

# 3-4: Content-Based Recommenders in Detail

# Learning Objectives

- To how to build content-based recommenders based on TFIDF concepts, including:
  - Computing vectors to describe items
  - Building profiles of user preference
  - Predicting user interest in items
- To understand key variants in implementing CBF recommenders, and their strengths and weaknesses

# Key concept: Keyword Vector

- The universe of possible keywords defines a content space
  - Each keyword is a dimension
  - Each item has a position in that space; that position defines a vector
  - Each user has a taste profile (or more than one) that is also a vector in that space
  - The match between user preference and items is measured by how closely the two vectors align
  - May want to limit/collapse keyword space (e.g., stem and stop)

# Useful reference

- Vector Space Model (originally conceived for queries, indexing)
  - [http://en.wikipedia.org/wiki/Vector\\_space\\_model](http://en.wikipedia.org/wiki/Vector_space_model)
  - Salton, Wong, and Yang (1995) “A Vector Space Model for Automatic Indexing,” *CACM* 18:11.

# Where choices come into play ...

- Representing an item through a keyword vector:
  - Simple 0/1 (keyword applies or doesn't)
  - Simple occurrence count
  - TFIDF, most commonly:
    - $\# \text{ occurrences in doc} * \log (\# \text{ docs} / \# \text{ docs with term})$
  - Other variants that include factors such as document length
  - Eventually, this vector is often normalized

# Consider the movie/tag case ...

- Do we consider tags to be yes-or-no?
  - Actor (we don't really get a measure for how much "Tom Hanks" a movie has)
  - Descriptive (is how often a tag is applied a proxy for how relevant/significant the feature is?)
- Do we care about IDF?
  - Actor (are infrequent actors more significant than stars?)
  - Descriptive (is "prison scene" more significant than "car chase" or "romance"?)

# From Items to User Profiles (1)

- Vector Space model conflates liking with importance
  - Works well in some applications (query terms), not as much in others (I like Hollandaise sauce a lot, but don't actually care much if it is in a dish I'm ordering)
  - Consider how this may play out with movie keywords ...

# ...to User Profiles (2)

- How do we accumulate profiles?
  - Add together the item vectors?
    - Do we normalize first?
      - Do we believe an item with a longer vector is more descriptive of preferences? Then maybe no.
      - Do we think all items should be the same weight? Then maybe yes.
    - Do we weigh the vectors somehow?
      - Could use ratings for weight ...
      - Could use currency, confidence ...



# ...to User Profiles (3)

- How do we factor in ratings?
  - Simply unary – aggregate profiles of items we rated without weights
  - Unary with threshold – only put items above a certain rating into our profile (but all likes are equal) – we often think 3.5 in MovieLens
  - Weight, but positive only – higher weight for things with higher scores
  - Weight, and include negative also – negative weight for low ratings (normalize rating scale)

# ...to User Profiles (3)

- How do we update profiles (new ratings)?
  - Don't – just recompute each time (wasteful)
  - Weight new/old similarly – keep track of total weight in profile and mix in new rating (linear combination)
    - Special case for changed rating; subtract old
  - Decay old profile and mix in new (e.g., may decide that profile should be dominated by most recent 50 movies – formula might be  $0.95 * \text{old} + 0.05 * \text{new}$ , in normalized form)

# Computing Predictions ...

- Prediction is the cosine of the angle between the two vectors (profile, item)
- This is the dot-product of normalized vectors, or if you prefer, the dot product divided by the product of the two lengths
- Cosine ranges between -1 and 1 (0 and 1 if all positive values in vectors) – closer to 1 is better.
- Top-n, or scale for rating-scale predictions.

# Strengths of this Approach

- Entirely content-based
- Understandable profile
- Easy computation
- Flexibility – can integrate with query-based systems, case-based approaches

# Challenges

- Figuring out the right weights and factors
  - Is more more, or just reiteration of the same
  - How to deal with ratings

# Limitations

- This is a highly simplified model, cannot handle interdependencies
  - I like Sandra Bullock in Action movies, but Meg Ryan in Romantic Comedy movies
  - I like comedies with violence, and historical documentaries, but not historical comedies or violent documentaries

# Take-aways

- Content-based filtering based on assessing the content profile of each item; can be done from metadata or user tagging
- User profiles built by aggregating profiles of items rated/consumed, possibly with a weighting scheme
- Evaluate unrated items by matching item profile against user profile: vector cosine

# Moving Forward

- Next Lecture
  - Survey of content filtering tools
- Assignments
  - Written: Using content filtering; design exercise
  - Programming: Adding a basic content filtering module to LensKit
- Looking further forward
  - After collaborative filtering, some of these concepts return when we look at SVD/LSI



# 3-4: Content-Based Recommenders in Detail