

Master on Foundations of Data Science



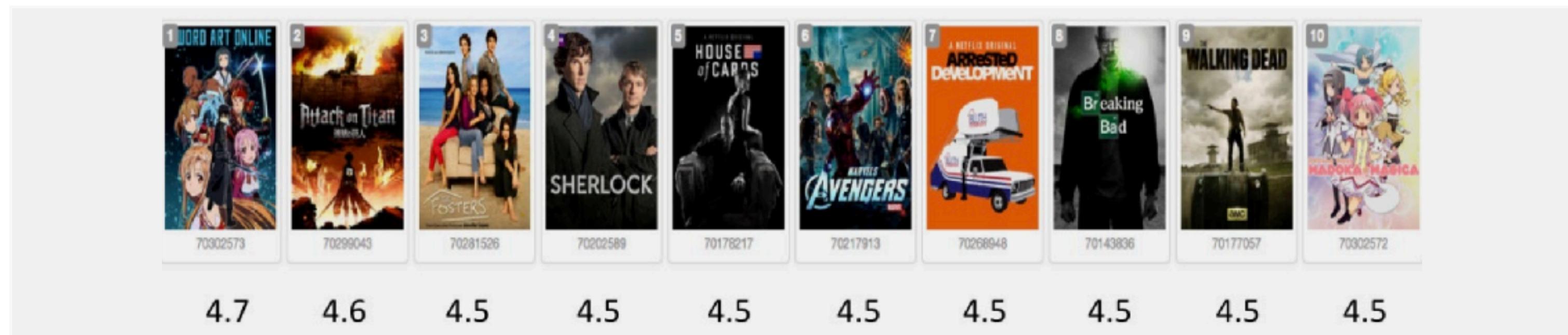
Recommender Systems

Learning to Rank

Santi Seguí | 2022-2023

Ranking

- Most of the recommendations are **presented in a sorted list**.



- Recommendation can be understood as a **ranking problem**.

Optimized to achieve the best

Accuracy

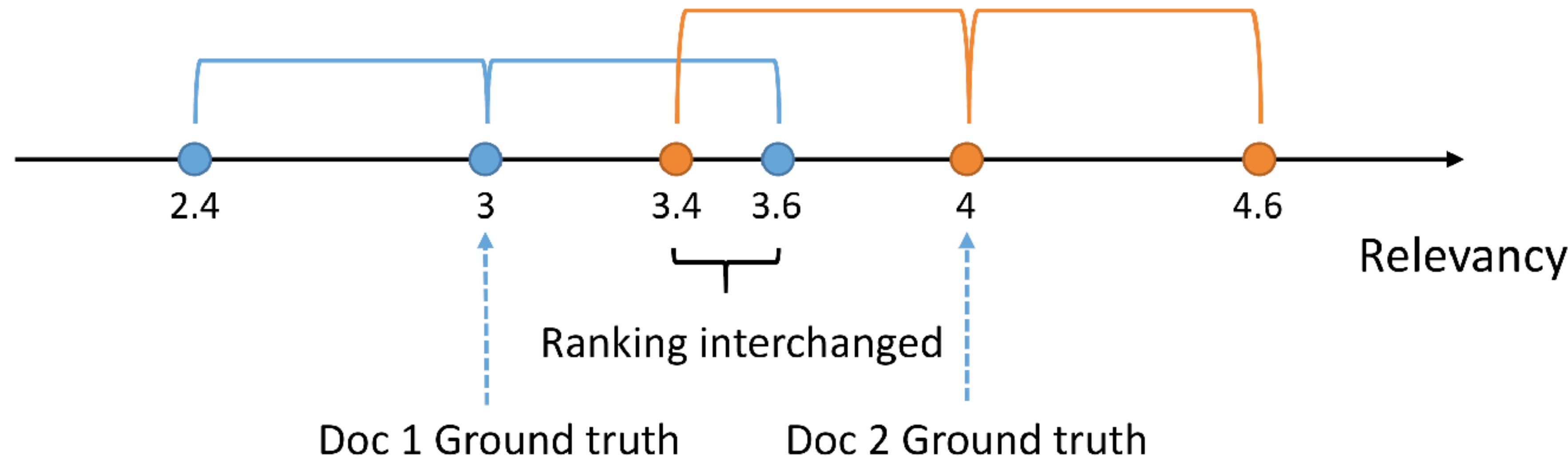
RMSE, MAE...

Accuracy



Any **problem** with RMSE or MAE?

Rating Prediction vs. Ranking



Same RMSE/MAE might lead to different rankings

Accuracy



Any **problem** with RMSE or MAE?

these metrics do **not** really **measure the user experience**

Top Picks for Joshua



Trending Now



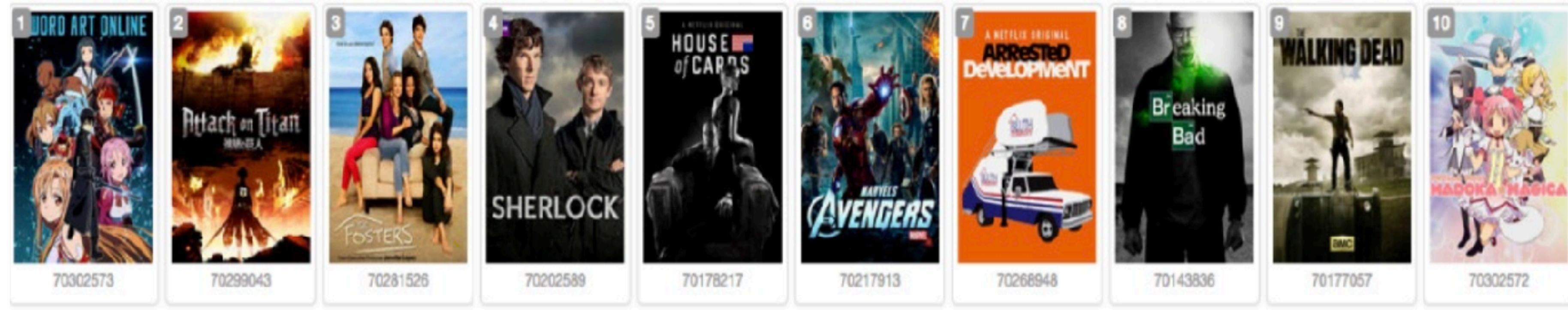
Because you watched Narcos



New Releases

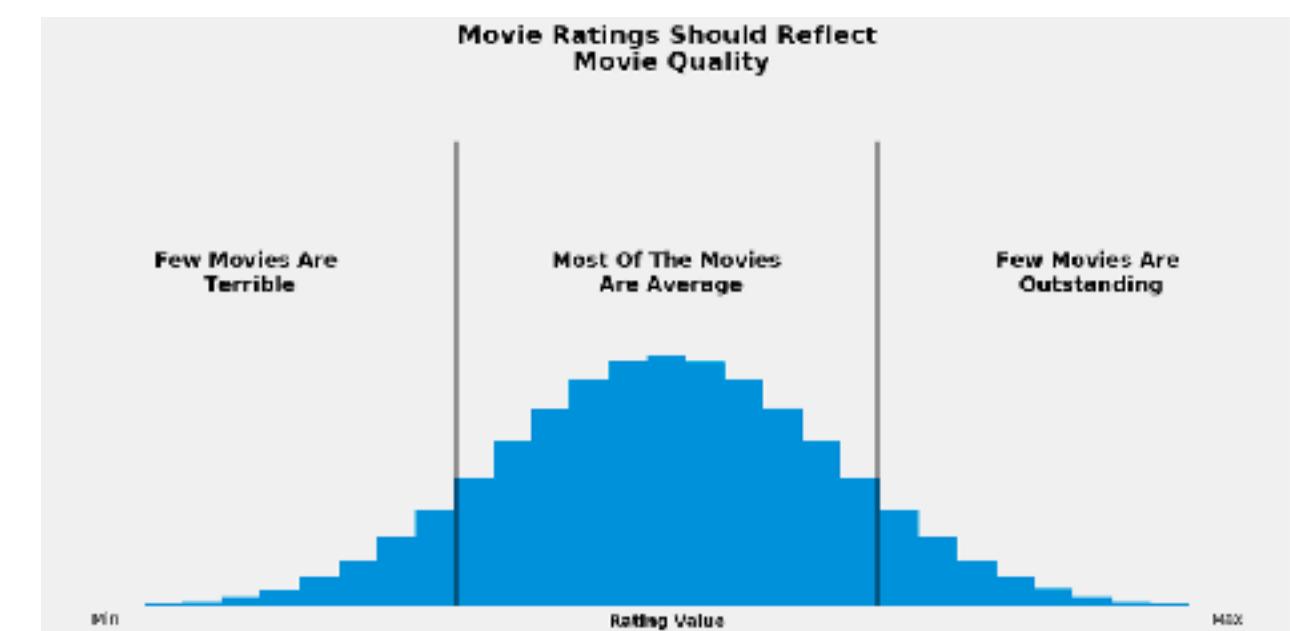


Top Ranking



4.7 4.6 4.5 4.5 4.5 4.5 4.5 4.5 4.5 4.5

- If we serve the recommendation as a top-ranked list **RMSE**, **MAE** and other similar metrics are not appropriate to evaluate how good is the system.



