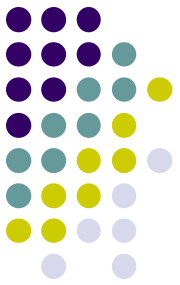


Recommender Systems - introduction

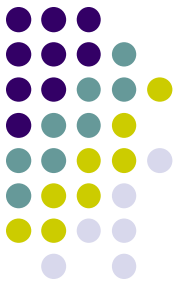
Ladislav Peška,

Department of Software Engineering,
Charles University in Prague,
Czech Republic





Background, disclaimer

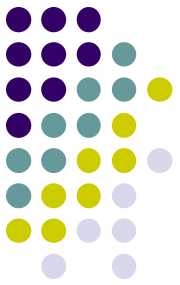


Content:

- What are recommender systems?
- Why we should love them?
- How they work?
- Challenges, problems, risks, practical deployment
- Discussion

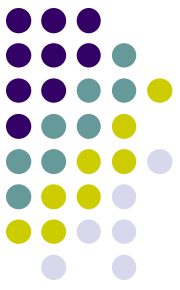
<http://www.recommenderbook.net/teaching-material/slides>

<http://ksi.mff.cuni.cz/~peska/ndbi037/recsysIntro.pdf>



You all know them, maybe you just didn't know that so far...;-)

RECOMMENDER SYSTEMS



Recommender Systems

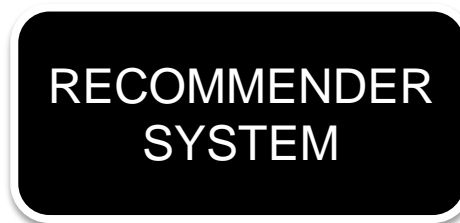
- Propose **relevant items** to the **right persons** at the **right time**
- Machine learning application
- Expose otherwise hard to find, unknown items
- Complementary to the catalogues, search engines etc.
- **„Win-win strategy“**

User Feedback
rating, clickstream,
time on page, buys...



User, Object Profiles
Object attributes

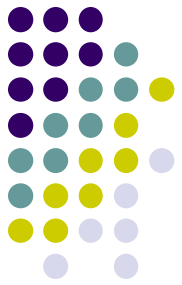
(Context)
Time, location,
Possible choices...



Top-K Recommended objects



Recommender Systems are everywhere 😊



- Movies, news, books, e-commerce, web/site-wide search, social networks...

Recent Automatic Recommendations

1-34 of 202 (next) titles | covers | shelf

Feb 2 Feb 2 Feb 2 Feb 2 Feb 2

Megnézem még egyszer

DÉMOPHOBIA - MLÝNEK
MIROSLAV NOVOTNÝ
2 918 megtekintés · 2 éve

Xindi X - Čecháček a totáček
XindiXOfficialVEVO
1 996 303 megtekintés · 1 éve

Žalman a spol Jantarová země
Monty z
192 484

Žalman & Sýr
Všech vandi

MovieLens recommends these movies

top picks

found 747 movies. show search tools

Star Wars 1977 PG 121 min

Star Wars: Episode 1980 PG 124 min

Raiders of the Lost 1981 PG 115 min

Return of the Jedi 1983 PG 135 min

Iron Man 2008 PG-13 126 min

Inside Out 2015 94 min

Band of Brothers 2001 705 min

Star Trek 2009 PG-13 127 min

People who viewed this item also viewed

Ajánlott

DÉMOPHOBIA - PLZEŇSKÉ POVĚSTI, PÍSNĚ A JINĚ...
MIROSLAV NOVOTNÝ
43:26

Vangelis - The Collection (2012) (CD 1)
1:17:04

Custom Gaming PC Desktop Computer...
125,594.41 HUF
Buy It Now + 16,892.31 HUF

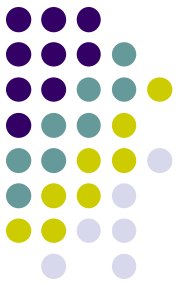
Xindi X
XindiXOfficialVEVO
268 182 megtekintés · 4 hete

Vangelis
1:17:04

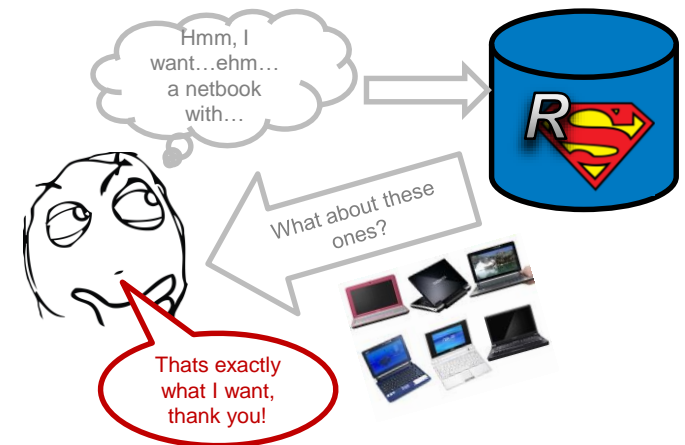
AMD Quad Core Gaming Desktop PC
100,696.50 HUF
Buy It Now + 16,892.31 HUF

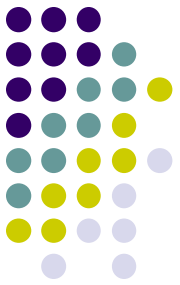
AMD Quad Core Gaming Desktop PC
104,895.10 HUF
Buy It Now + 16,892.31 HUF

Motivation



- Recommender systems
 - Are complementary to classical GUI elements (*cats. hierarchy, search queries,...*)
 - Most effective when users don't know exactly what they want
 - Or their intent is hard to express
 - Serendipity (i.e. a lucky hit)
 - Should both increase user satisfaction and provider's success metrics



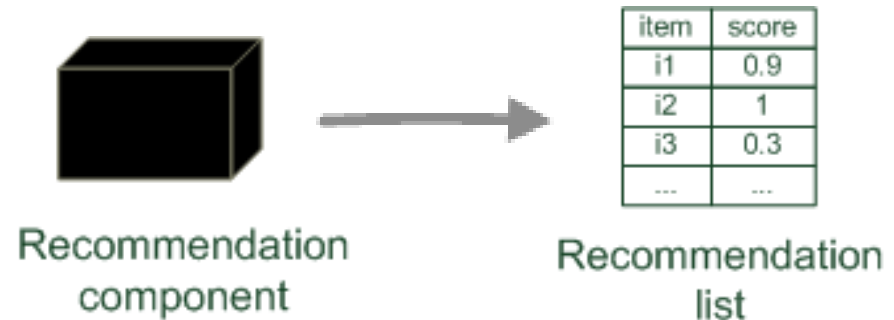


How do they work?

