### Collaborative Computing

### Information Filtering (IF)

- Today's lecture:
  - Information Filtering (and Information Retrieval)
  - o Collaborative Filtering and Content based filtering
  - Filtering Techniques
  - o User Modelling and levels of adaptation
  - Adaptive Systems in general
  - Basic Evaluation of Information Filtering Systems
  - o Research Trends

Andrea Andrenucci

- The advent of the Web has exposed users to a huge amount of information.
- Information Filtering is a technique that tries to reduce the information overload and filter information that is relevant to users.
- Information is filtered with the help of a profile of the user, also called User Model (UM) or User Profile (UP)

Andrea Andrenucci

### Differences between Information Filtering and Information Retrieval

- Old Definitions:
  - Information Retrieval is concerned with retrieving information to a user on the basis of user questions/queries
  - Information Filtering is concerned with building a long term profile of the user information needs and sort out incoming information to the user
  - (Belkin and Croft 1992)
- Techniques in IF are similar to the techniques utilized in IR.
- Today: Personalized Information Retrieval is also considered in the domain of IF (Waern 04)

Andrea Andrenucci

# Where is Information Filtering utilized? Recommender Systems, i.e. systems that make a gersonalized selection of information items or products. Bows filtering, i.e. systems filter incoming streams of news information The information items or products. Bows filtering (e.g. filters out SPAM)

# Information Filtering main categorization

- Balabanovic and Shoham (97) categorize IF into two major topics: Content Based Filtering and Collaborative Filtering
- Content Based Filtering: representations of the information items are compared to the representation of the user (user model) in order to find the information items that are relevant.
- Collaborative Filtering aims at predicting user preferences, based on the preferences of a group of users (with similar interests).
- Information can be gathered from both restricted domains or open domains.

### Other Categorizations of IF systems

- Initiative of Operation: Active vs Passive
- Location of Filtering Operation: at the information source, at a filtering server (3-tier architecture), or at the user's site (locally)
- Methods for acquiring information about the user: explicit, implicit or a combination of both (more about this later)

5

6

### More on Collaborative Filtering

- The focus is on the opinions of user groups rather than on the content of a document or item.
- Can be defined as a Social Navigation Technology (Munro, Hook, Benyon 99).
- Social Navigation: human beings are social animals and tend to follow other people's advice or judgment when looking for information or buying items.

### Andrea Andrenucci

### Techniques utilized in Information Filtering

- Filtering techniques are usually divided in:
- Knowledge-based techniques: e.g. rules and semantic nets
- Statistical techniques: data based (e.g. user profiles are weighted vectors of terms that are compared to weighted vector of terms of information items or other user profiles)

Andrea Andrenucci



































### Content-based filtering research example – Persival (McKeown et al., 2001)

- Personalized Search Engine for HealthCare articles
- Extract user profiles from patient records and utilizes the information in the UP to rank the documents retrieved from online medical resources
- Represent user profiles and documents as term-value vectors and utilize cosine similarity
- Utilize Natural Language Processing techniques (syntactic parsing) to parse medical documents and create summaries that are tailored to the patients background
- Helps practitioners to find evidence for treatment or diagnosis of patient diseases

Andrea Andrenu

### Content-based Filtering - Drawbacks Collaborative Filtering - Drawbacks Require machine-readable/parsable items, e.g. Requires bootstrapping: recommendations cannot be text-based documents, since it creates a formal done if there is not sufficient amount of data, i.e. user representation of the content of the information items ratings · It is more difficult to automatically create a Sparsity problem: users may rate small sets of all representation of images, speech or sound available items or different sets of items, which make comparison of user preferences more difficult. News filtering poses real-time constraints (the system) cannot process the information too long) Early rater problem: prediction cannot be made for an ٠ item when it first appears, since there are not user It is difficult to judge the quality of the information items ratings for that item. Changing interestes problem: what happens if our interests change? Do we have to re-rate all the items?

# How to overcome those problems? Combine! Claypool et al. (99) combines both approaches in a system called P-Tango, that filters news articles. The prediction of articles relevance is based on the average of the content-based predictions and the collaborative predictions. Grouplens (Sarwar, Konstan et al. 98) also combines the approaches. The system provides a content-based evaluation of news articles and computer ratings are treated just like ratings of human beings. ProfBuilder (Wasfi, 99) recommends Web pages in two lists: one generated by content-based Filt. and another generated by Collaborative Filtering.



Andrea Andrenucci

Andrea Andrenucci

2







"Content adaptation" Example Text for doctors and patients: Opade system from (De Carolis et al. 96)

Patient	Doctor
Comments to the drug prescription of Mr Fictif.	As you certainly remember, Mr Fictif is a 62 years old man. He is overweight
You have been diagnosed as suffering from a mild of what we call "angina pectoris", that is a spasm of chest resulting from overexertion when heart is diseased. In addition you have elevated cholesterol	He is suffering from a mild form of angina and he has got elevated cholesterol











## 

# Exampels that combine both explicit and implicit UM

- P-Tango (Claypool et al.) recommends news articles for an online newspaper. The UM both *explicit*, with keywords entered by the user, and *implicit*, with keywords gathered from articles the users rated as interesting.
- ConfCall (Waern et al. 04) recommends relevant Conference calls. Through a profile editor, users submit keywords about their interests. The system then monitors which incoming documents users read or discard, updating the profile.

<section-header><text><image>

