Higher **accuracy** does **not** always mean higher **satisfaction**

UNDERSTANDING ONLINE STAR RATINGS:



Solution: Consider other behaviors

Being accurate is not enough: how accuracy metrics have hurt
recommender systems

SM McNee, J Riedl, JA Konstan

CHI'06 extended abstracts on Human factors in computing systems, 1097-1101

680

2006

**How to evaluate the
ranking?**

Rank Evaluation



Top-N Recommendations

- Popular measures to evaluate the **quality of top-N rankings** are the following:
 - Precision@K
 - Mean Average Precision (MAP)
 - Normalized Discounted cumulative gain (NDCG)

Rank Evaluation



- Recall@K

$$\text{recall}@K = \frac{\text{number of items that the user likes among the top K}}{\text{total number of items that the user likes}}$$

Rank Evaluation



- Precision@K

$$precision@K = \frac{\text{relevant items in the top } k \text{ results}}{k}$$

Rank Evaluation



- **Mean Average Precision (MAP)** calculates the precision at the position of every corrected item in the ranked results list

$$AP = \frac{\sum_k P@k \cdot y_i(k)}{\text{number of relevant documents}}$$

$$MAP = \frac{\sum_{q=1}^Q AP(q)}{Q}$$

Query 1: AP= $\frac{1}{3}(1/1 + 2/3 + 3/6) = 0.72$

Query 2: AP= $\frac{1}{3}(1/1 + 2/2 + 3/4) = 0.917$

$$MAP = (0.72 + 0.917) / 2 = 0.8185$$

Rank Evaluation



Case: Kaggle Challenge Expedia

Submissions are evaluated according to the [Mean Average Precision @ 5 \(MAP@5\)](#):

$$MAP@5 = \frac{1}{|U|} \sum_{u=1}^{|U|} \sum_{k=1}^{\min(5,n)} P(k)$$

where $|U|$ is the number of user events, $P(k)$ is the precision at cutoff k , n is the number of predicted hotel clusters.

Rank Evaluation



- Normalized Discounted cumulative gain (**NDCG**).
 - If there is items which are more relevant to others (e.g. ratings)

$$DCG_p = \sum_{i=1}^p \frac{2^{rel_i} - 1}{\log_2(i + 1)}$$

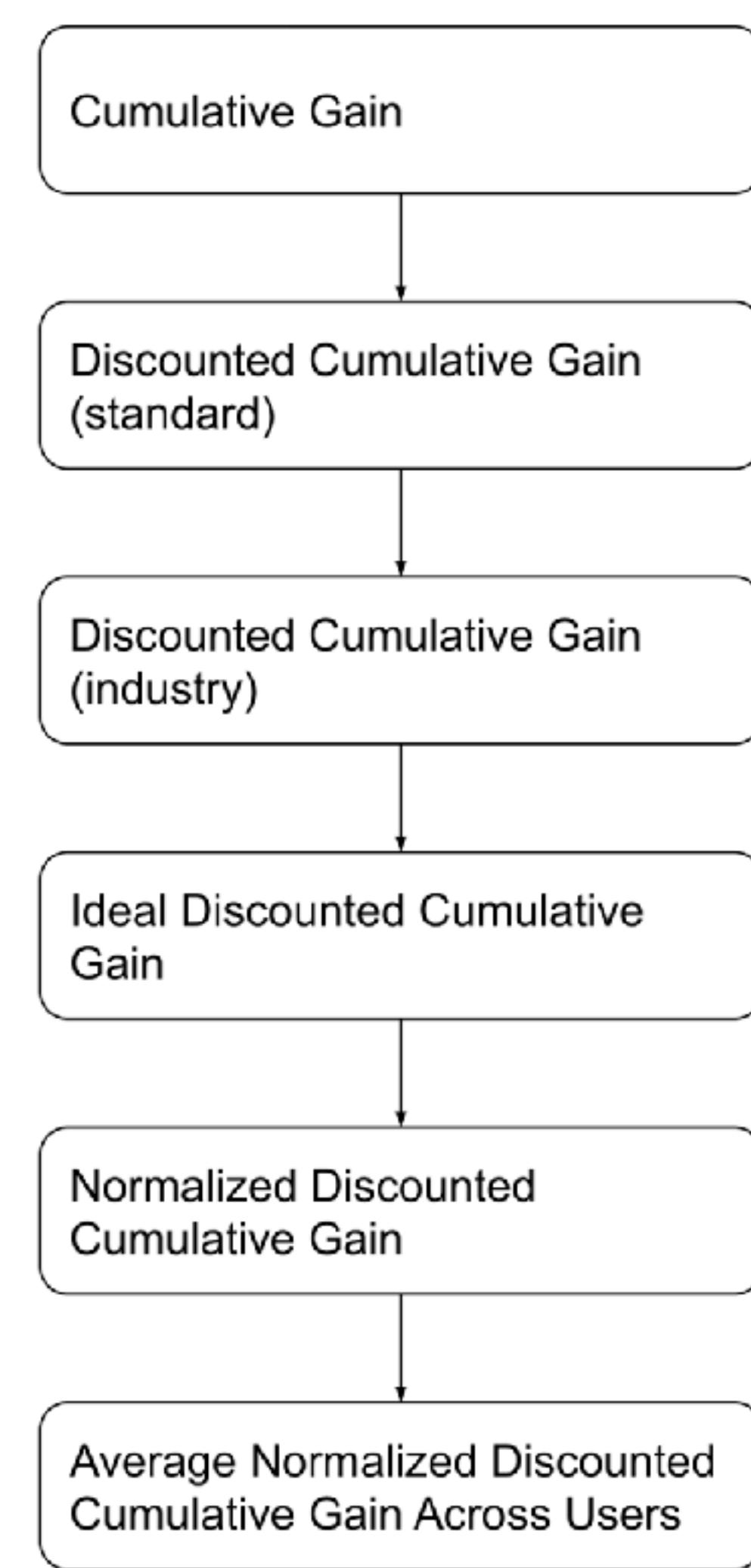
graded relevance to get a stronger emphasis on retrieving relevant documents

$$IDCG_p = \sum_{i=1}^{|REL_p|} \frac{2^{rel_i} - 1}{\log_2(i + 1)}$$

list of relevant documents (ordered by their relevance) in the corpus up to position p

$$nDCG_p = \frac{DCG_p}{IDCG_p}$$

Rank Evaluation



$$CG_p = \sum_{i=1}^p rel_i$$

p elements in the recommended ranking
graded relevance of the result at position i (gain)

$$DCG_p = \sum_{i=1}^p \frac{rel_i}{\log_2(i + 1)}$$

logarithmic reduction factor to penalize proportionally to the position of the result

