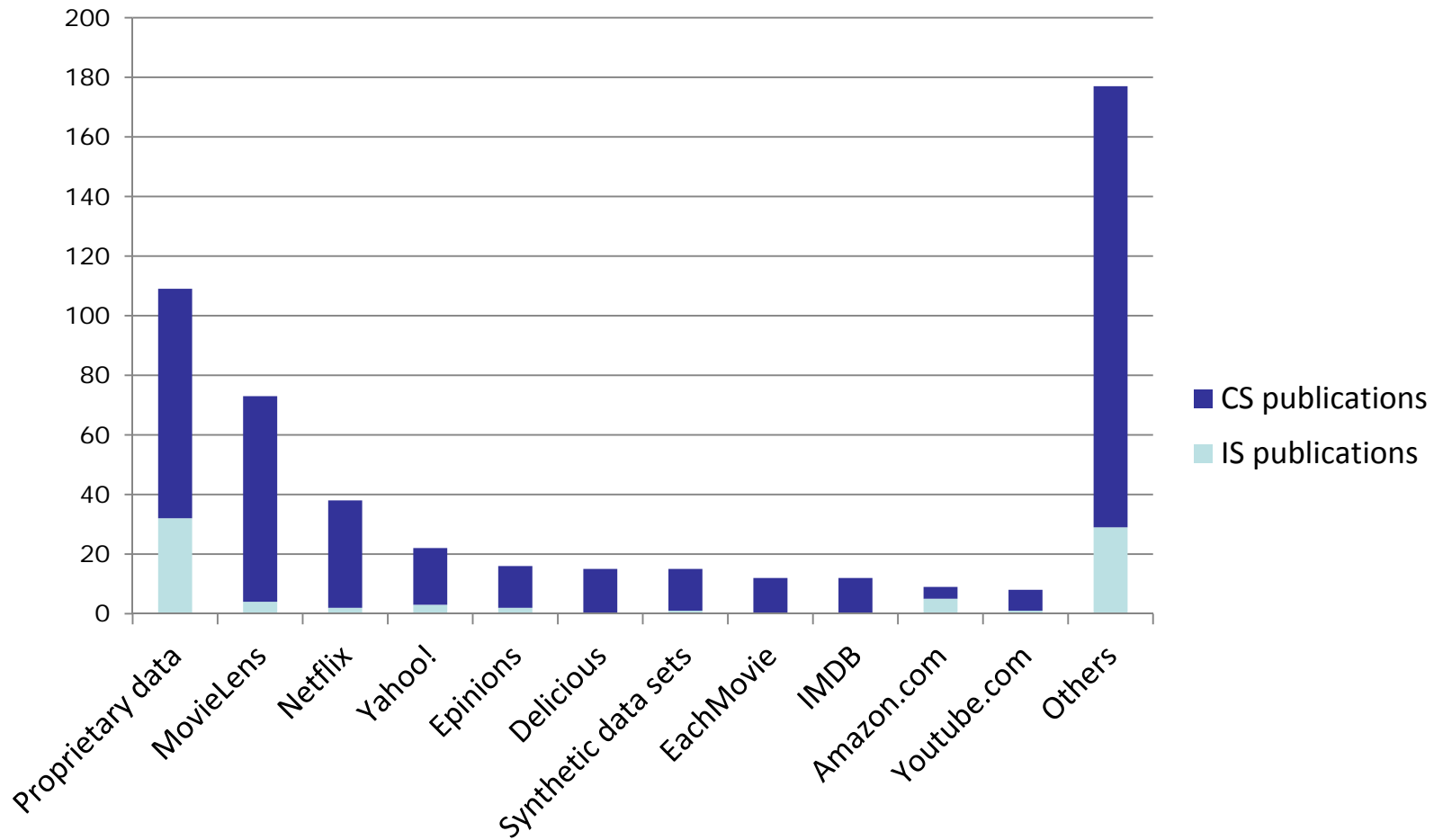
Recommendation paradigms



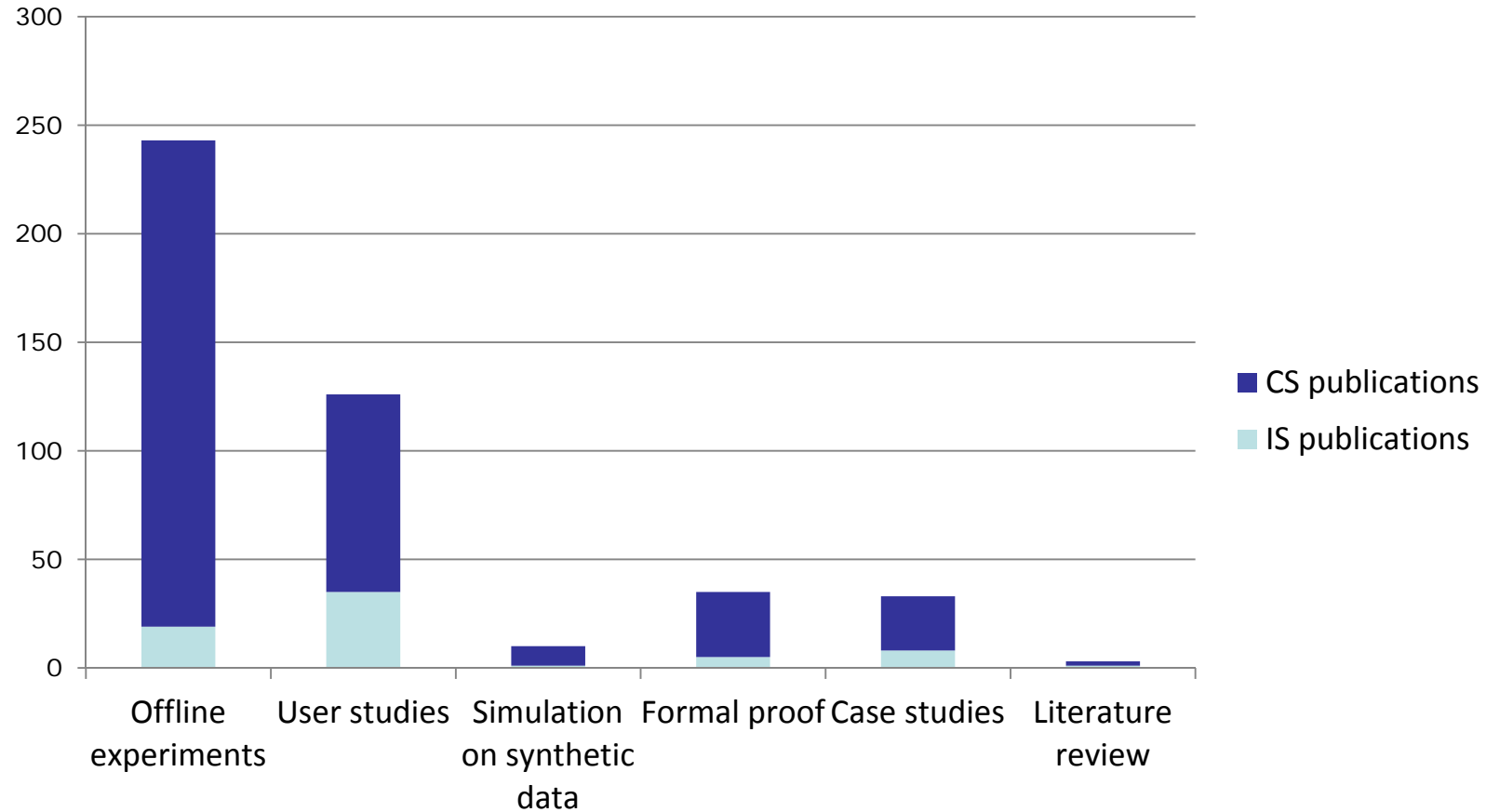
Most popular CF techniques



Most popular datasets



Methodologies



Publication outlets 2/2

	IS outlets	CS outlets
IR measures		
Precision and Recall	12	115
F1	2	20
Rank measures (e.g. NDCG)	9	27
ROC curve	1	11
Area under ROC	0	8
ML measures		
Mean absolute Error	6	57
Root Mean Squared Error	0	49
Application quality		
Computation time	2	28
Coverage metrics	2	28
Decision support quality		
Perceived utility or user satisfaction	11	7
Online conversion	3	12
Diversity metrics	0	10

Agenda

- What is the current state-of-practice?

- „How to“ from different perspectives:
 - **Empirical research principles**
 - Information Retrieval
 - Machine Learning
 - HCI and Decision Support

- Outlook

Empirical research

- **Characterizing dimensions:**
 - Who is the **subject** that is in the focus of research?
 - What **research methods** are applied?
 - In which **setting** does the research take place?

Subject	Online customers, students, historical online sessions, computers, ...
Research method	Experiments, quasi-experiments, non-experimental research
Setting	Lab, real-world scenarios