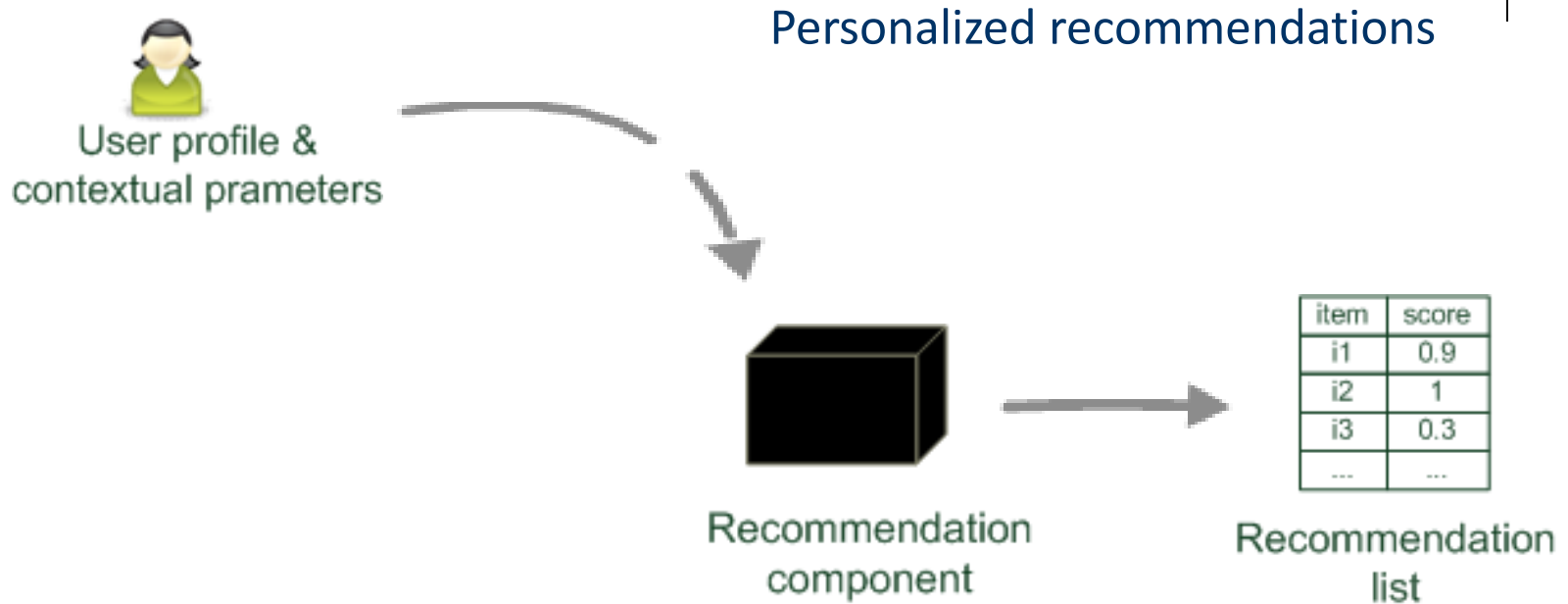
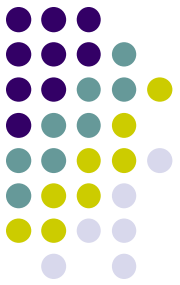
Paradigms of recommender systems



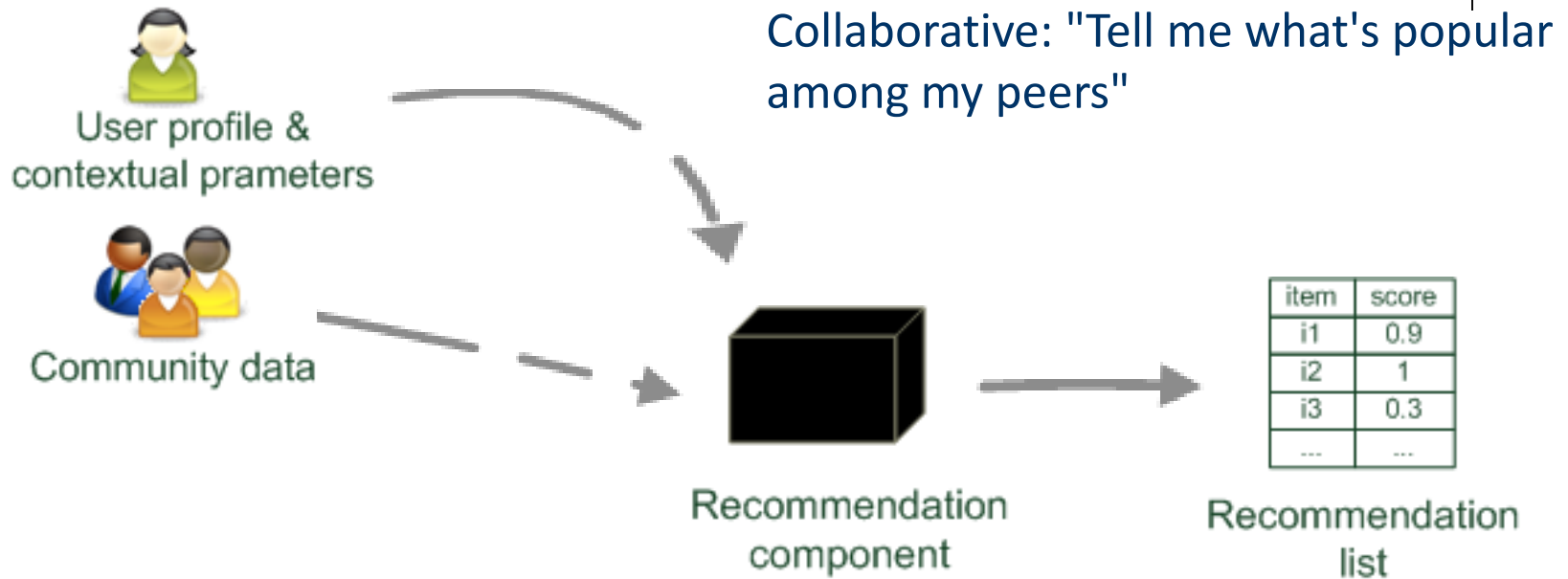
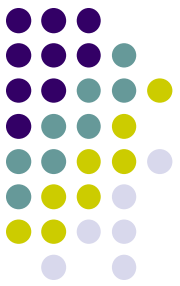
Recommender systems reduce information overload by estimating relevance



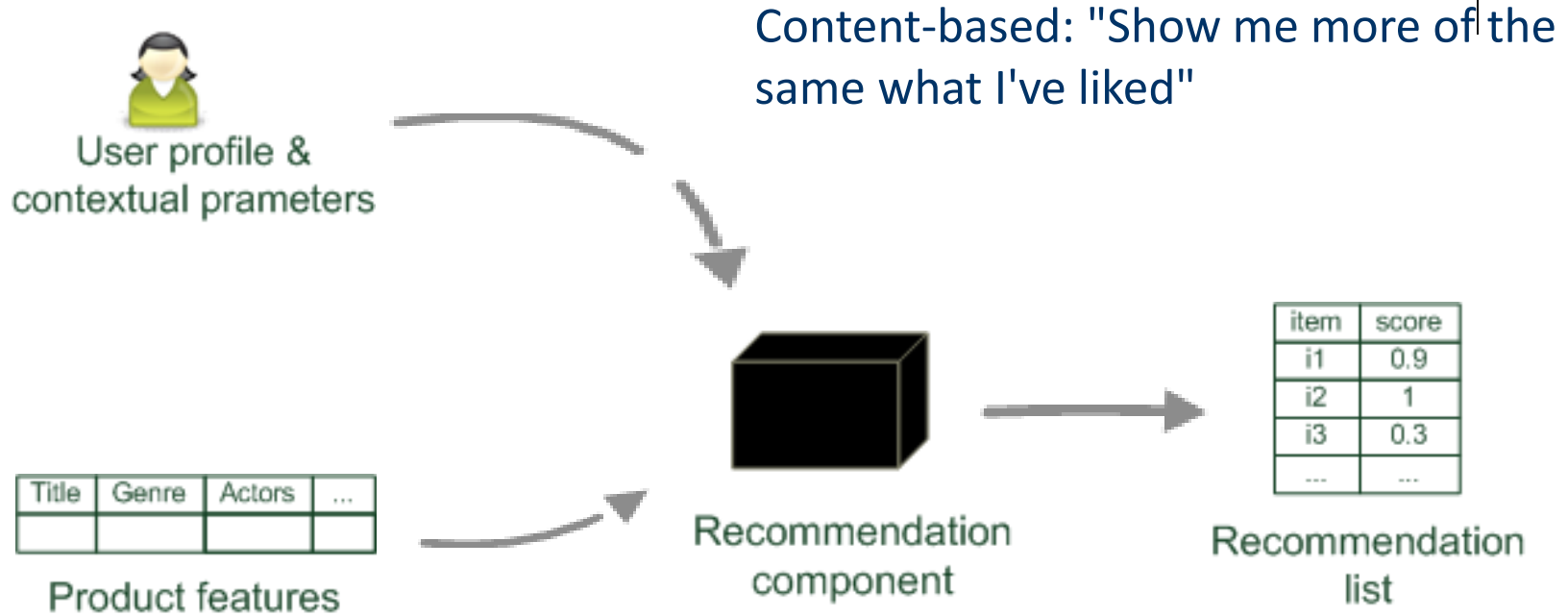
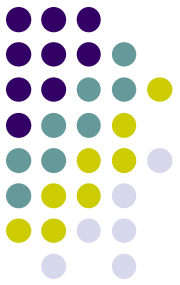
Paradigms of recommender systems