$$DCG_p = \sum_{i=1}^p \frac{2^{rel_i} - 1}{\log_2(i + 1)}$$

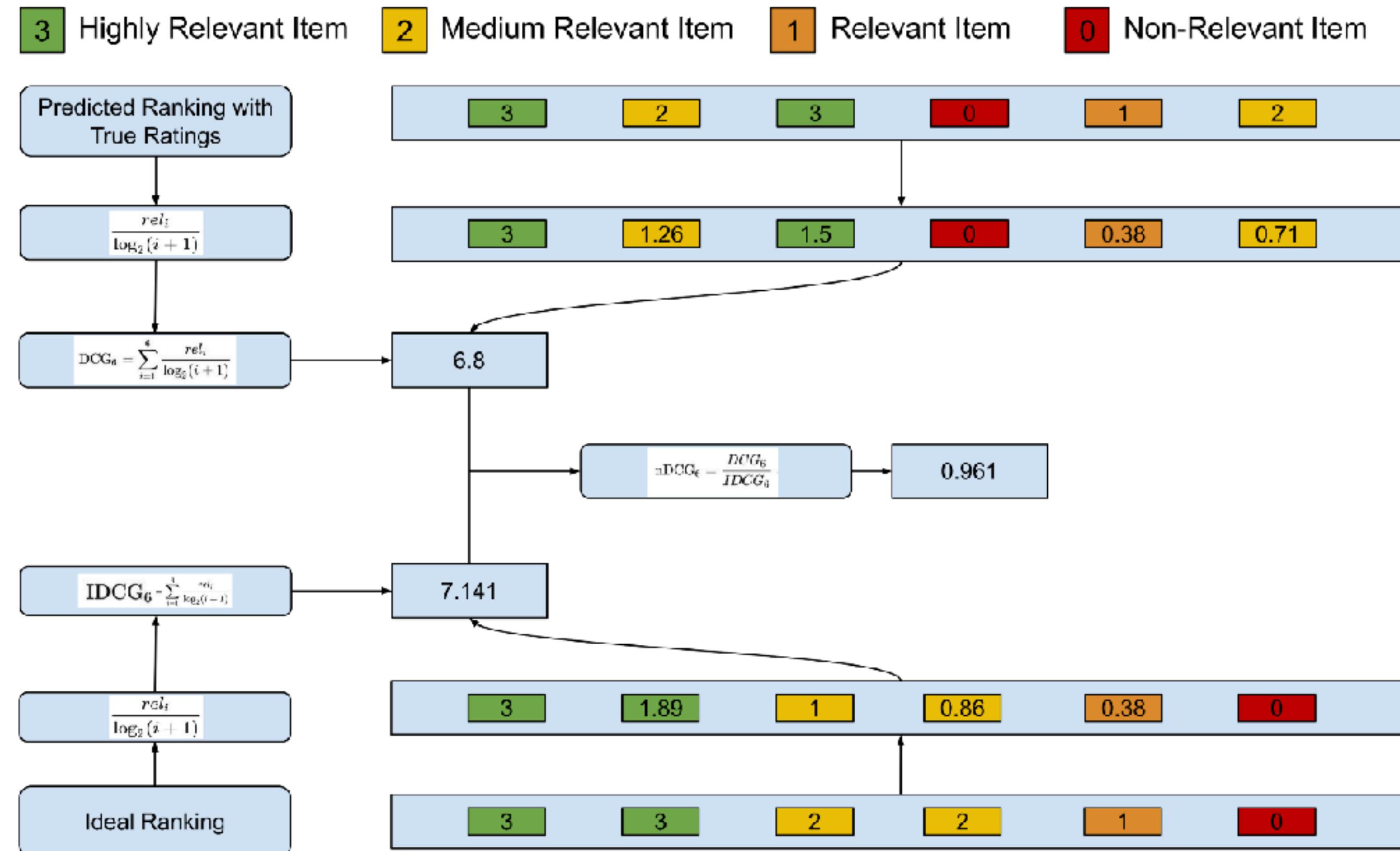
graded relevance to get a stronger emphasis on retrieving relevant documents

$$IDCG_p = \sum_{i=1}^{|REL_p|} \frac{2^{rel_i} - 1}{\log_2(i + 1)}$$

list of relevant documents (ordered by their relevance) in the corpus up to position p

$$nDCG_p = \frac{DCG_p}{IDCG_p}$$

Rank Evaluation



Rank Evaluation



- We can use the strategy followed by P. Cremonesi et.al.
- In order to measure the precision recall, first the models is trained using the training data, and then, for each item i rated with 5 stars in the test data set:
 - A set of 100 random unseen movies for the user of the item i are selected. We assume that these random movies will not be at the same interest than the 5 star movie
 - We predict the rating of the movie of item i and 100 random unseen movies.
 - We form a rank list by ordering all the 101 item according to the predicted rating. Let denote p the rank of the test item i within the list. The best results corresponds to the case the test item i precedes all the random items (i.e., $p=1$).
- A top-N recommendation list by picking the N top ranked items from the list. If $p \leq N$ we have a hit. Otherwise we have a miss. Chances of hit increases as N is higher.

Rank Evaluation

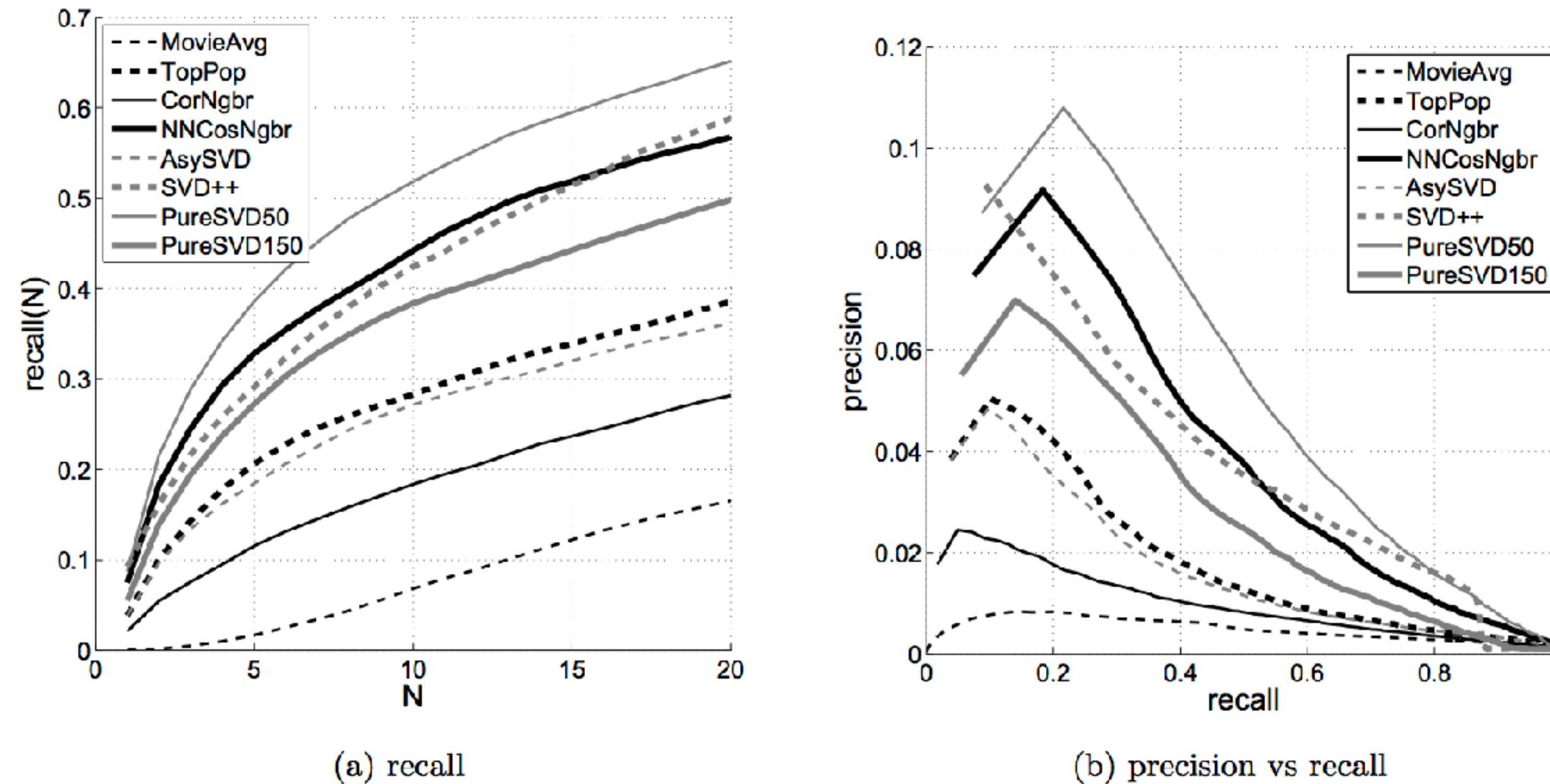


Figure 2: MovieLens: (a) recall-at-N and (b) precision-versus-recall on all items.