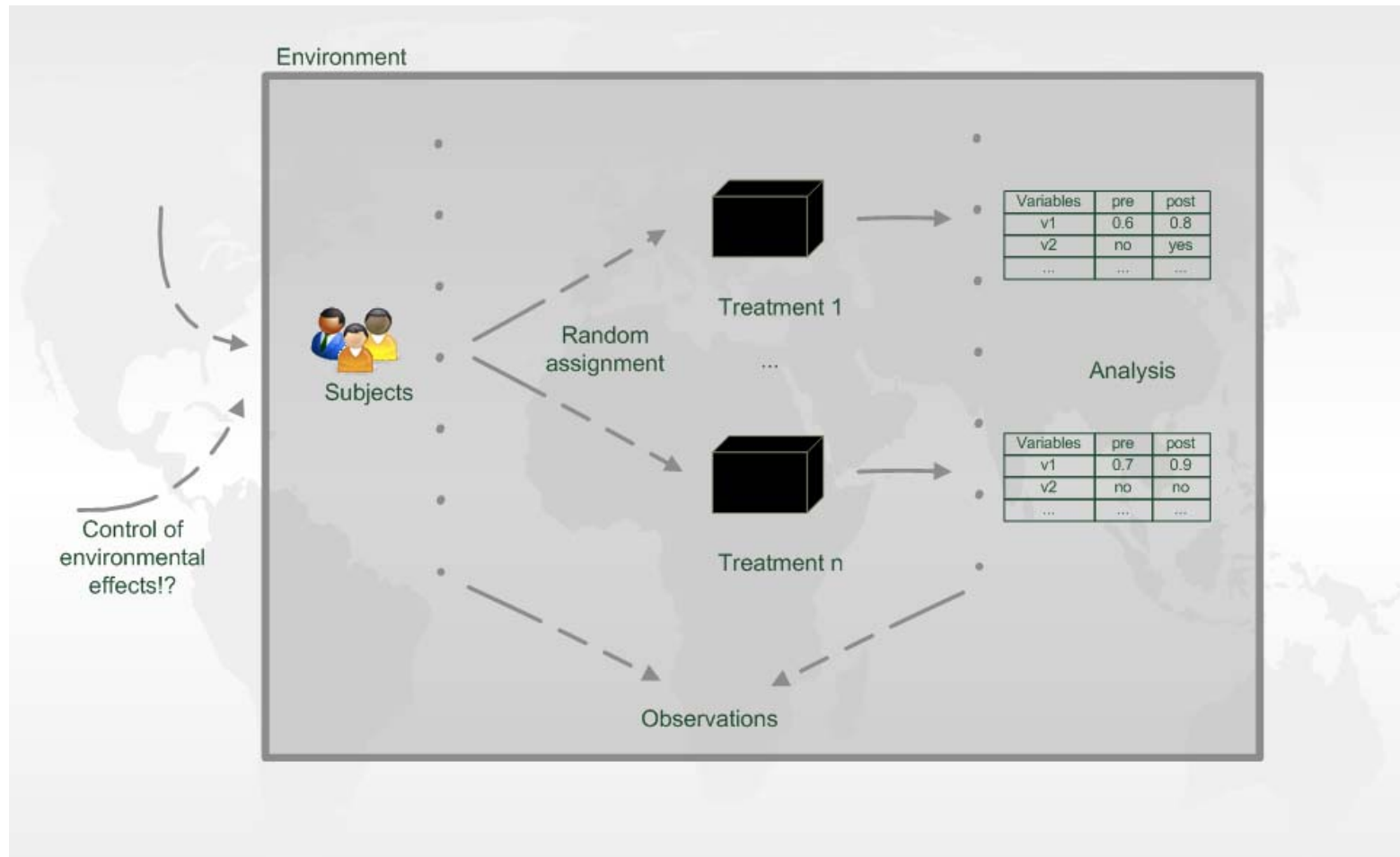
Evaluation settings

- **Lab studies**
 - Expressly created for the purpose of the study
 - Extraneous variables can be controlled more easy by selecting study participants
 - But doubts may exist about participants motivated by money or prizes
- **Participants should behave as they would in a real-world enviroment**
- **Field studies**
 - Conducted in an preexisting real-world enviroment
 - Users are intrinsically motivated to use a system

Research methods

- **Experimental vs. non-experimental (observational) research methods**
 - Experiment (test, trial):
 - *"An experiment is a study in which at least one variable is manipulated and units are randomly assigned to different levels or categories of manipulated variable(s)."*
 - Units: users, historic sessions, ...
 - Manipulated variable: type of RS, groups of recommended items, explanation strategies ...
 - Categories of manipulated variable(s): content-based RS, collaborative RS

Experiment designs



Agenda

- What is the current state-of-practice?

- „How to“ from different perspectives:
 - Empirical research principles
 - **Information Retrieval**
 - Machine Learning
 - HCI and Decision Support

- Outlook

Evaluation in information retrieval (IR)

- **Historical Cranfield collection (late 1950s)**
 - 1,398 journal article abstracts
 - 225 queries
 - Exhaustive relevance judgements (over 300K)
- **Ground truth established by human domain experts**

		Reality	
		Actually Good	Actually Bad
Prediction	Rated Good	True Positive (tp)	False Positive (fp)
	Rated Bad	False Negative (fn)	True Negative (tn)

All recommended items

All good items

Metrics: Precision and Recall

- **Recommendation is viewed as information retrieval task:**
 - Retrieve (recommend) all items which are predicted to be “good”.
- **Precision: a measure of exactness, determines the fraction of relevant items retrieved out of all items retrieved**
 - E.g. the proportion of recommended movies that are actually good

$$Precision = \frac{tp}{tp + fp} = \frac{|good\ movies\ recommended|}{|all\ recommendations|}$$



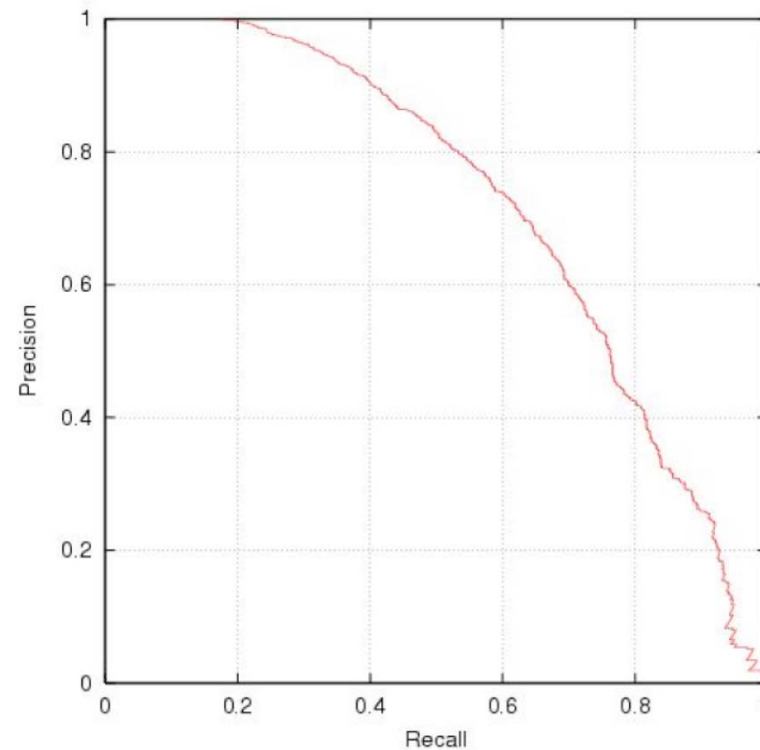
- **Recall: a measure of completeness, determines the fraction of relevant items retrieved out of all relevant items**
 - E.g. the proportion of all good movies recommended

$$Recall = \frac{tp}{tp + fn} = \frac{|good\ movies\ recommended|}{|all\ good\ movies|}$$



Precision vs. Recall

- E.g. typically when a recommender system is tuned to increase precision, recall decreases as a result (or vice versa)