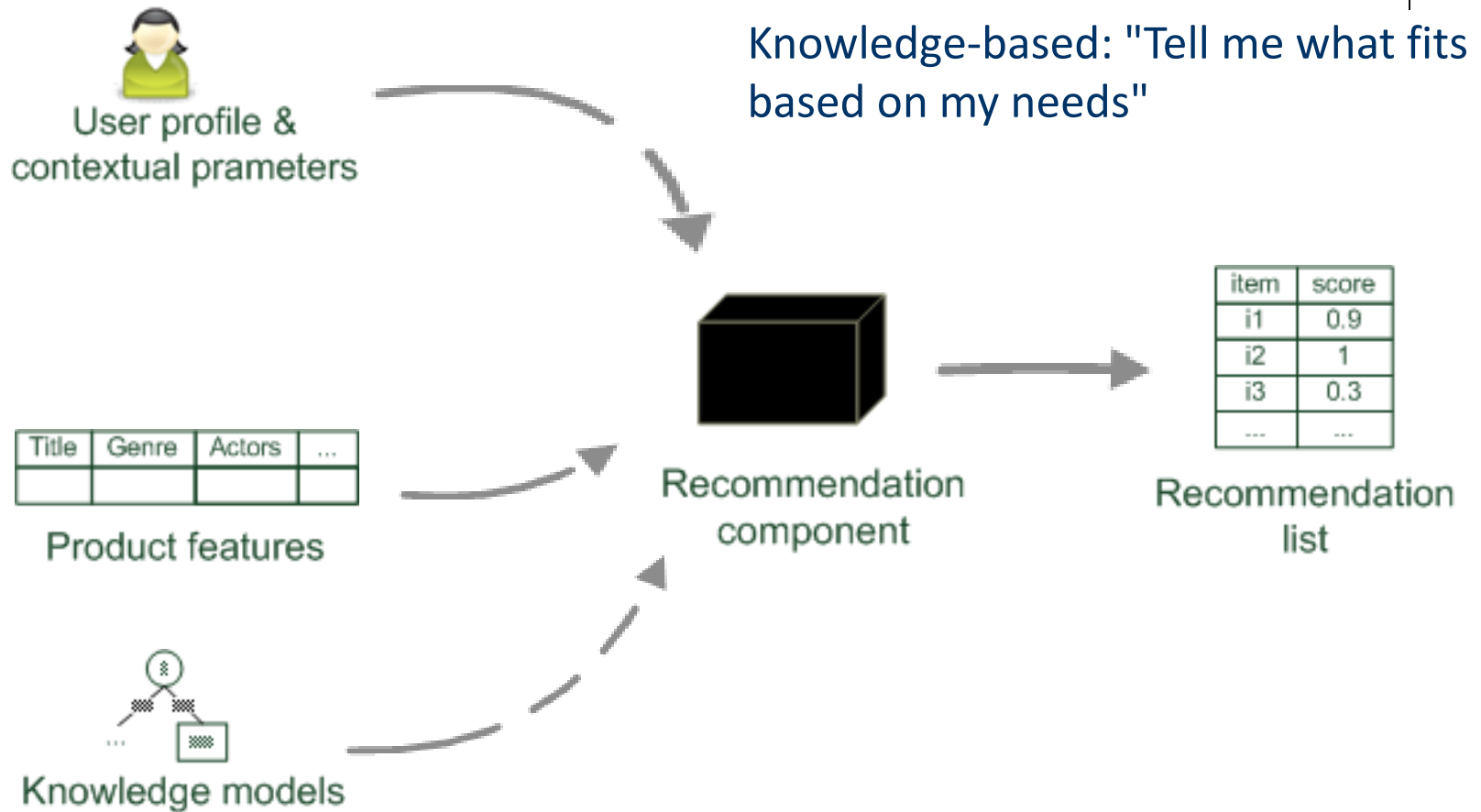
Paradigms of recommender systems



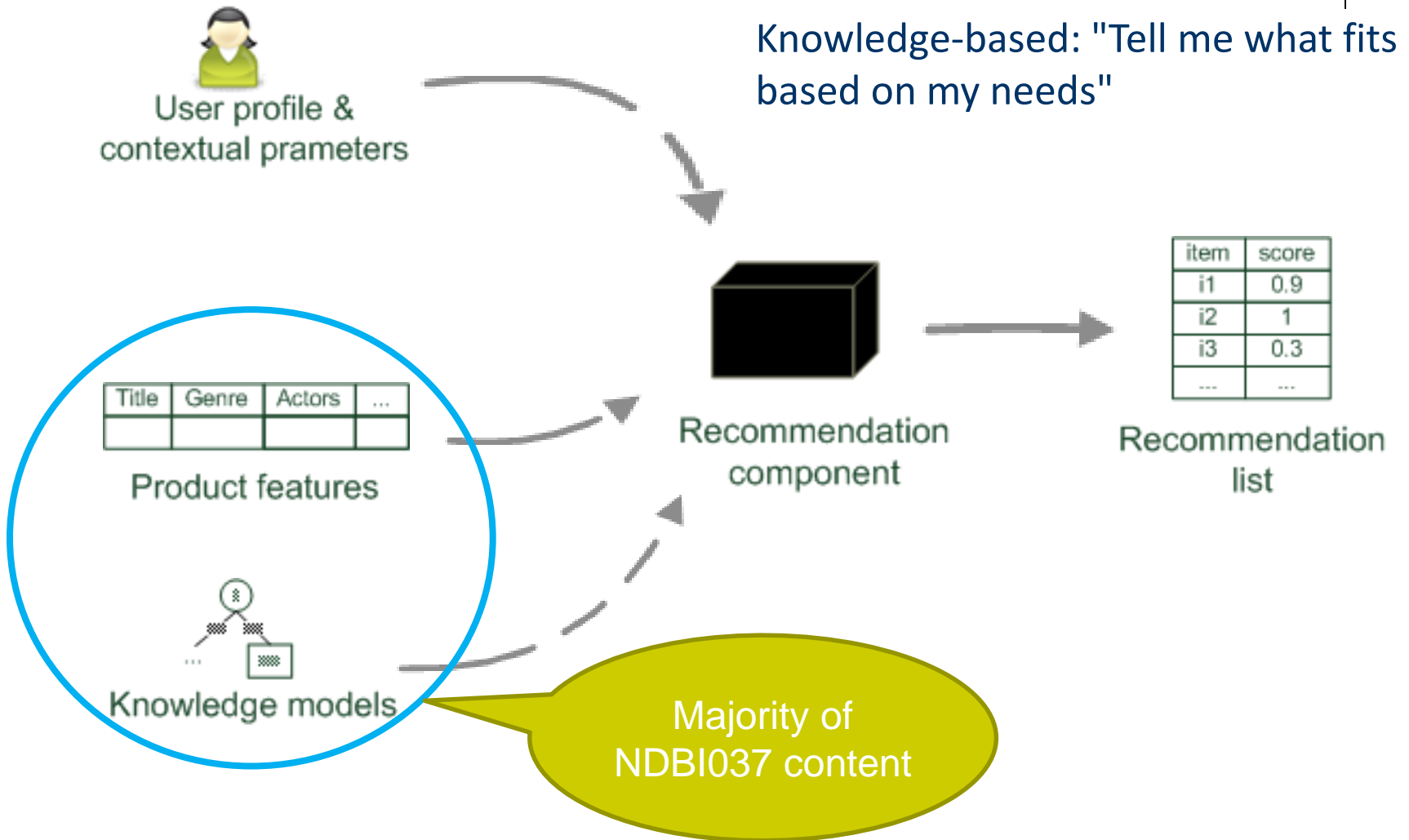
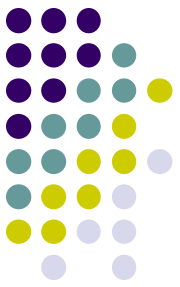
Paradigms of recommender systems