Performance of recommender algorithms on top-n recommendation tasks

P Cremonesi, Y Koren, R Turrin

Proceedings of the fourth ACM conference on Recommender systems, 39-46

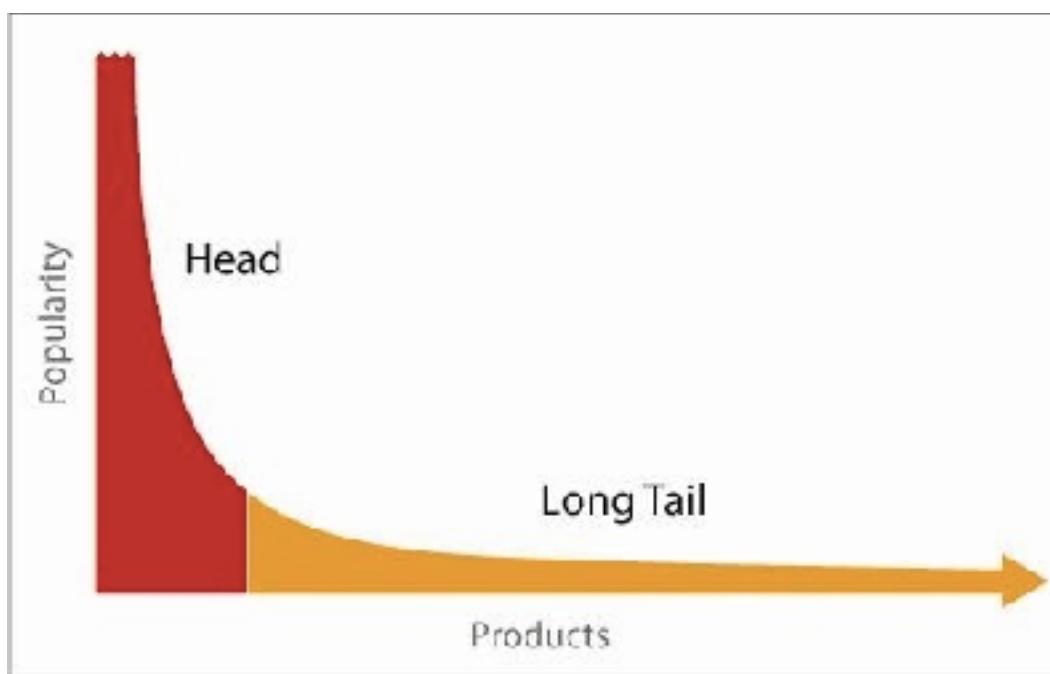
612

2010

Rank Evaluation



- The accuracy measures are usually heavily influence by rating from populars item. Item that received very few items are practically ignored.
- However, the non-popular items are typically the ones which provides higher profit to the companies.



Trick: weight items according to the profit, or remove the those items which are too popular.

How to rank?

The problem:

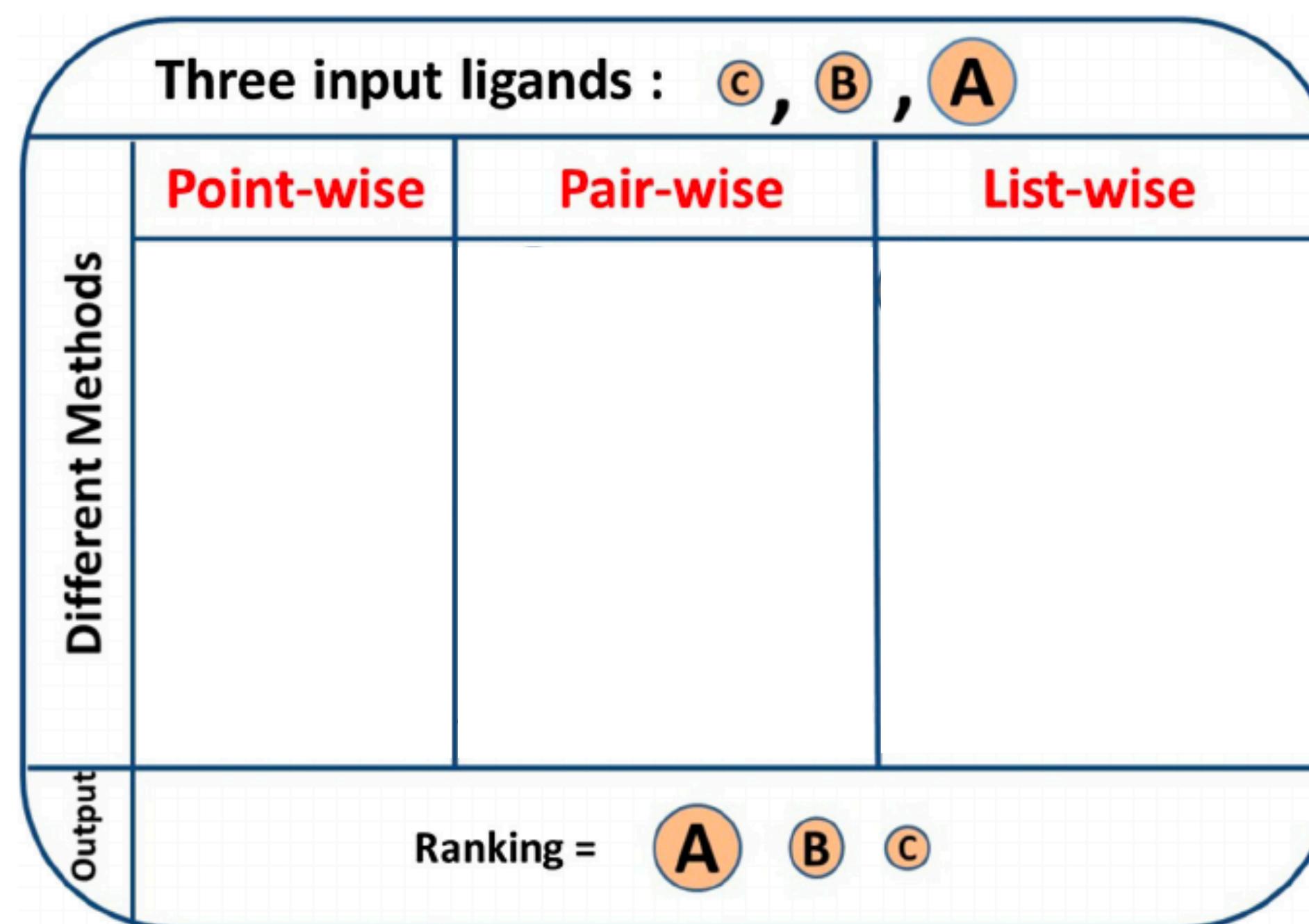
It is **hard to optimize** a machine learning model using the previous **measures** since they are **not differentiable**.

Learning to rank

- Optimization techniques that **learn rank-order directly instead of minimizing prediction error** are referred to as Learning-to-Rank (LTR).
- Three main approaches:
 - **Pointwise:**
 - Predict the absolute relevance (e.g. RMSE)
 - **Pairwise:**
 - Predict the ranking of a document pair (e.g. AUC)
 - **Listwise:**
 - Predict the ranking of a document list (e.g. cross entropy)