F₁ Metric

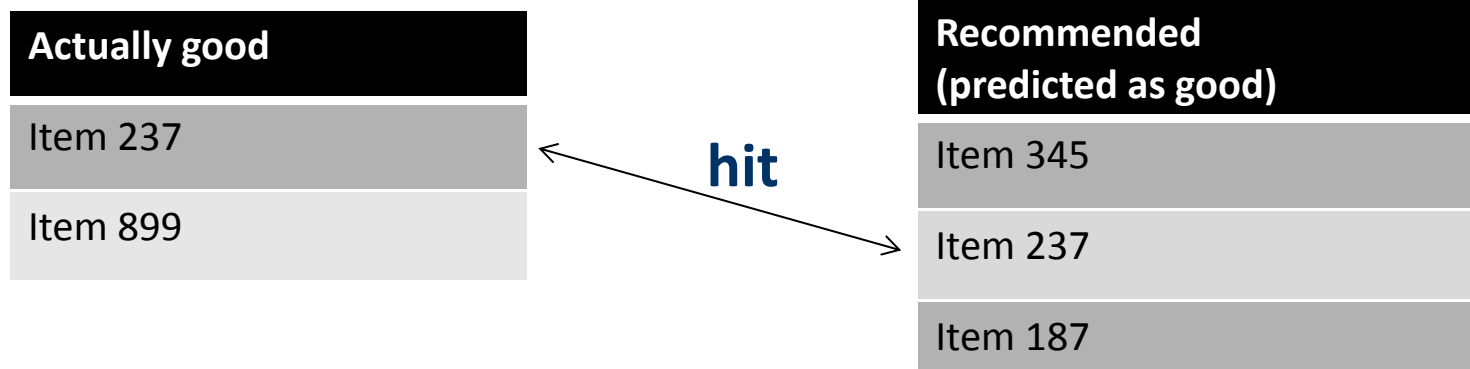
- **The F₁ Metric attempts to combine Precision and Recall into a single value for comparison purposes.**
 - May be used to gain a more balanced view of performance

$$F_1 = 2 \cdot \frac{\textit{precision} \cdot \textit{recall}}{\textit{precision} + \textit{recall}}$$

- **The F₁ Metric gives equal weight to precision and recall**
 - Other F_β metrics weight recall with a factor of β.

Metrics: Rank position matters

For a user:



- **Rank metrics extend recall and precision to take the positions of correct items in a ranked list into account**
 - Relevant items are more useful when they appear earlier in the recommendation list
 - Particularly important in recommender systems as lower ranked items may be overlooked by users

Metrics: Rank Score

- **Rank Score extends the recall metric to take the positions of correct items in a ranked list into account**
 - Particularly important in recommender systems as lower ranked items may be overlooked by users
- **Rank Score is defined as the ratio of the Rank Score of the correct items to best theoretical Rank Score achievable for the user, i.e.**

$$\text{rankscore} = \frac{\text{rankscore}_p}{\text{rankscore}_{\max}}$$

$$\text{rankscore}_p = \sum_{i \in h} 2^{-\frac{\text{rank}(i)-1}{\alpha}}$$

$$\text{rankscore}_{\max} = \sum_{i=1}^{|T|} 2^{-\frac{i-1}{\alpha}}$$

Where:

- h is the set of correctly recommended items, i.e. hits
- rank returns the position (rank) of an item
- T is the set of all items of interest
- α is the *ranking half life*, i.e. an exponential reduction factor

Metrics: Liftindex

- Assumes that ranked list is divided into 10 equal deciles S_i , where

$$\sum_{i=1}^{10} S_i = |h|$$

- Linear reduction factor

- Liftindex:**

$$\text{liftindex} = \begin{cases} \frac{1 \times S_1 + 0.9 \times S_2 + \dots + 0.1 \times S_{10}}{\sum_{i=1}^{10} S_i} & : \text{ if } |h| > 0 \\ 0 & : \text{ else} \end{cases}$$

» h is the set of correct hits

Metrics: Normalized Discounted Cumulative Gain

- **Discounted cumulative gain (DCG)**

- Logarithmic reduction factor

$$DCG_{pos} = rel_1 + \sum_{i=2}^{pos} \frac{rel_i}{\log_2 i}$$

Where:

- pos denotes the position up to which relevance is accumulated
- rel_i returns the relevance of recommendation at position i

- **Idealized discounted cumulative gain (IDCG)**

- Assumption that items are ordered by decreasing relevance

$$IDCG_{pos} = rel_1 + \sum_{i=2}^{|h|-1} \frac{rel_i}{\log_2 i}$$

- **Normalized discounted cumulative gain (nDCG)**

- Normalized to the interval [0..1]

$$nDCG_{pos} = \frac{DCG_{pos}}{IDCG_{pos}}$$

Example

- **Assumptions:**

- $|T| = 3$
- Ranking half life (alpha) = 2

Rank	Hit?
1	
2	X
3	X
4	X
5	

$$rankscore = \frac{rankscore_p}{rankscore_{max}} \approx 0.71$$

$$nDCG_5 \frac{DCG_5}{IDCG_5} \approx 0.81$$

$$liftindex = \frac{0.8 \times 1 + 0.6 \times 1 + 0.4 \times 1}{3} = 0.6$$

$$rankscore_p = \frac{1}{2^2} + \frac{1}{2^3} + \frac{1}{2^4} = 1.56$$

$$rankscore_{max} = \frac{1}{2^1} + \frac{1}{2^2} + \frac{1}{2^3} = 2.21$$

$$DCG_5 = \frac{1}{\log_2 2} + \frac{1}{\log_2 3} + \frac{1}{\log_2 4} = 2.13$$

$$IDCG_5 = 1 + \frac{1}{\log_2 2} + \frac{1}{\log_2 3} = 2.63$$

Example cont.

- Reducing the ranking half life (alpha) = 1

Rank	Hit?
1	
2	X
3	X
4	X
5	

$$rankscore = \frac{rankscore_p}{rankscore_{max}} = 0.5$$

$$rankscore_p = \frac{1}{2^1} + \frac{1}{2^1} + \frac{1}{2^1} = 0.875$$

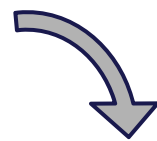
$$rankscore_{max} = \frac{1}{2^1} + \frac{1}{2^1} + \frac{1}{2^1} = 1.75$$

Rankscore (exponential reduction) < Liftscore (linear red.) < NDCG (log. red.)