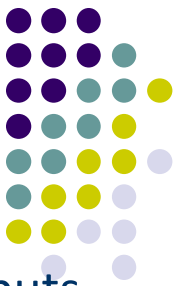
Paradigms of recommender systems



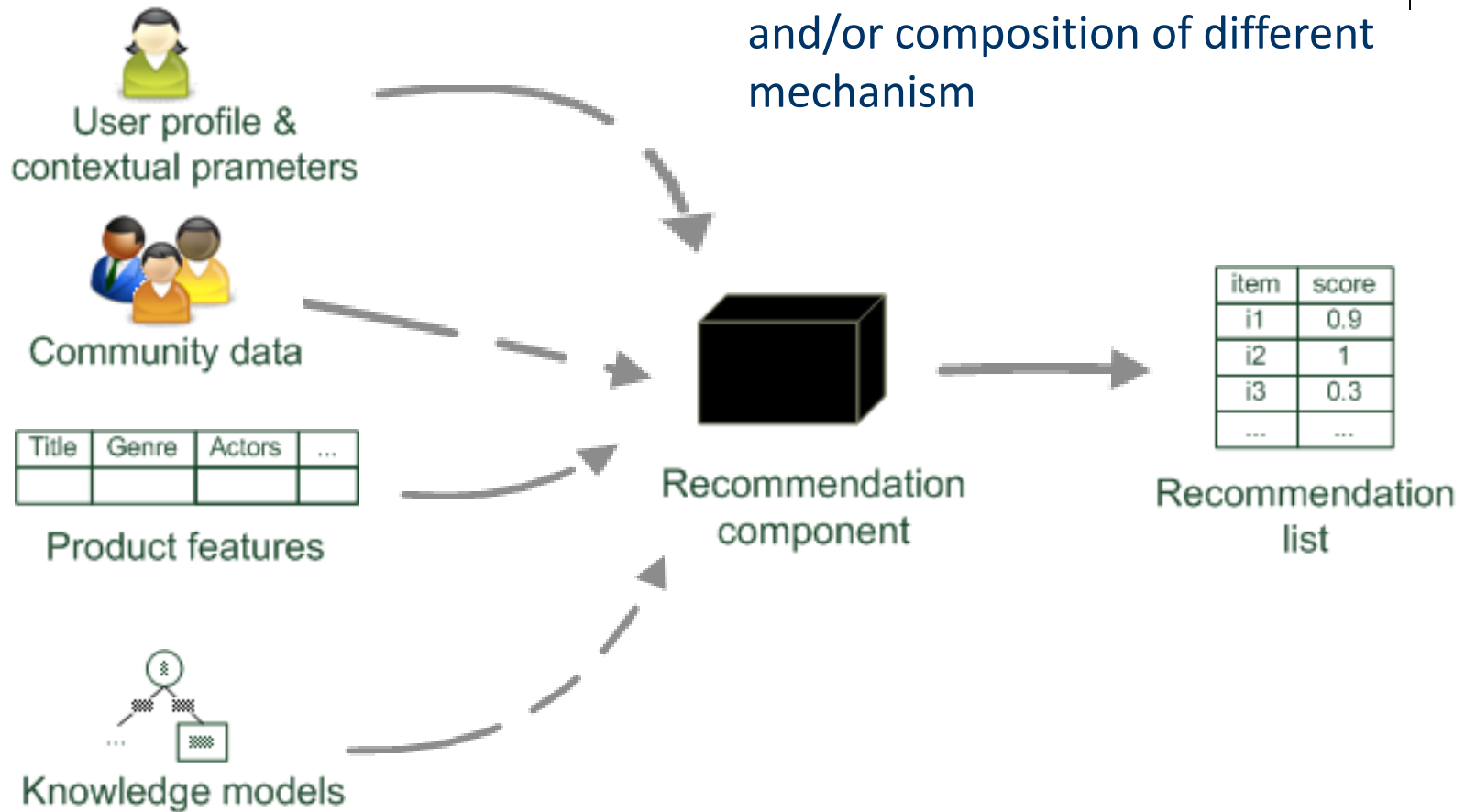
Paradigms of recommender systems



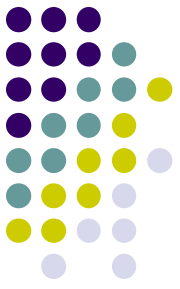
Paradigms of recommender systems



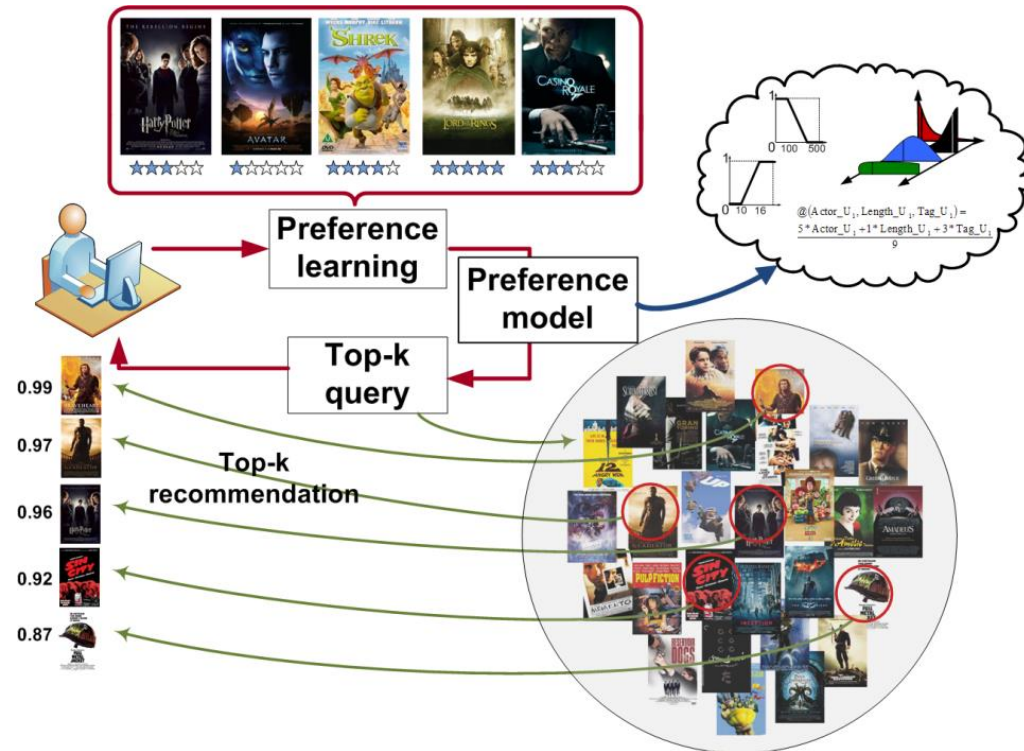
Hybrid: combinations of various inputs and/or composition of different mechanism



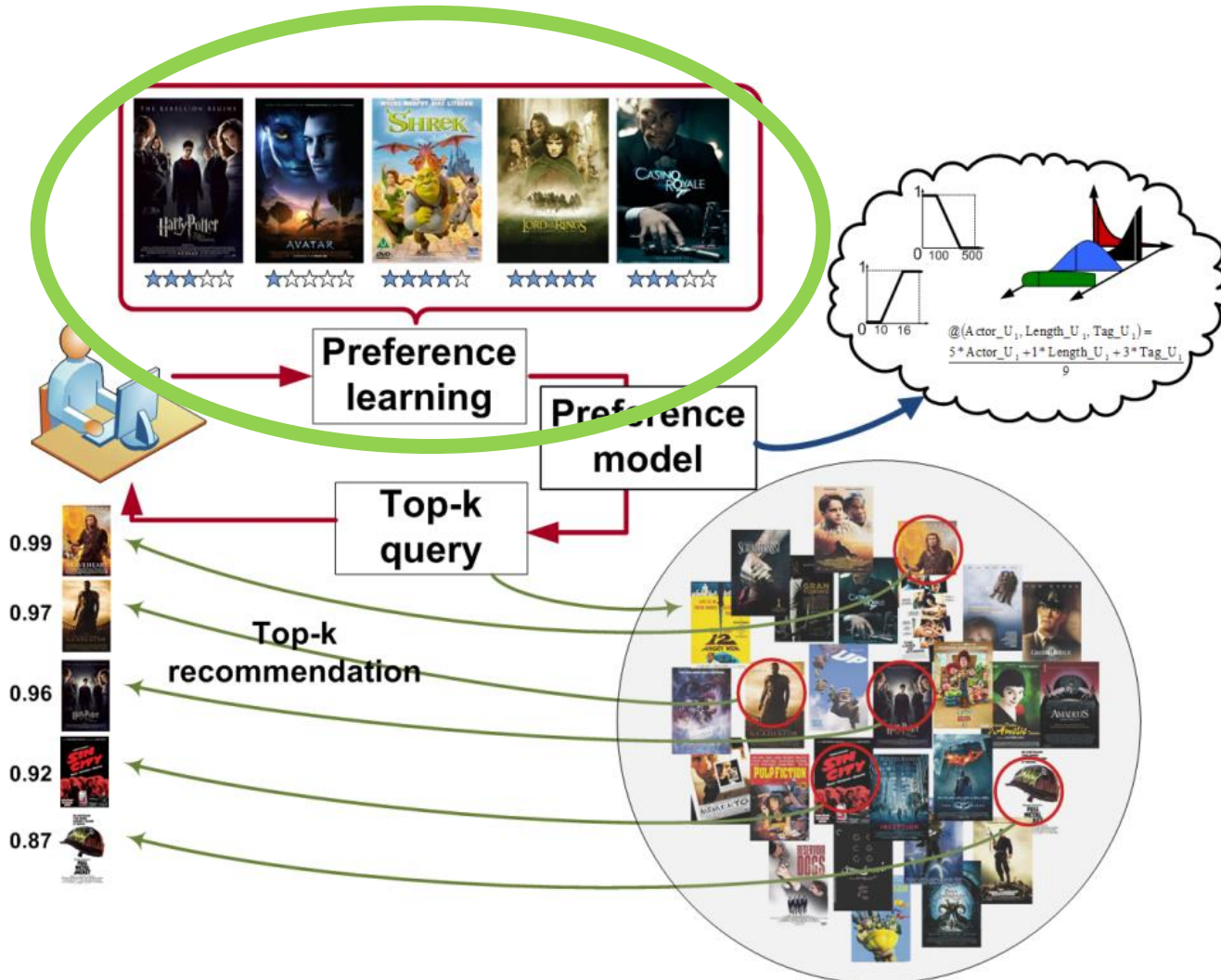
How Does Recommender Systems Work?



1. Get feedback from users (+ additional data)
2. Learn model of user's preferences