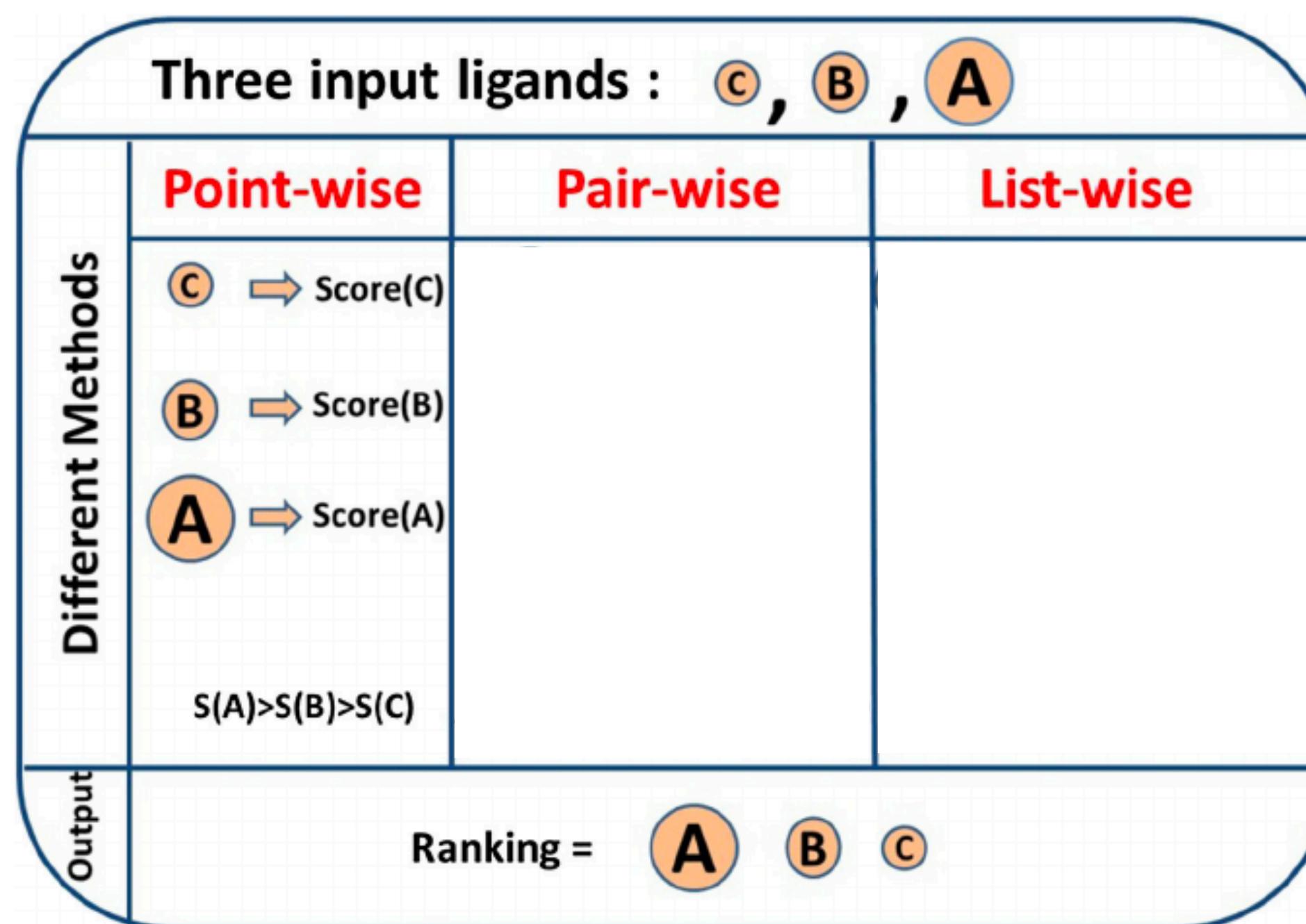
Learning to rank

- Using Machine Learning, the goal is to **construct a ranking model** from the training data.
- The problem can be treat as a standard classification problem



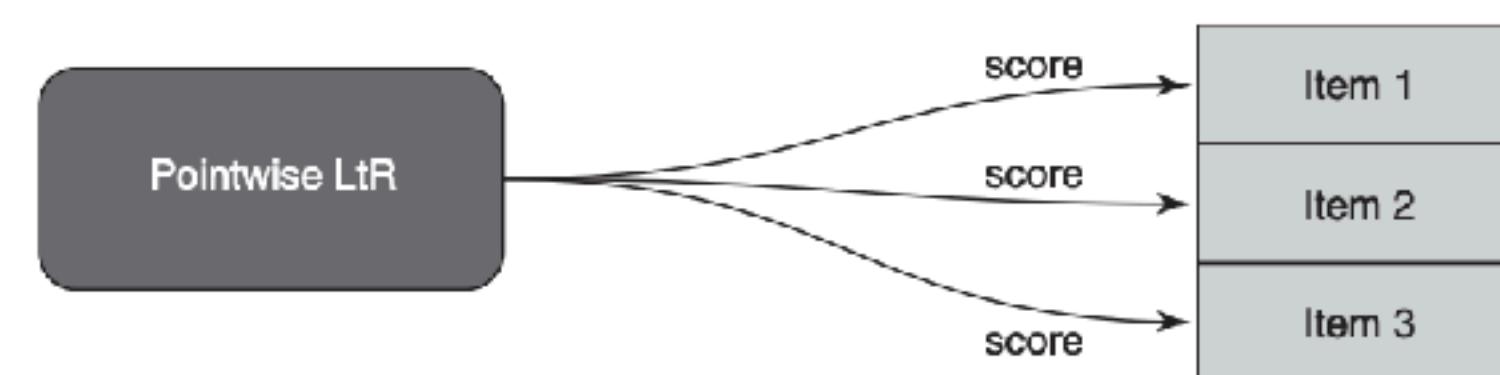
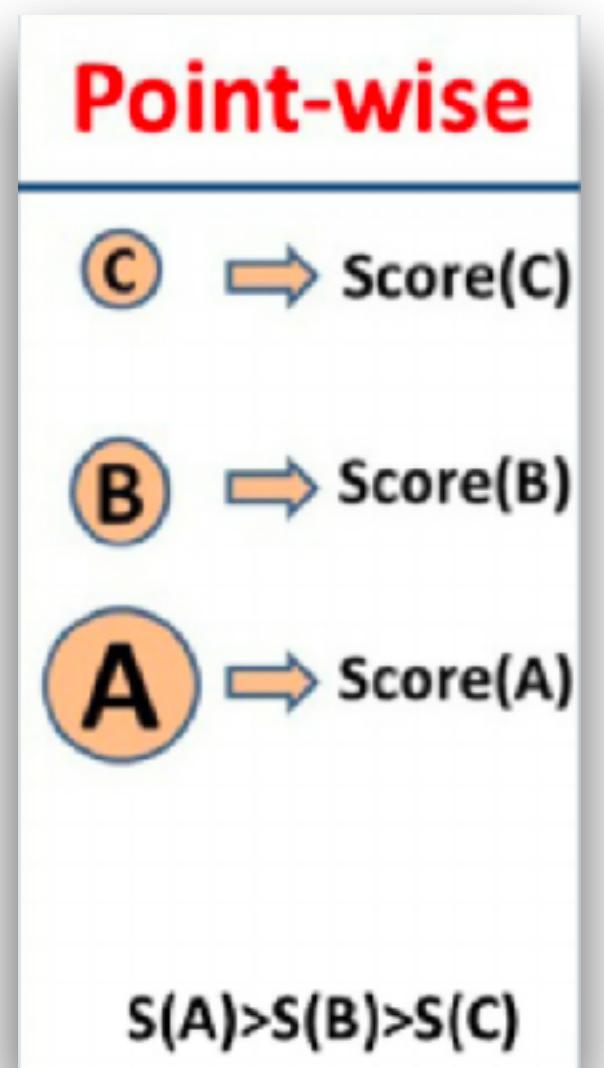
Learning to rank

- Using Machine Learning, the goal is to **construct a ranking model** from the training data.
- The problem can be treat as a standard classification problem

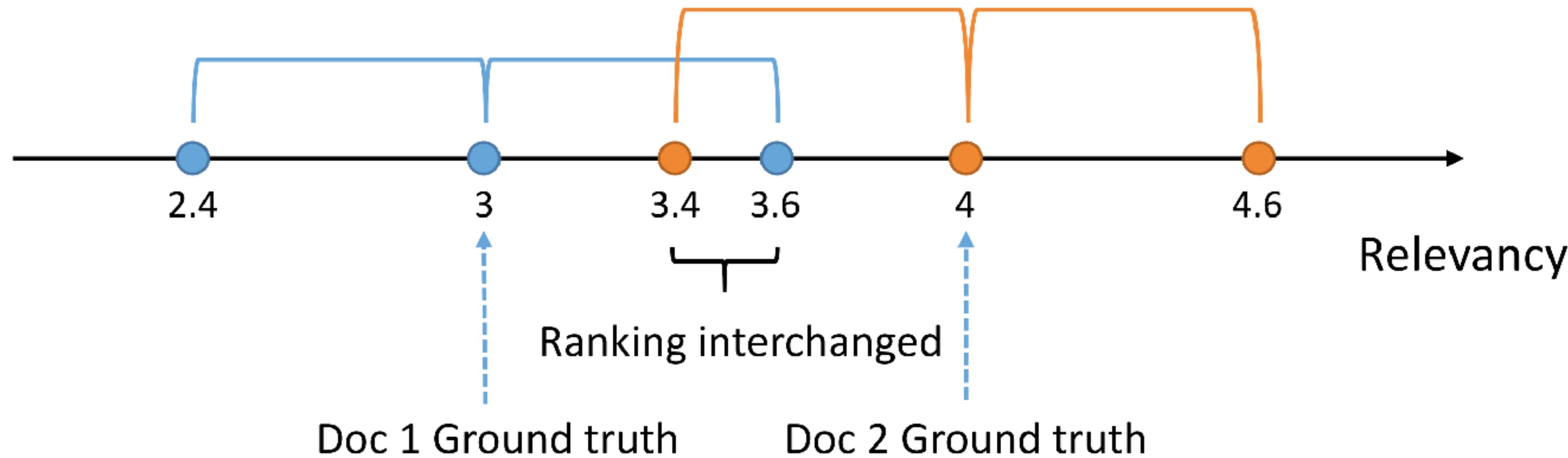


Learning to rank

- Pointwise approaches look at a **single item** at a time in the loss function. They essentially take a single item and train a classifier / regressor on it to predict how relevant it is for the current query/ user. The final ranking is achieved by simply sorting the result list by these document scores.
- For pointwise approaches, the score for each item is independent of the other items that are in the ranking list.
- All the standard regression and classification algorithms (SVM, XGBOOST, NN,...) can be directly used for pointwise learning to rank.
 - These methods tries are trained trying to minimize the error on the rating prediction... Usually RMSE, CE, or other similar loss function...
 - **Problem:** Obtained accuracy is highly biased to popular items.