Average Precision

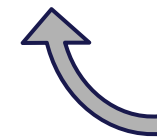
- **Average Precision (AP)** is a ranked precision metric that places emphasis on highly ranked correct predictions (hits)
- Essentially it is the average of precision values determined after each successful prediction, i.e.

Rank	Hit?
1	
2	X
3	X
4	X
5	



$$AP = \frac{1}{3} \left(\frac{1}{2} + \frac{2}{3} + \frac{3}{4} \right) = \frac{23}{36} \approx 0.639$$

$$AP = \frac{1}{3} \left(\frac{1}{1} + \frac{2}{4} + \frac{3}{5} \right) = \frac{21}{30} = 0.7$$



Rank	Hit?
1	X
2	
3	
4	X
5	X

Agenda

- What is the current state-of-practice?

- „How to“ from different perspectives:
 - Empirical research principles
 - Information Retrieval
 - **Machine Learning**
 - HCI and Decision Support

- Outlook

RS from a ML perspective

- Recommendation is concerned with learning from noisy observations (x,y) , where $f(x) = \hat{y}$ has to be determined such that $\sum_{\hat{y}} (\hat{y} - y)^2$ is minimal.
- A huge variety of different learning strategies have been applied trying to estimate $f(x)$
 - Non parametric neighborhood models
 - MF models, SVMs, Neural Networks, Bayesian Networks,...

Error measures

- **Datasets with items rated by users**

- MovieLens datasets 100K-10M ratings
- Netflix 100M ratings



- **Historic user ratings constitute ground truth**

- **Metrics measure error rate**

- Mean Absolute Error (*MAE*) computes the deviation between predicted ratings and actual ratings

$$MAE = \frac{1}{n} \sum_{i=1}^n |p_i - r_i|$$

- Root Mean Square Error (*RMSE*) is similar to *MAE*, but places more emphasis on larger deviation

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (p_i - r_i)^2}$$

Data sparsity

- **Natural datasets include historical interaction records of real users**
 - Explicit user ratings
 - Datasets extracted from web server logs (implicit user feedback)

- **Sparsity of a dataset is derived from ratio of empty and total entries in the user-item matrix:**
 - Sparsity = $1 - |R|/|I| \cdot |U|$
 - R = ratings
 - I = items
 - U = users

Example

Nr.	UserID	MovieID	Rating (r_i)	Prediction (p_i)	$ p_i - r_i $	$(p_i - r_i)^2$
1	1	134	5	4.5	0.5	0.25
2	1	238	4	5	1	1
3	1	312	5	5	0	0
4	2	134	3	5	2	4
5	2	767	5	4.5	0.5	0.25
6	3	68	4	4.1	0.1	0.01
7	3	212	4	3.9	0.1	0.01
8	3	238	3	3	0	0
9	4	68	4	4.2	0.2	0.04
10	4	112	5	4.8	0.2	0.04
					4.6	5.6

- MAE = 0.46

- RMSE = 0.75

Removing line nr. 4

- MAE = 0.29

- RMSE = 0.42

Removing lines 1,2,4,5

- MAE = 0.1

- RMSE = 0.13

Dilemma of establishing ground truth

- IR measures are frequently applied, however:

Offline experimentation	Online experimentation
Ratings, transactions	Ratings, feedback
Historic session (not all recommended items are rated)	Live interaction (all recommended items are rated)
Ratings of unrated items unknown, but interpreted as "bad" (default assumption, user tend to rate only good items)	"Good/bad" ratings of not recommended items are unknown
If default assumption does not hold: True positives may be too small False negatives may be too small	False/true negatives cannot be determined
Precision may increase Recall may vary	Precision ok Recall questionable

Results from offline experimentation have limited predictive power for online user behavior.

Offline experimentation

- **Netflix competition**

- Web-based movie rental
- Prize of \$1,000,000 for accuracy improvement (RMSE) of 10% compared to own Cinematch system.

- **Historical dataset**

- ~480K users rated ~18K movies on a scale of 1 to 5
- ~100M ratings
- Last 9 ratings/user withheld
 - Probe set – for teams for evaluation
 - Quiz set – evaluates teams' submissions for leaderboard
 - Test set – used by Netflix to determine winner

Methodology

- **Setting to ensure internal validity:**
 - One randomly selected share of known ratings (**training set**) used as input to train the algorithm and build the model
 - Model allows the system to compute recommendations at runtime
 - Remaining share of withheld ratings (**testing set**) required as ground truth to evaluate the model's quality
 - To ensure the reliability of measurements the random split, model building and evaluation steps are repeated several times

 - **N-fold cross validation is a stratified random selection procedure**
 - N disjunct fractions of known ratings with equal size ($1/N$) are determined
 - N repetitions of the model building and evaluation steps, where each fraction is used exactly once as a testing set while the other fractions are used for training
 - Setting N to 5 or 10 is popular
-

Analysis of results

- **Are observed differences statistically meaningful or due to chance?**
 - Standard procedure for testing the statistical significance of two deviating metrics is the pairwise analysis of variance (ANOVA)
 - Null hypothesis H_0 : observed differences have been due to chance
 - If outcome of test statistics rejects H_0 , significance of findings can be reported

- **Practical importance of differences?**
 - Size of the effect and its practical impact
 - External validity or generalizability of the observed effects

Agenda

- What is the current state-of-practice?