3. On demand provide recommendations




User Feedback



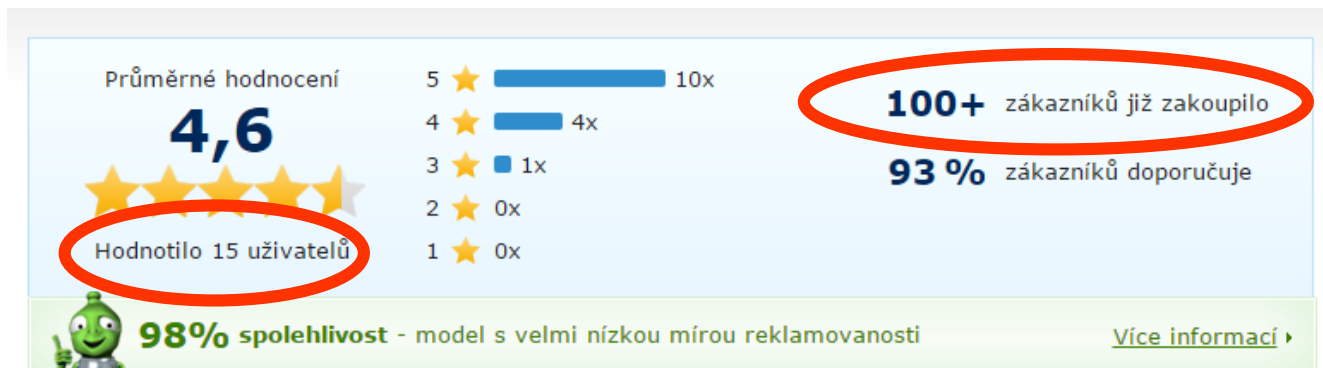
Zpětná vazba



- Explicit feedback (*uživatelské hodnocení*) 
 - Ne od všech uživatelů, ne ke každému objektu, vyžaduje úsilí uživatele
 - Problém v e-commerce (chybí motivace a nutnost nejprve objekt vyzkoušet)
 - + Poměrně přesně popisuje preferenci uživatele
- Implicit feedback (*sledování chování uživatele*)
 - Šum, těžko se interpretuje
 - Není jasné co všechno sledovat
 - + **Můžeme mít feedbacku „kolik chceme“**