Point Accuracy != Ranking Accuracy



Same square error might lead to different rankings

Learning to rank

- **Pairwise/Listwise vs. Pointwise**
 - Train on **pair or lists of training samples** instead of individual observations.
 - The loss functions are based on the relative ordering of items instead of the raw orders.
 - doesn't care much about the exact score that each item gets, but cares more about the relative ordering among all the items.

Pairwise approach

Instead of focusing on recommendations as a rating prediction problem, it sometimes makes more sense to look at **how the items should be stacked -> relative preference**

Pairwise approach

- **General Criteria:** The ranking function f learns to rank pairs of items (i.e. for $\{x_i, x_j\}$, is y_i greater than y_j ?).



Pairwise approach

- **General Criteria:** The ranking function f learns to rank pairs of items (i.e. for $\{x_i, x_j\}$, is y_i greater than y_j ?).
 - predict for every **pair of items** based on feature vector x
 - users' **relative preference** regarding the documents.
 - training determines the **parameter** θ based on a **loss function** (e.g. the number of inverted pairs)

Advantage:

Models relative order

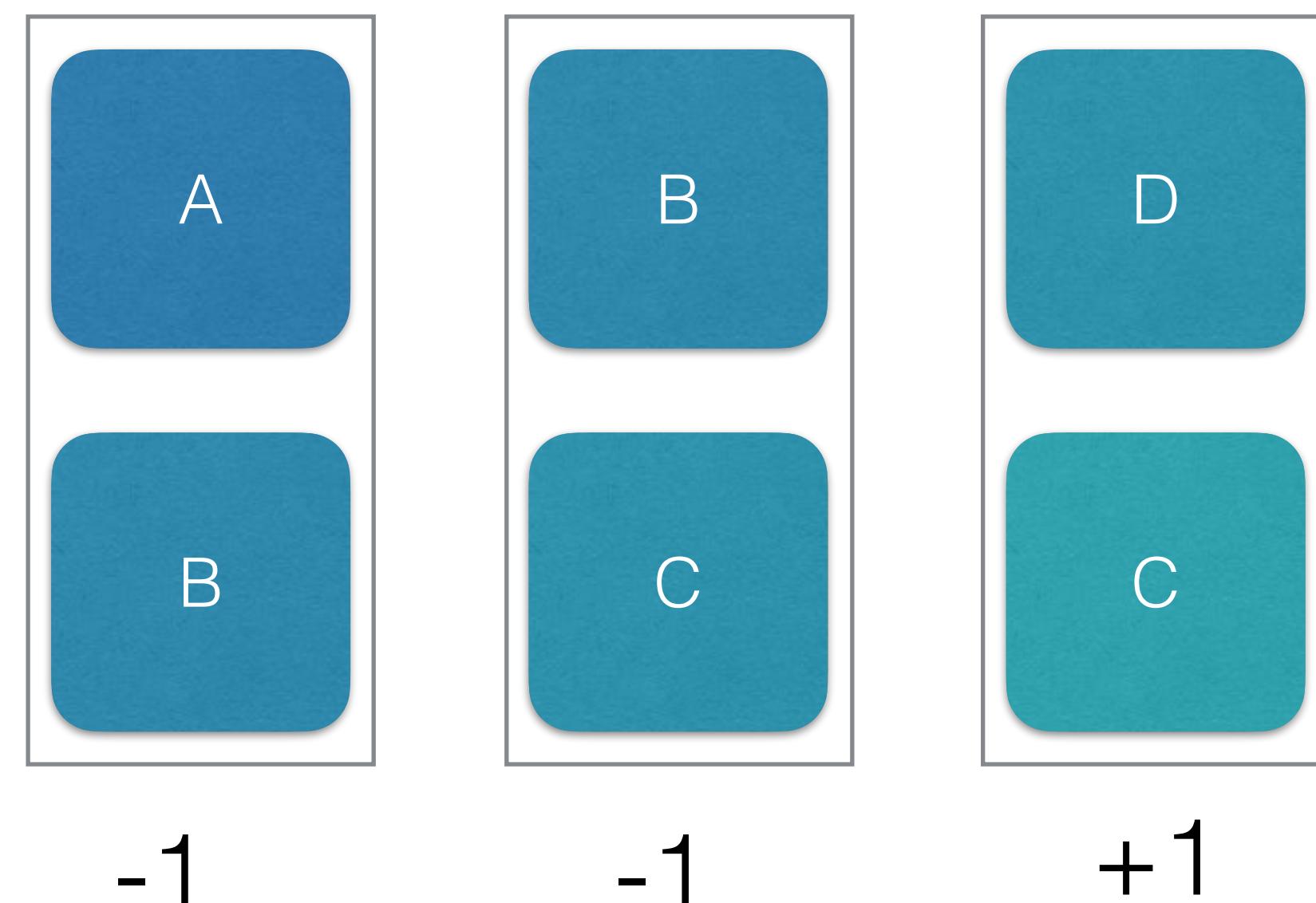
Main disadvantage:

- no distinction between excellent-bad and fair-bad
- sensitive to noise labels.

Pairwise approach

| Data | Item | A | B | C | D | E |
|------|-------|---|---|---|---|---|
| | Score | 1 | 2 | 3 | 4 | 5 |

- Training data instances are item pairs in learning



With labels:
+1 ($i > j$) or **-1** ($j < i$)

Pairwise approach

- **Several Methods:**
 - RANK SVM
 - RankBoost
 - RANK NET
 - Lambda Rank
 - Lambda Mart
 - Bayesian Personalized Ranking (BPR)

Pairwise approach

RankNet

- RankNet was originally developed using neural nets.
- **The cost function for RankNet aims to minimize the number of *inversions* in ranking.** Here an inversion means an incorrect order among a pair of results, i.e. when we rank a lower rated result above a higher rated result in a ranked list. RankNet optimizes the cost function using Stochastic Gradient Descent.

Learning to Rank using Gradient Descent

Keywords: ranking, gradient descent, neural networks, probabilistic cost functions, internet search

Chris Burges

CRURGES@MICROSOFT.COM

Tal Shaked*

TAL.SHAKED@GMAIL.COM

Erin Renshaw

ERINREN@MICROSOFT.COM

Microsoft Research, One Microsoft Way, Redmond, WA 98052-6399

Ari Lazier

ARIEL@MICROSOFT.COM

Matt Deeds

MADEEDS@MICROSOFT.COM

Nicole Hamilton

NICHAM@MICROSOFT.COM

Greg Hullender

GREGHULL@MICROSOFT.COM

Microsoft, One Microsoft Way, Redmond, WA 98052-6399

Pairwise approach

RankNet

- RankNet [Burges et al., 2005] is a pairwise loss function—popular choice for training neural L2R models and also an industry favorite [Burges, 2015]

- Predicted probabilities:

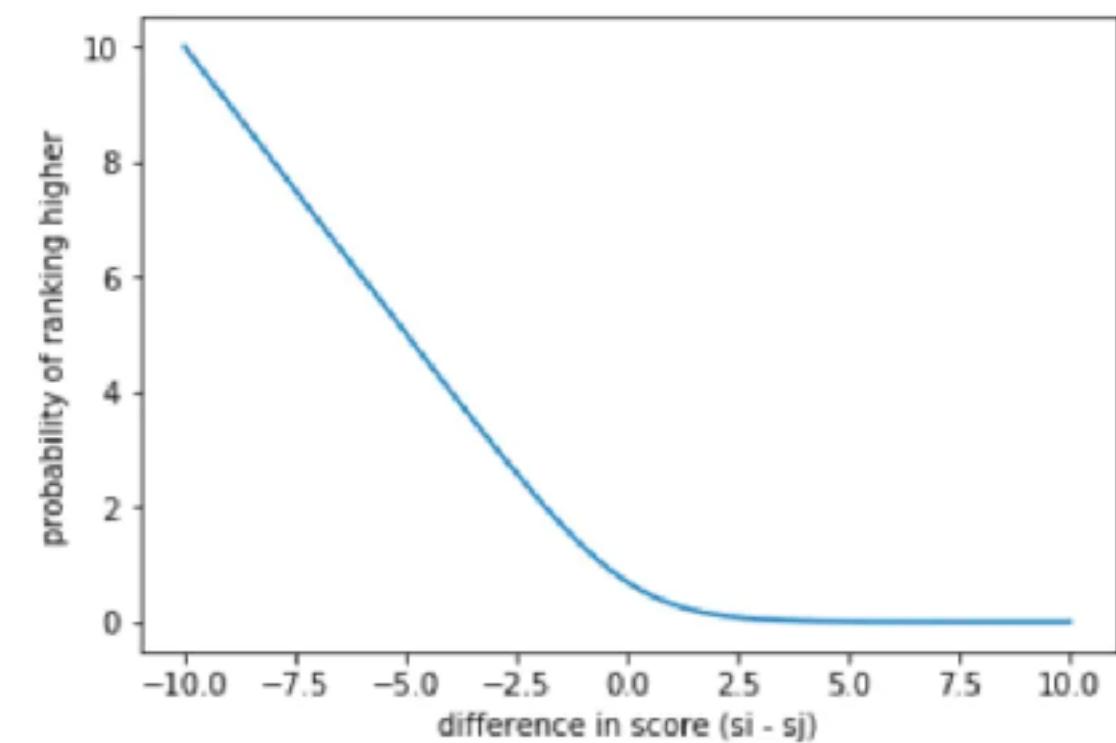
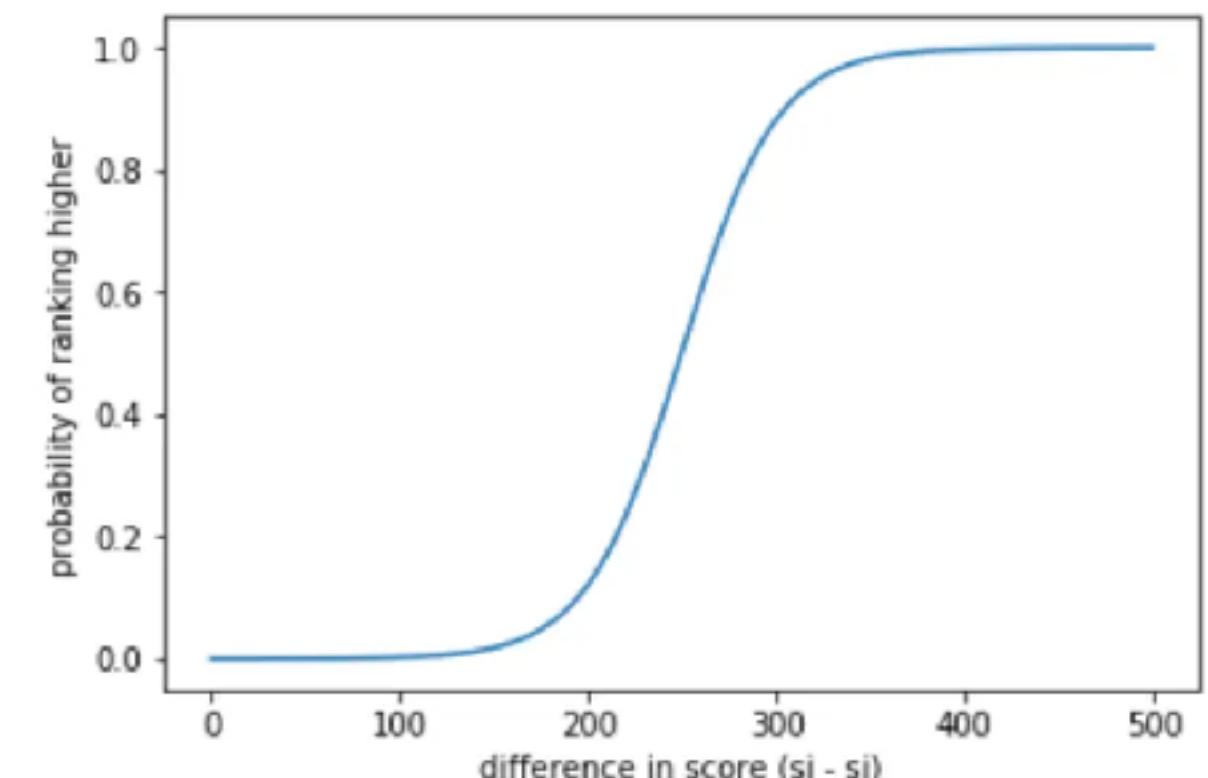
$$p_{i,j} = p(s_i > s_j) = \frac{1}{1 + e^{-\gamma(s_i - s_j)}}$$
$$p_{j,i} = \frac{1}{1 + e^{-\gamma(s_j - s_i)}}$$

- Desired probabilities:

$$p_{ij} = 1 \text{ and } p_{ji} = 0$$

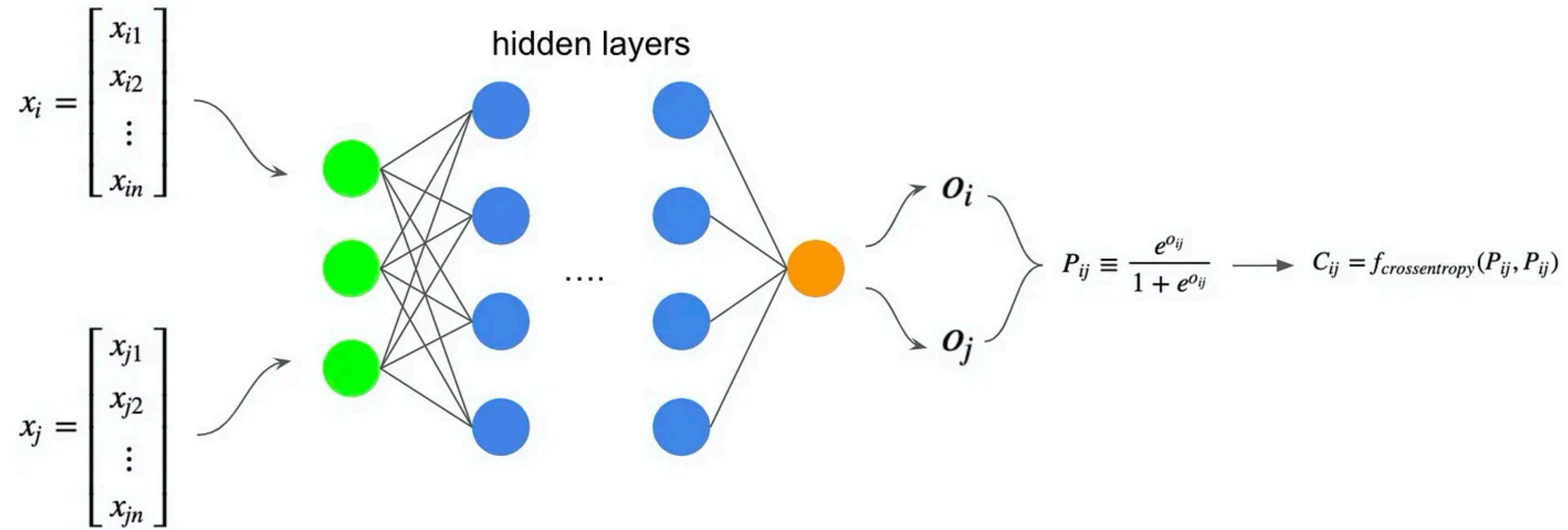
- Computing cross-entropy between y and p

$$L_{RankNet} = y_{ij} \log(p_{ij}) - y_{ji} \log(p_{ji}) = y_{ij} \log(p_{ij}) - (1 - y_{ij}) \log(1 - p_{ij})$$