- „How to“ from different perspectives:
 - Empirical research principles
 - Information Retrieval
 - Machine Learning
 - **HCI and Decision Support**

- Outlook

Online experimentation

- Effectiveness of different algorithms for recommending cell phone games
[Jannach, Hegelich 09]
- Involved 150,000 users on a commercial mobile internet portal
- Comparison of recommender methods
- Random assignment of users to a specific method



Experimental Design

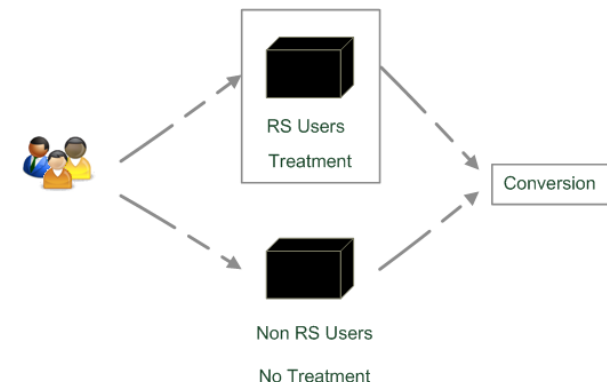
- **A representative sample 155,000 customers were extracted from visitors to site during the evaluation period**
 - These were split into 6 groups of approximately 22,300 customers
 - Care was taken to ensure that customer profiles contained enough information (ratings) for all variants to make a recommendation
 - Groups were chosen to represent similar customer segments
 - **A catalog of 1,000 games was offered**
 - **A five-point ratings scale ranging from -2 to +2 was used to rate items**
 - Due to the low number of explicit ratings, a click on the “details” link for a game was interpreted as an implicit “0” rating and a purchase as a “1” rating
 - **Hypotheses on personalized vs. non-personalized recommendation techniques and their potential to**
 - Increase conversion rate (i.e. the share of users who become buyers)
 - Stimulate additional purchases (i.e. increase the average shopping basket size)
-

Non-experimental research

- **Quasi-experiments**
 - Lack random assignments of units to different treatments
 - **Non-experimental / observational research**
 - Surveys / Questionnaires
 - Longitudinal research
 - Observations over long period of time
 - E.g. customer life-time value, returning customers
 - Case studies
 - Focus on answering research questions about how and why
 - E.g. answer questions like: *How recommendation technology contributed to Amazon.com's becomes the world's largest book retailer?*
 - Focus group
 - Interviews
 - Think aloud protocols
-

Quasi-experimental

- **SkiMatcher Resort Finder introduced by Ski-Europe.com to provide users with recommendations based on their preferences**
- **Conversational RS**
 - question and answer dialog
 - matching of user preferences with knowledge base
- **Delgado and Davidson evaluated the effectiveness of the recommender over a 4 month period in 2001**
 - Classified as a quasi-experiment as users decide for themselves if they want to use the recommender or not



SkiMatcher Results

	July	August	September	October
Unique Visitors	10,714	15,560	18,317	24,416
• SkiMatcher Users	1,027	1,673	1,878	2,558
• Non-SkiMatcher Users	9,687	13,887	16,439	21,858
Requests for Proposals	272	506	445	641
• SkiMatcher Users	75	143	161	229
• Non-SkiMatcher Users	197	363	284	412
Conversion	2.54%	3.25%	2.43%	2.63%
• SkiMatcher Users	7.30%	8.55%	8.57%	8.95%
• Non-SkiMatcher Users	2.03%	2.61%	1.73%	1.88%
Increase in Conversion	359%	327%	496%	475%

[Delgado and Davidson, ENTER 2002]

Interpreting the Results

- **The nature of this research design means that questions of causality cannot be answered (lack of random assignments), such as**
 - Are users of the recommender systems more likely convert?
 - Does the recommender system itself cause users to convert?

Some hidden exogenous variable might influence the choice of using RS as well as conversion.

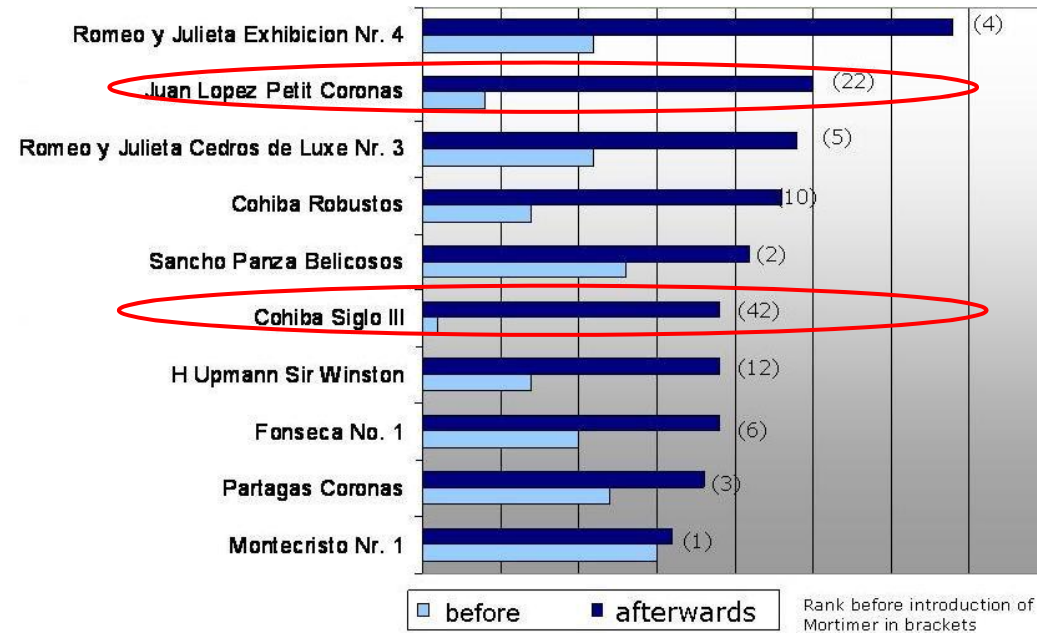
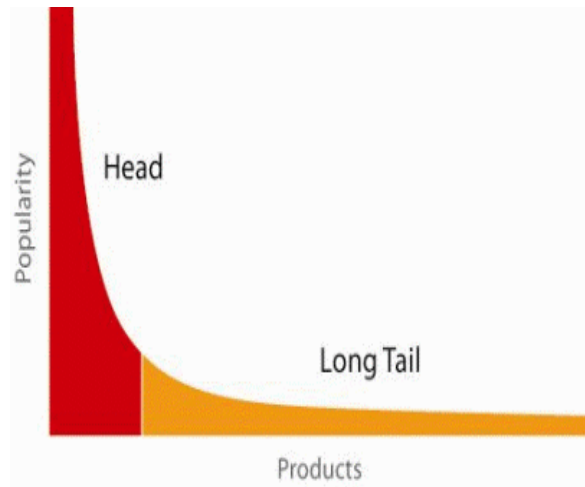
- **However, significant correlation between using the recommender system and making a request for a proposal**