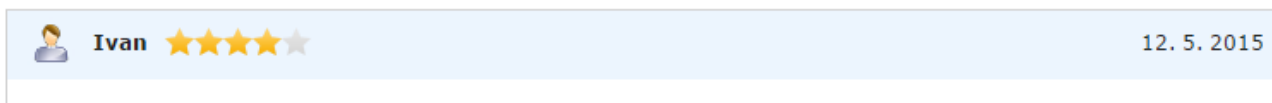


Zpětná vazba



 [Napište uživatelskou recenzi](#)

Uživatelské hodnocení Lenovo IdeaPad G500 Black



Implicit Feedback v E-commerce



Kategorie

- Seznam / výběr objektů
- Parametry kategorie, vyhledávané slovní spojení,...
- Evaluace jednotlivých objektů (visible time, mouseover,...)
- Evaluace kategorie jako celku

Detail objektu

- Počet zobrazení
- Akce myši, scrolování
- Čas na stránce
- Nákup
- (bookmark, tisk, eyetracking, kopírování textu ...)

TOP Nejprodávější Od nejdražšího Od nejlevnějšího 119 položek

HP Officejet Pro 8000 Enterprise	Canon PIXMA MG5350	Epson Stylus SX435W
Inkoustová tiskárna A4, 15str.mono, 14str.color, 600x600dpi, LCD, 256MB, duplex, USB+LAN	Inkoustová tiskárna multifunkční, A4, skener/ kopírka, ESAT 12.5str.mono, 9.5str.color, LCD, 960x2400dpi, duplex, USB+WiFi	Inkoustová tiskárna multifunkční, A4 tiskárna/ skener/ kopírka, 33str.mono, 15str. color, LCD, čtečka karet, 5760x1400dpi, USB 2.0 + WiFi
 -10% 3.461,- 2 325,- s DPH 2 790,-	 -41% 3.880,- 1 908,- s DPH 2 290,-	 -10% 2.719,- 1 249,- s DPH 1 499,-
Koupit	Koupit	Koupit
Skladem > 5 ks	Skladem > 5 ks	Skladem > 5 ks

Below the grid, three more products are visible: Epson Expression Home XP-405WH, Epson Expression Home XP-305, and HP DeskJet 2050A.

HP Officejet Pro 8000 Enterprise

Inkoustová tiskárna A4, 15str.mono, 14str.color, 600x600dpi, LCD, 256MB, duplex, USB+LAN

Cena bez DPH: 2 325,-
Cena s DPH: **2 790,-**
Běžná cena: 3.461,-
Ušetříte: 10% / 311,-

Garantujeme doručení
2* Zbývá u vás nebo ještě dnes na pobočce

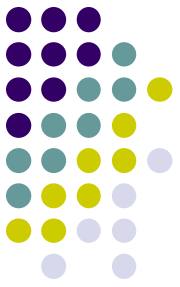
Skladem > 5 ks **Koupit**

Kód: P4222a3
Záruka: 24 měsíců
Prod. číslo: CQ514A#BEN
Odkazy: Stránky výrobce

HP Officejet Pro 8000 Enterprise

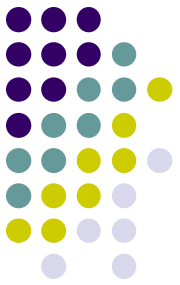
Spolehlivá a tichá inkoustová tiskárna HP Officejet Pro 8000 Enterprise s automatickým sbourávacím systémem jako výhodné tiskové zařízení v kancelářích a se snadnou integrací do stávajícího tiskového prostředí. Tiskárna

Jak ze zpětné vazby získám preferenci?

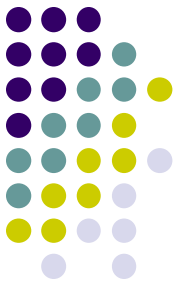


- Explicitní feedback je relativně přímočarý
 - Normalizace hodnocení mezi uživateli
- Implicitní může být větší problém
 - Time on page: 120 sec
 - Number of visits: 2
 - Scrolling distance: 500px
 - Mouse distance: 150px
 - ...

