

MAPR[®] Academy

Introduction to Machine Learning

Introduction to Recommendation

Course Name

© 2014 MapR Technologies **MAPR** 1





Introduction to Recommendation

- ▶ Define recommendation terms and concepts
- ▶ Describe the K-nearest neighbors algorithm
- ▶ Evaluate a recommendation model
- ▶ Discuss the Mahout implementations for recommendation



Course Name

© 2014 MapR Technologies  2





Recommender examples in the wild

The collage consists of three separate screenshots:

- Top right:** A map showing a residential area with several points of interest labeled A through G. Points D, B, C, and G are highlighted with red circles.
- Middle left:** A Netflix interface. A red circle highlights the text "Because you watched American Crime Story: Part One". Below it, another red circle highlights the section "Customers Who Bought This Item Also Bought".
- Bottom center:** A product listing for "Getting Started with D3" by Mike Dewar. It shows a 4-star rating from 18 reviews. To its right are two other book covers: "Data Visualization with D3.js Cookbook" by Nick Qi Zhu and "Data Points: Visualization That ..." by Nathan Yau.

ologies MAPR. 3





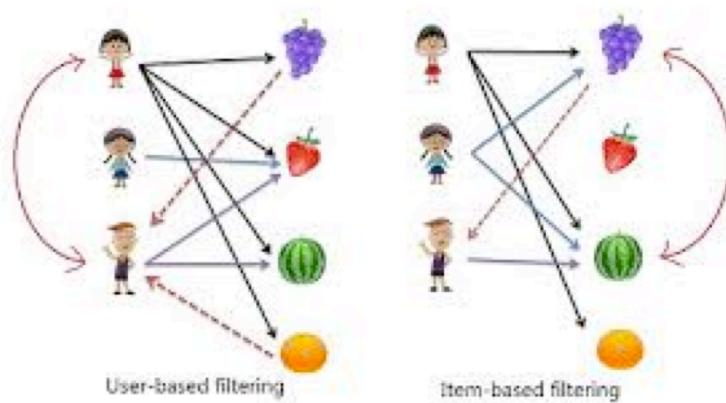
What is recommendation?

- **Recommendation** is a class of ML that seeks to predict a user's *preference* for or *rating* of an item
- **Recommender systems are used in industry to recommend:**
 - Books and other products (e.g Amazon)
 - Music (e.g. Pandora)
 - Movies (e.g. Netflix)
 - Restaurants (e.g. Yelp)
 - Jobs (e.g. LinkedIn)
 - ... LOTS ...
- **Main approaches to recommendation**
 - Collaborative filtering
 - Content-based filtering





User-based vs item-based filtering



Course Name

© 2014 MapR Technologies **MAPR** 5

Users are associated with item ***preferences*** through history

Similarity is constructed between users (or items)

Recommendations based on ***similarity to other users (or items)***





Similarity metrics

- **Pearson correlation**

- Ratio of covariance to product of standard deviations
- $-1 \rightarrow$ inversely proportional; $0 \rightarrow$ no correlation; $1 \rightarrow$ directly proportional

- **Euclidean distance**

- Coordinates represent item preferences (i_1, i_2, i_3, \dots)
- Smaller distance \rightarrow more similar (so return $1/(1+d)$)

- **Tanimoto coefficient**

- Ratio of intersection to union
- Between 0 and 1 \rightarrow bigger is more similar





Differences between collaborative and content-based filtering

- **Collaborative-based** filtering is *deductive* (needs user-item history to start with to learn associations)
- → example: last.fm
- **Content-based** filtering is *inductive* (needs domain knowledge to construct associations between users and items)
- → example: Pandora





Challenges with collaborative filtering

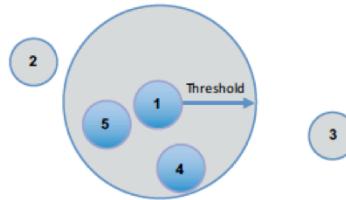
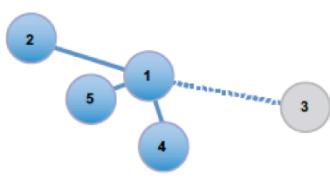
Challenge	Description
Cold start	No user history means no associations day 1
Scale	Huge number of products and users requires a lot of computation
Sparsity	Most users express very little behavior with very few items and no behavior on the vast majority of items





User neighborhood

- A neighborhood (of users) is a group of similar users
- 2 types of neighborhoods:
 - **Fixed-size**: cardinality of neighborhood is fixed a priori (e.g. top ten)
 - **Threshold-based**: degree of similarity is fixed a priori (e.g. $\geq 90\%$)



Q: why no *item* neighborhood in collaborative filtering?