- **Size of effect has been replicated in other domains**
 - Tourism [Jannach et al., JITT 2009]
 - Electronic consumer products

Observational research

- Increased demand in niches/long tail products

- Books ranked above 250.000 represent >29% of sales at Amazon, approx. 2.3 million books [Brynjolfsson et al., Mgt. Science, 2003]
- Ex post from webshop data [Zanker et al., EC-Web, 2006]



Agenda

- What is the current state-of-practice?

- „How to“ from different perspectives:
 - Empirical research principles
 - Information Retrieval
 - Machine Learning
 - HCI and Decision Support

- **Outlook**

Reality check regarding F_1 and accuracy measures for RS

- **Real value lies in increasing conversions**
 - ...and satisfaction with bought items, low churn rate
- **Some reasons why it might be a fallacy to think F_1 on historical data is a good estimate for real conversion:**
 - Recommendation can be self-fulfilling prophecy
 - Users' preferences are not invariant, but can be constructed [ALP03]
 - Position/Rank is what counts (e.g. serial position effects)
 - Actual choices are heavily biased by the item's position [FFG+07]
 - Smaller recommendation sets increase users' confidence in decision making
 - Effect of choice overload - large sets at the same time increase choice difficulty and reduce choice satisfaction [BKW+10]
 - Inclusion of weak (dominated) items increases users' confidence
 - Replacing some recommended items by *decoy* items fosters choice towards the remaining options [TF09]

Bounded rationality



- **Framing and reference dependence, e.g.**
 - Presentation of the decision problem and its recommendations to the user (e.g. gains and losses)
 - Bias towards initial anchor point (in conversational RS)
- **Cognitive consistency theory**
 - Preferences are re-constructed in the course of decision making in order to avoid conflicts
- **Serial position effects**
 - Primacy and recency
- **Decoy effects**
 - Items below pareto frontier
 - Dominance relationships

Implications for RS



- **Preferences cannot be assumed to be stable during RS interaction**
 - Sequence, content and wording matter
- **Preference reasoning**
 - Independence of preferences may not be assumed, impacts the computation of the preference score
- **Presentation/recommendation of items**
 - Serial positions significantly influence the decision
 - Inclusion of decoy items

Which one will the majority select?



Primacy and Decoy effect

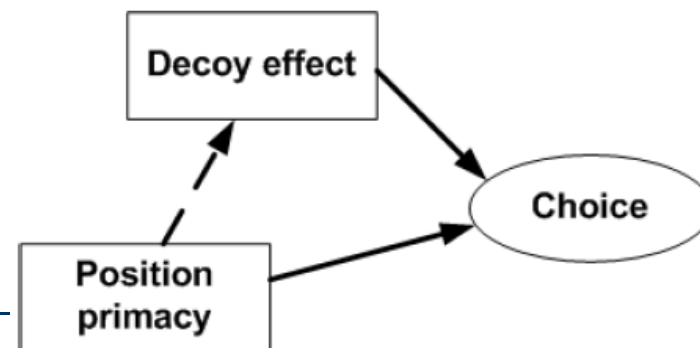
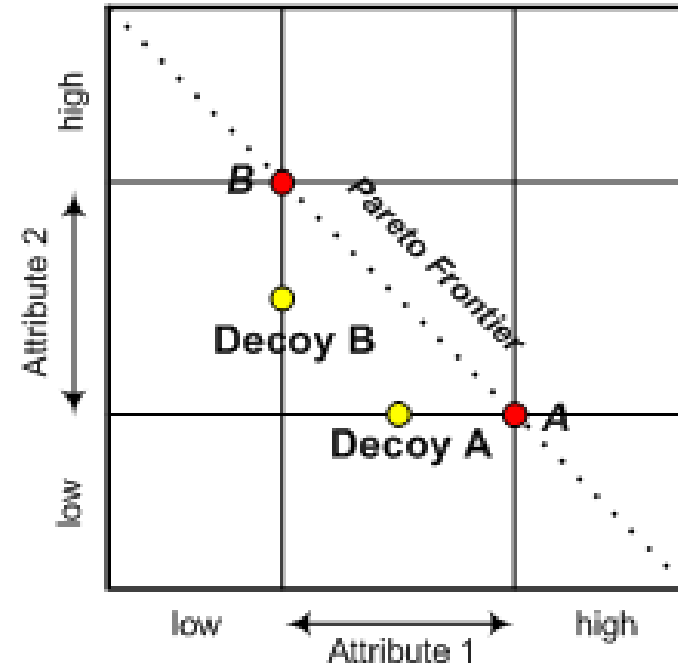
	Item A	Item B	Decoy A
Attribute 1	High 	Low	Medium
Attribute 2	Low 	High	Low

Primacy effect and opposite Decoy effect

	Item B	Item A	Decoy B
Attribute 1	Low 	High 	Low
Attribute 2	High	Low	Medium

Primacy and Decoy effect

	Item B	Item A	Decoy B
Attribute 2	High 	Low	Medium
Attribute 1	Low 	High	Low



Discussion & summary

- General principles of empirical research and current state of practice in evaluating recommendation techniques were presented
- Focus on how to perform empirical evaluations on historical datasets
- Discussion about different methodologies and metrics for measuring the accuracy or coverage of recommendations.
- Overview of which research designs are commonly used in practice.
- From a technical point of view, measuring the accuracy of predictions is a well accepted evaluation goal
 - but other aspects that may potentially impact the overall effectiveness of a recommendation system remain largely under developed.