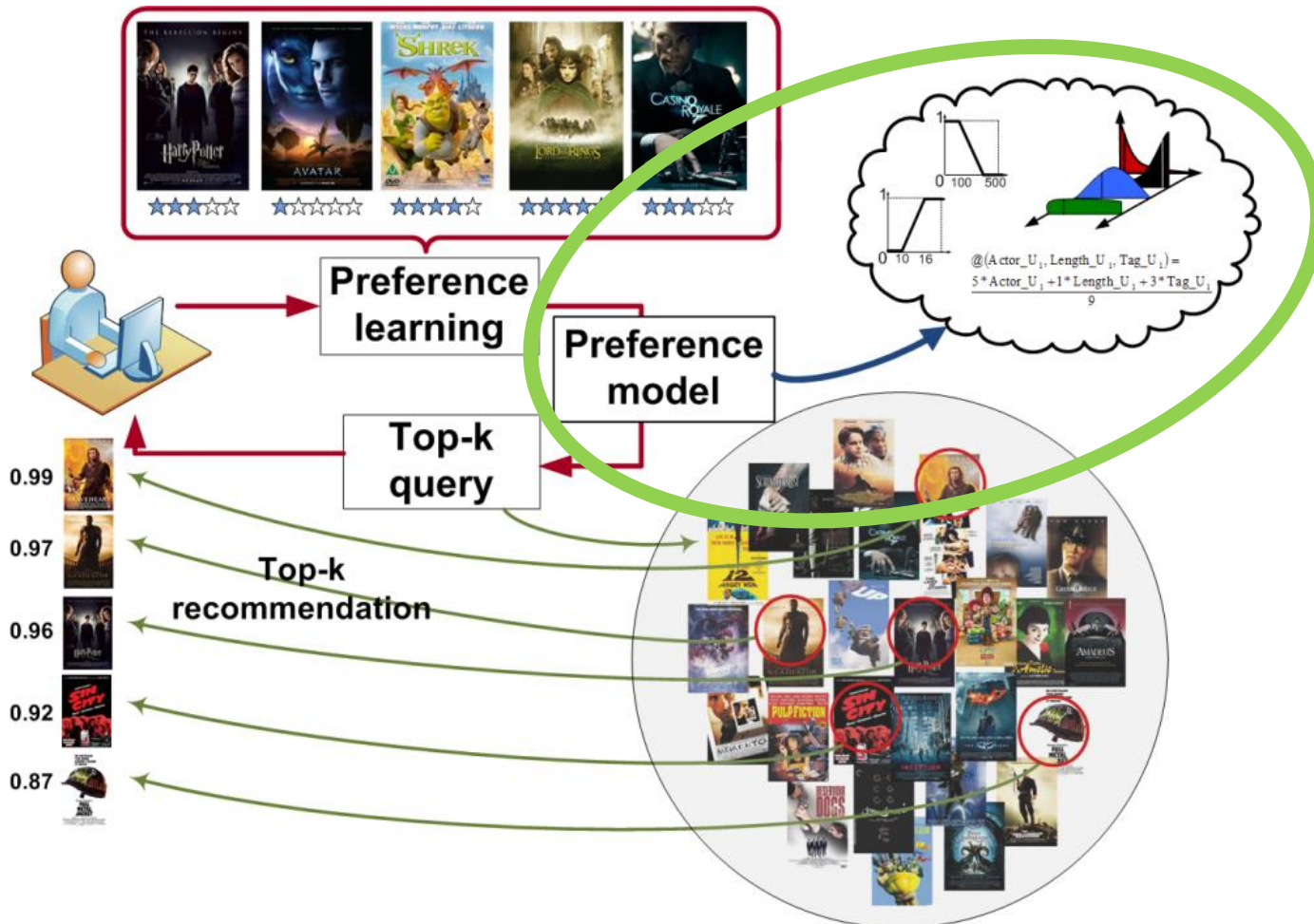




Recommending Algorithms



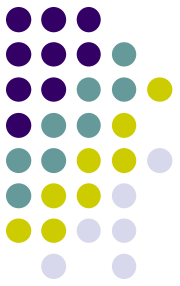
Doporučovací algoritmy



DISCLAIMER



- Whenever there is a **human factor** (client, user) in evaluation
 - ...and this THE case of RecSys*
- You cannot say, what is **right or wrong**
 - You may just evaluate and interpret usage statistics or ask users for their opinion
- *Algorithms materialize some hypothesis about user's behavior*
 - *Sometimes work, sometimes not*
 - *Never cover all aspects of each user*



Recommending Algorithms

No personalization

- „**Most popular**“
- „**People who bought this also bought that**“
- „**Similar objects to this one**“
- ...

Personalized

- Similarity of objects
• „**Content-based**“
- Similarity of behavior
• „**Collaborative filtering**“
- Hybrid, context-based
- Session based...

People were also interested in



NEW Acer ES1-512-C96S
Intel...
7,261.39 CZK

Buy It Now
Free shipping
Popular



Asus Laptop PC
X551MAV-UB01 15...
5,440.90 CZK

Buy It Now
Free shipping



Dell Inspiron 11.6"
Convertible Laptop...
6,777.29 CZK

Buy It Now
Free shipping
Popular

ěška, Doporučovací systémy,
14.5.2015

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**



Collaborative-filtering algorithms



- K-Nearest Neighbors
 - Simplest solution, 90' way of thinking, performance problems, easy implementation
- Matrix factorization (latent factor models)
 - SOTA before deep learning emerged
- Graph based algorithms
 - Social networks
- Market Basket Analysis
- ...

User-based nearest-neighbor collaborative filtering

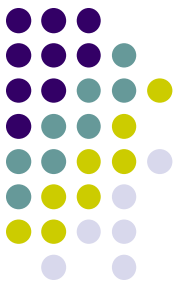


- You can write it after reading 3-sentence description
- Get most similar users to the current one.
- Aggregate their feedback on not-yet-visited items.
- Recommend items with best aggregated feedback.

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**
 - *item*-based CF / **matrix factorizations** are examples for model-based approaches

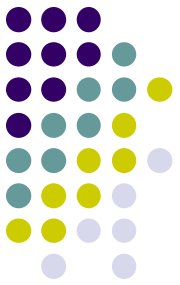


Matrix Factorization

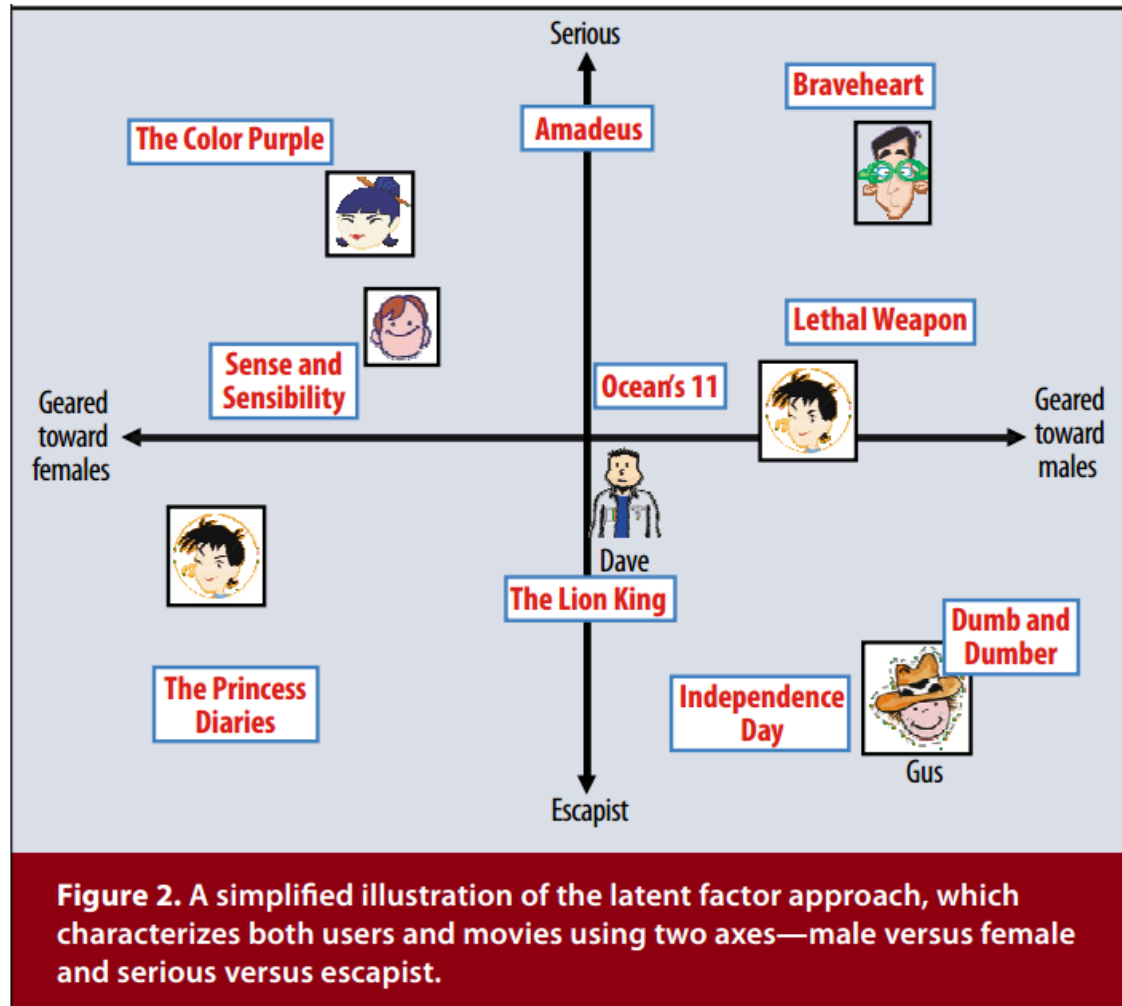
- Matrix of ratings (implicit feedback) of users on items
- Decompose this matrix into latent factors of users and items
- Decomposition should provide similar results on known data

$$\mathbf{R} \approx \mathbf{U}\mathbf{O}^T = \underbrace{\begin{bmatrix} \mu_1^T \\ \mu_2^T \\ \vdots \end{bmatrix}}_{n \times f} \times \underbrace{[\sigma_1 \quad \sigma_2 \quad \dots]}_{f \times m}$$

- Define error of this representation
- Learn to minimize this error
- <http://www2.research.att.com/~volinsky/papers/ieeecomputer.pdf>