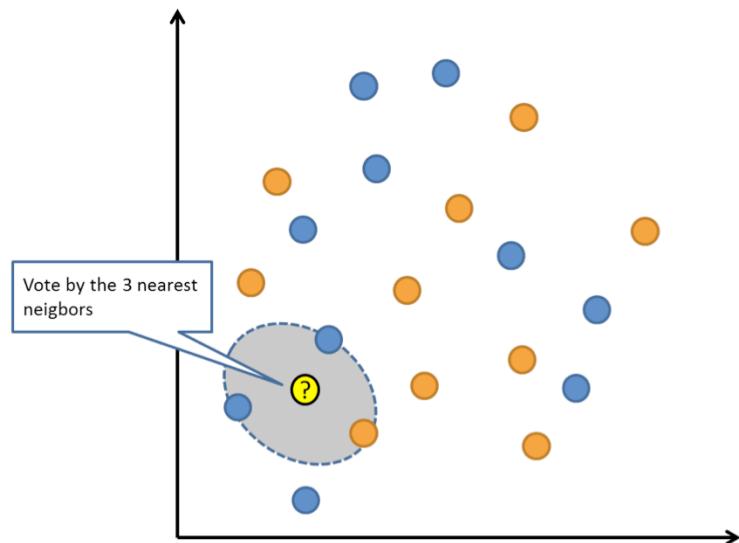
Introduction to Recommendation

- ▶ Define recommendation terms and concepts
- ▶ Describe the K-nearest neighbors algorithm
- ▶ Evaluate a recommender model
- ▶ Discuss the Mahout implementations for recommendation





K-Nearest Neighbor (KNN)



Course Name

© 2014 MapR Technologies **MAPR** 11



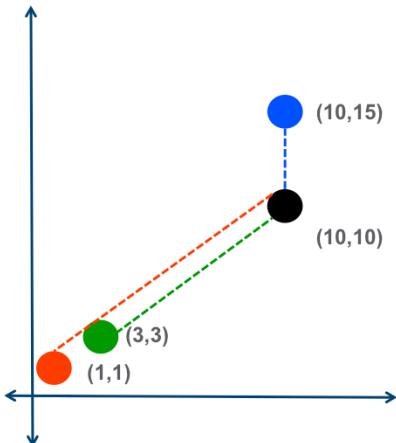
KNN algorithm

- **Distance**
 - Euclidean
 - Pearson
 - Spearman
- **Majority voting**
 - Problematic if class distribution is skewed
 - Could use distance as weight to compensate
 - Pick odd number for k to avoid ties





Demo: MyKNearestNeighbor.java



```
java MyKNearestNeighbor 1 10 10 euclidean
input: [10 10]
[1 1]
distance is 12.727922061357855
[3 3]
distance is 9.899494936611665
[10 15]
distance is 5.0
k nearest neighbors are:
[10 15]
```

```
java MyKNearestNeighbor 1 10 10 cosine
input: [10 10]
[1 1]
distance is 0.9999999999999998
[3 3]
distance is 1.0
[10 15]
distance is 0.9805806756909201
k nearest neighbors are:
[3 3]
```

Course Name

© 2014 MapR Technologies MAPR 13





Introduction to Recommendation

- ▶ Define recommendation terms and concepts
- ▶ Describe the K-nearest neighbors algorithm
- ▶ Evaluate a recommender model
- ▶ Discuss the Mahout implementations for recommendation



Course Name

© 2014 MapR Technologies  14





A search engine is actually a recommendation engine

Google

[Web](#) [Images](#) [Videos](#) [Shopping](#) [News](#) [More](#) [Search tools](#)

About 527,000 results (0.28 seconds)

Scholarly articles for recommender system algorithms

[... of dimensionality reduction in recommender system-a ... - Sarwar](#) - Cited by 1045
[... next generation of recommender systems: A survey of ... - Adomavicius](#) - Cited by 5083
[Introduction to recommender systems: Algorithms and ... - Konstan](#) - Cited by 108

Recommender system - Wikipedia, the free encyclopedia
[en.wikipedia.org/wiki/Recommender_system](#) - Wikipedia
Recommender systems are a useful alternative to search **algorithms** since they help users discover items they might not have found by themselves. Interestingly ...
Collaborative filtering - Information filtering system - Cold start

[\[PDF\] Recommendation Systems - The Stanford University Infol...](#)
[infolab.stanford.edu/~ullman/mmds/ch9.pdf](#) - client, and there are some new algorithms that have proven effective for **algorithm** that could beat its own **recommendation system** by 10%.1. The prize was ...