Outlook

- **Additional topics covered by the book “Recommender Systems - An Introduction”**
 - Case study on the Mobile Internet
 - Attacks on CF Recommender Systems
 - Recommender Systems in the next generation Web (Social Web, Semantic Web)
 - Consumer decision making
 - Recommending in ubiquitous environments

- **“RS research will become much more diverse”**
 - Various forms of feedback mechanisms and preference representation
 - More focus on interfaces, interaction processes, explaining and trust-building
 - Plurality of evaluation methods complementing offline experiments

- **More focus on causal relationships**
 - When, where and how to recommend?
 - Consumer / sales psychology
 - Consumer decision making theories

Thank you for your attention!

Questions?

Questions?

Questions?

Dietmar Jannach

e-Services Research Group

Department of Computer Science

TU Dortmund, Germany

M: dietmar.jannach@tu-dortmund.de

P: +49 231 755 7272

Markus Zanker

Intelligent Systems and Business Informatics

Institute of Applied Informatics

University Klagenfurt, Austria

M: markus.zanker@aau.at

P: +43 463 2700 3705

<http://recsys.acm.org>



<http://www.recommenderbook.net>



Recommender Systems – An Introduction by

Dietmar Jannach, Markus Zanker, Alexander Felfernig and
Gerhard Friedrich

Cambridge University Press, 2010

References

-
- [Adomavicius & Tuzhilin, IEEE TKDE, 2005] Adomavicius G., Tuzhilin, A. Toward the next generation of recommender systems: a survey of the state-of-the-art and possible extensions, *IEEE TKDE*, 17(6), 2005, pp.734-749.
- [ALP03] Ariely, D., Loewenstein, G., Prelec, D. (2003) “Coherent Arbitrariness”: Stable Demand Curves Without Stable Preferences. *The Quarterly Journal of Economics*, February 2003, 73-105.
- [BKW+10] Bollen, D., Knijnenburg, B., Willemsen, M., Graus, M. (2010) Understanding Choice Overload in Recommender Systems. *ACM Recommender Systems*, 63-70.
- [Brynjolfsson et al., Mgt. Science, 2003] Brynjolfsson, E., Hu, Y., Smith, M.: Consumer Surplus in the Digital Economy: Estimating the Value of Increased Product Variety at Online Booksellers, *Management Science*, Vol 49(11), 2003, pp. 1580-1596.
- [BS97] Balabanovic, M., Shoham, Y. (1997) Fab: content-based, collaborative recommendation, *Communications of the ACM*, Vol. 40(3), pp. 66-72.
- [FFG+07] Felfernig, A., Friedrich, G., Gula, B. et al. (2007) Persuasive recommendation: serial position effects in knowledge-based recommender systems. 2nd international conference on Persuasive technology, Springer, 283–294.
- [Friedrich& Zanker, AIMag, 2011] Friedrich, G., Zanker, M.: A Taxonomy for Generating Explanations in Recommender Systems. *AI Magazine*, Vol. 32(3), 2011.
- [Jannach et al., CUP, 2010] Jannach D., Zanker M., Felfernig, A., Friedrich, G.: *Recommender Systems an Introduction*, Cambridge University Press, 2010.
- [Jannach et al., JITT, 2009] Jannach, D., Zanker, M., Fuchs, M.: Constraint-based recommendation in tourism: A multi-perspective case study, *Information Technology & Tourism*, Vol 11(2), pp. 139-156.
- [Jannach, Hegelich 09] Jannach, D., Hegelich K.: A case study on the effectiveness of recommendations in the Mobile Internet, *ACM Conference on Recommender Systems*, New York, 2009, pp. 205-208
- [Ricci et al., JITT, 2009] Mahmood, T., Ricci, F., Venturini, A.: Improving Recommendation Effectiveness by Adapting the Dialogue Strategy in Online Travel Planning. *Information Technology & Tourism*, Vol 11(4), 2009, pp. 285-302.
- [Teppan& Felfernig, CEC, 2009] Teppan, E., Felfernig, A.: Asymmetric Dominance- and Compromise Effects in the Financial Services Domain. *IEEE International Conference on E-Commerce and Enterprise Computing*, 2009, pp. 57-64
- [TF09] Teppan, E., Felfernig, A. (2009) Impacts of decoy elements on result set evaluations in knowledge-based recommendation. *Int. J. Adv. Intell. Paradigms* 1, 358–373.
-

References

- [Xiao & Benbasat, MISQ, 2007] Xiao, B., Benbasat, I.: E-Commerce Product Recommendation Agents: Use, Characteristics, and Impact, *MIS Quarterly*, Vol 31(1), pp. 137-209.
- [Zanker et al., EC-Web, 2006] Zanker, M., Bricman, M., Gordea, S., Jannach, D., Jessenitschnig, M.: Persuasive online-selling in quality & taste domains, *7th International Conference on Electronic Commerce and Web Technologies*, 2006, pp. 51-60.
- [Zanker, RecSys, 2008] Zanker M., A Collaborative Constraint-Based Meta-Level Recommender. *ACM Conference on Recommender Systems*, 2008, pp. 139-146.
- [Zanker et al., UMUI, 2009] Zanker, M., Jessenitschnig, M., Case-studies on exploiting explicit customer requirements in recommender systems, *User Modeling and User-Adapted Interaction*, Springer, 2009, pp.133-166.
- [Zanker et al., JITT, 2009] Zanker M., Jessenitschnig M., Fuchs, M.: Automated Semantic Annotations of Tourism Resources Based on Geospatial Data, *Information Technology & Tourism*, Vol 11(4), 2009, pp. 341-354.
- [Zanker et al., Constraints, 2010] Zanker M., Jessenitschnig M., Schmid W.: Preference reasoning with soft constraints in constraint-based recommender systems. *Constraints*, Springer, Vol 15(4), 2010, pp. 574-595.