Faktorizace Matic

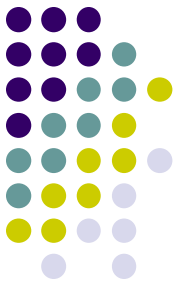




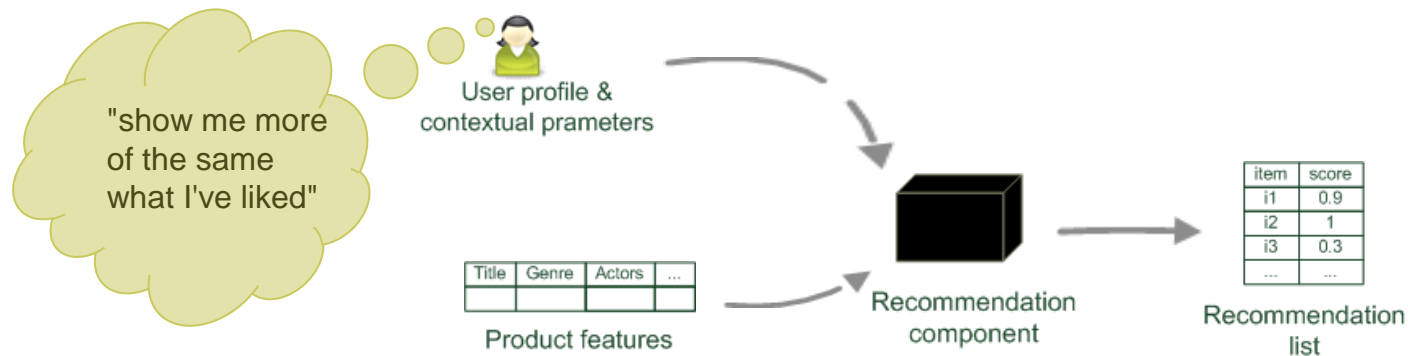
Data sparsity problems

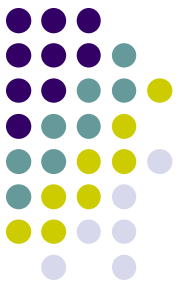
- Cold start problem
 - How to recommend new items? What to recommend to new users?
- Straightforward approaches
 - Ask/force users to rate a set of items
 - Use another method (e.g., content-based, hybrid or simply non-personalized) in the initial phase
 - Default voting: assign default values to items that only one of the two users to be compared has rated (Breese et al. 1998)

Content-based recommendation



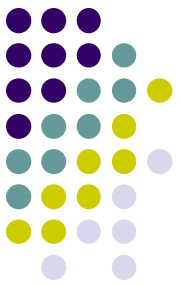
- **While CF – methods do not require any information about the items,**
 - it might be reasonable to exploit such information; and
 - recommend fantasy novels to people who liked fantasy novels in the past
- **What do we need:**
 - some information about the available items such as the genre ("content")
 - 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





Content-based Algorithms

- Vector Space Model
 - Shared space of users (a.k.a. query) and objects, define similarity w.r.t. content-based attributes
 - Often TF-IDF weighting
- Common machine learning
 - Decision trees / forests / GBT
 - Graph-based algorithms (Linked Data)
 - SVM...



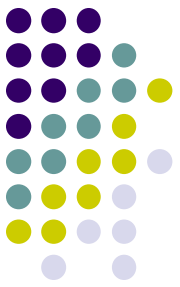
What is the "content"?

- Most CB-recommendation techniques were applied to recommending text documents.
 - Like web pages or newsgroup messages for example.
- Content of items can also be represented as text documents.
 - With textual descriptions of their basic characteristics.
 - Structured: Each item is described by the same set of attributes
Unstructured: free-text description.



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

Content representation and item similarities



Item representation

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

User profile

Title	Genre	Author	Type	Price	Keywords
...	Fiction	Brunonia, Barry, Ken Follett	Paperback	25.65	Detective, murder, New York

$keywords(b_j)$ describes Book b_j with a set of keywords



Simple approach

- Compute the similarity of an unseen item with the user profile based on the keyword overlap (e.g. using the Dice coefficient)



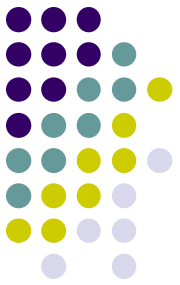
$$\frac{2 \times |keywords(b_i) \cap keywords(b_j)|}{|keywords(b_i)| + |keywords(b_j)|}$$

- Or use and combine multiple metrics

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 II

- Given a keyword i and a document j
- $TF(i, j)$
 - term frequency of keyword i in document j
- $IDF(i)$
 - inverse document frequency calculated as $IDF(i) = \log \frac{N}{n(i)}$
 - N : number of all recommendable documents
 - $n(i)$: number of documents from N in which keyword i appears
- $TF - IDF$
 - is calculated as: $TF-IDF(i, j) = TF(i, j) * IDF(i)$



Cosine similarity

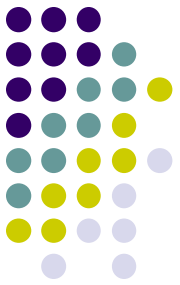
- Usual similarity metric to compare vectors: Cosine similarity (angle)
 - Cosine similarity is calculated based on the angle between the vectors

- $sim(\vec{a}, \vec{b}) = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| * |\vec{b}|}$

- Adjusted cosine similarity

- take average user ratings into account (\bar{r}_u), transform the original ratings
 - U: set of users who have rated both items a and b

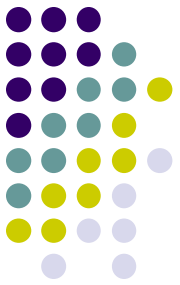
- $sim(\vec{a}, \vec{b}) = \frac{\sum_{u \in U} (r_{u,a} - \bar{r}_u)(r_{u,b} - \bar{r}_u)}{\sqrt{\sum_{u \in U} (r_{u,a} - \bar{r}_u)^2} \sqrt{\sum_{u \in U} (r_{u,b} - \bar{r}_u)^2}}$



How to know what to use?

- Experiment as much as you can
 - *If You want to double your success rate, you should double your failure rate.*

Experiments



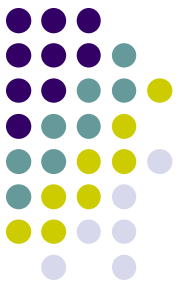
Online

- Production servers
- Hard to repeat
- Expensive (time + money)
- Selected methods
- Real metrics
- GUI changes etc. can be also evaluated

Offline

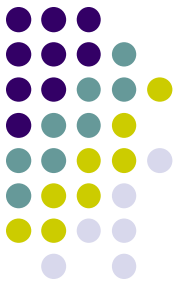
- Data-based simulation
- Easily repeatable
- Fast and cheap
- Artificial metrics (RMSE, MAE, diversity...)
- Causality problems; Only algorithms can be compared

Success in offline do not imply success in online...
...however the lack of success usually holds.



Conclusions

- Recommender systems slowly became standard in web applications.
- There is always problem with insufficient data
 - Tradeoff between complexity and train-ability
 - Multiple pathways to explore
 - Domain dependent best practices
- Basic algorithms are easy to handle (devil in details)
- Important research topic (Bc, Mgr, PhD thesis etc.)



Co jsem vynechal:

- Deep learning
- Jak zobrazovat doporučení
- Jak vysvětlovat doporučení (user trust)
- Vývoj preferencí uživatele v čase, dlouhodobé / krátkodobé preference
- Kontext (místo, čas, nabídka,...)

Dotazy?