Recommender Systems :: Algorithms - Computer Science
[www.cs.carleton.edu/cs_comps/0607.../recommender/algorithms.html](#) - Algorithms. We experimented with a number of different types of **algorithms** to build **recommender systems**. To learn more about them, please click on a link ...

Course Name © 2014 MapR Technologies  15

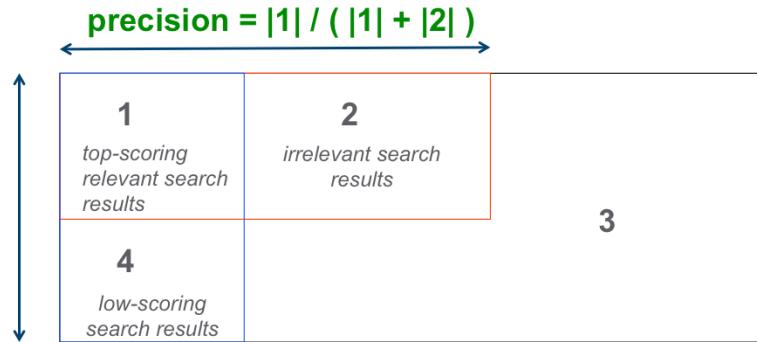




Precision and recall in document corpus search

All docs in corpus	= 1 + 2 + 3 + 4
All results from search	= 1 + 2
Top results from search	= 1
Relevant documents	= 1 + 4

$$\text{Recall} = |1| / (|1| + |4|)$$



Precision and recall are used to ***evaluate search results***

Precision = proportion of top-scoring results that are relevant

Recall = proportion of relevant results that are top-scoring





Precision and recall in collaborative filtering is challenging

- In collaborative filtering, a user history is crucial
 - So the rec evaluator will pick an item from the user history (which isn't necessarily the best rec)
 - Conversely, the rec evaluator should pick an item not in the user history (which would be penalized in the test)
- In boolean recommenders, there's no rating





What to tune for better performance?

- **Model type**
 - User-based
 - Item-based
 - Content-based
- **Distance metric**
 - Euclidean
 - Tanimoto
 - Loglikelihood
- **Neighborhood size**
 - 2
 - 3
 - 200



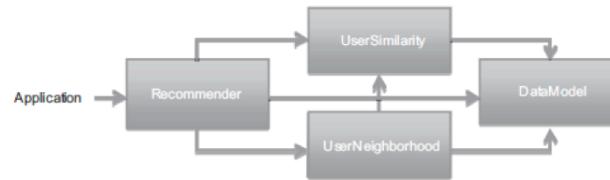


Introduction to Recommendation

- ▶ Define recommendation terms and concepts
- ▶ Describe the K-nearest neighbors algorithm
- ▶ Evaluate a recommender model
- ▶ Discuss the Mahout implementations for recommendation 



Mahout Recommender Architecture



- **Recommender**
 - Uses data and algorithms to make recommendation
- **UserSimilarity**
 - Configurable measure of how similar users are to each other
- **UserNeighborhood**
 - Configurable calculation of set of users who are similar to each other
- **DataModel**
 - Provides interface to store and retrieve users, items, and preferences



User-item preferences

```
1,101,5.0  
1,102,3.0  
1,103,2.5  
2,101,2.0  
2,102,2.5  
2,103,5.0  
2,104,2.0  
3,101,2.5  
3,104,4.0  
3,105,4.5  
3,107,5.0  
4,101,5.0  
4,103,3.0
```

- can use TSV or JDBC
- can use boolean preference

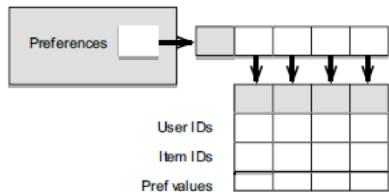
user 4 has preference 3.0 for item 103





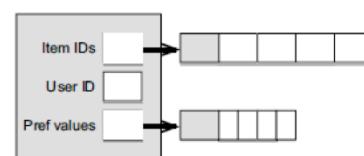
How preferences are stored in memory

Array of Preference objects



inefficient

PreferenceArray



efficient (4x smaller)





Mahout recommender scale considerations

- Mahout recommenders are primarily *memory-bound*
- Noisy data can be *pruned* to reduce memory requirement
- Data may be *sampled* to reduce memory requirement
- Use *database backend* to scale up







- What are the two types of collaborative filtering?
- Which one is deductive: collaborative or content-based filtering?
- How are user preferences for items represented to Mahout?