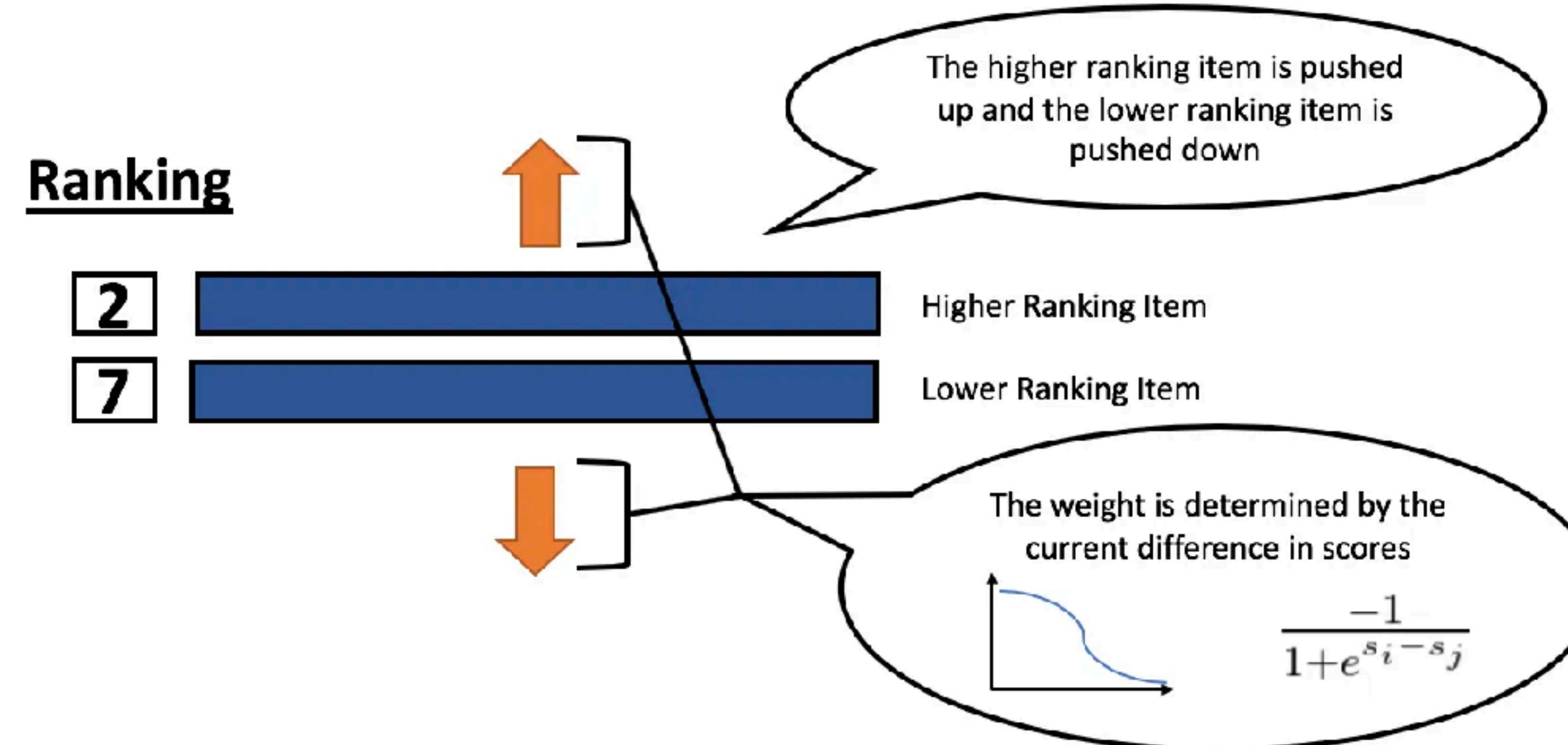
Pairwise approach

RankNet



Pairwise approach

RankNet



Pairwise approach

RankNet

RankNet and is pretty good in many situations. There is one major pitfall to the RankNet approach though: the loss is agnostic to the **actual ranking** of the item.

The strength of the push is determined only by the current difference in scores. So it doesn't matter if the items that rank below are 1, 2, or 500 ranks below

| Documents | Rating | |
|-----------|--------|---|
| — | 2 | — |
| — | 4 | — |
| — | 3 | |
| — | 2 | — |
| — | 4 | — |

Same pair-level error
but different list-level
error

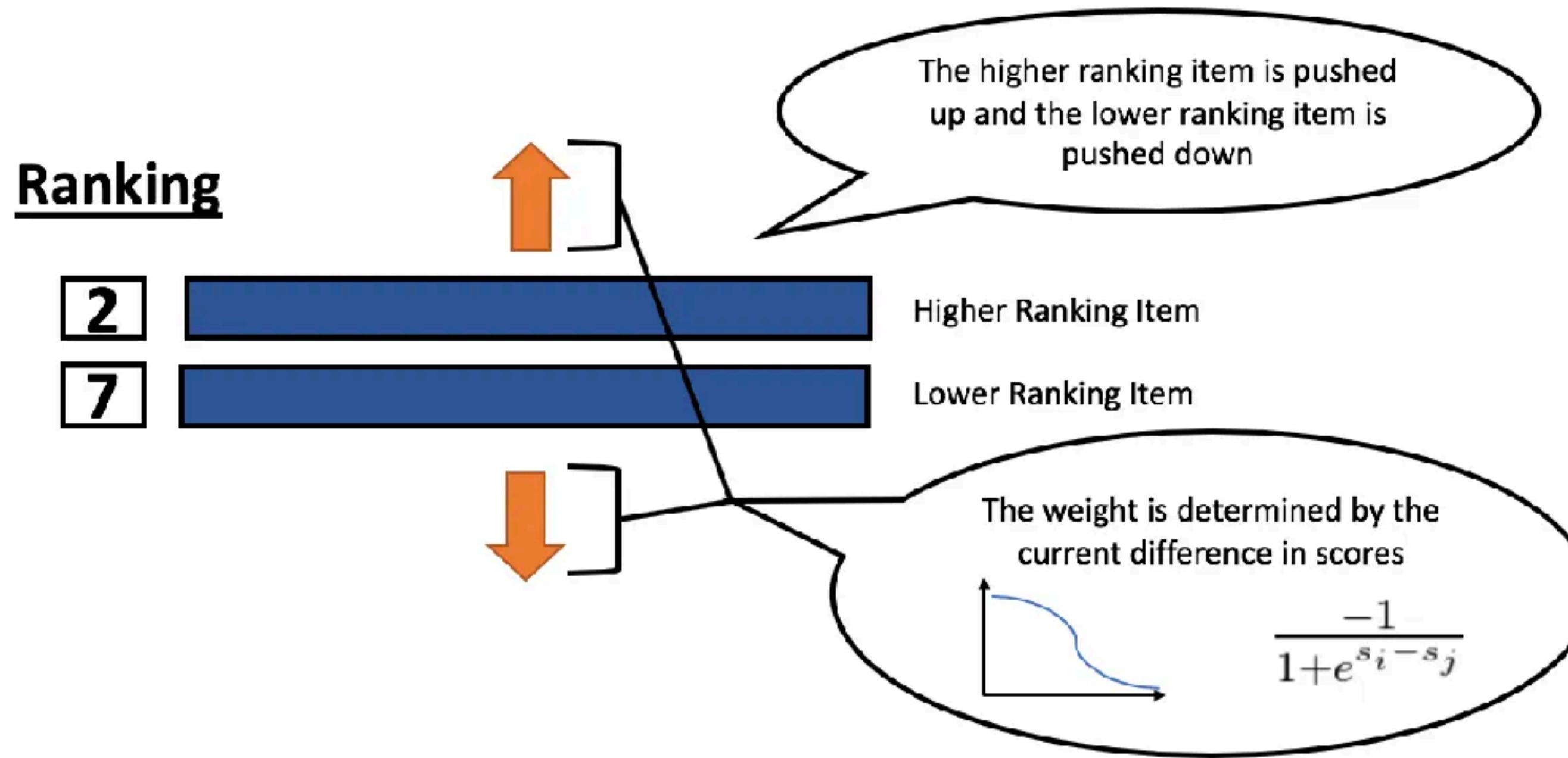
Pairwise approach

LambdaRank

- Burgess et. al. found that during RankNet training procedure, you don't need the costs, only need the gradients (λ) of the cost with respect to the model score. You can think of these gradients as little arrows attached to each document in the ranked list, indicating the direction we'd like those documents to move.



- Further they found that scaling the gradients by the change in NDCG found by swapping each pair of documents gave good results. The core idea of LambdaRank is to use this new cost function for training a RankNet. On experimental datasets, this shows both speed and accuracy improvements over the original RankNet.



$$\text{Now, } \lambda_{ij} = \frac{-\Delta(i,j)}{1 + e^{(s_i - s_j)}}$$

where $\Delta(i,j)$ is a penalty corresponding to how “bad” it is to rank items i and j in the wrong order.

Multiple metrics can be used to compute Δ , but the most common is the NDCG (normalized discounted cumulative gain)

Pairwise approach

Bayesian Personalized Ranking (BPR)

- With implicit feedback, the notion of “ratings” become less precise. As a result, matrix factorization models moved away from estimating a specific set of ratings to, instead, model the relative preferences (or orders) between different items.
- Let D_i be a set of item pairs (j, k) where user i has interacted with item j but not item k , assuming user i might be more interested in item j than item k , BPR minimizes the pairwise ranking loss

$$\min_{\mathbf{u}_*, \mathbf{v}_*} \sum_{i \in \mathcal{I}} \sum_{(j, k) \in \mathcal{D}_i} -\log \sigma(\mathbf{u}_i^T \mathbf{v}_j - \mathbf{u}_i^T \mathbf{v}_k) + \lambda_u \|\mathbf{u}_i\|^2 + \lambda_v \|\mathbf{v}_j\|^2,$$

- where σ is the sigmoid function.
- Problem: BPR loss does not sufficiently penalize the items that are at a lower rank

BPR: Bayesian personalized ranking from implicit feedback
[S Rendle, C Freudenthaler, Z Gantner... - arXiv preprint arXiv ..., 2012 - arxiv.org](#)
... Item recommendation is the task of predicting a **personalized ranking** ... **implicit feedback** (eg clicks, purchases). There are many methods for item recommendation from **implicit feedback** ...
☆ Save ⚡ Cite Cited by 5233 Related articles All 15 versions ⚡

Pairwise approach

Bayesian Personalized Ranking (BPR)

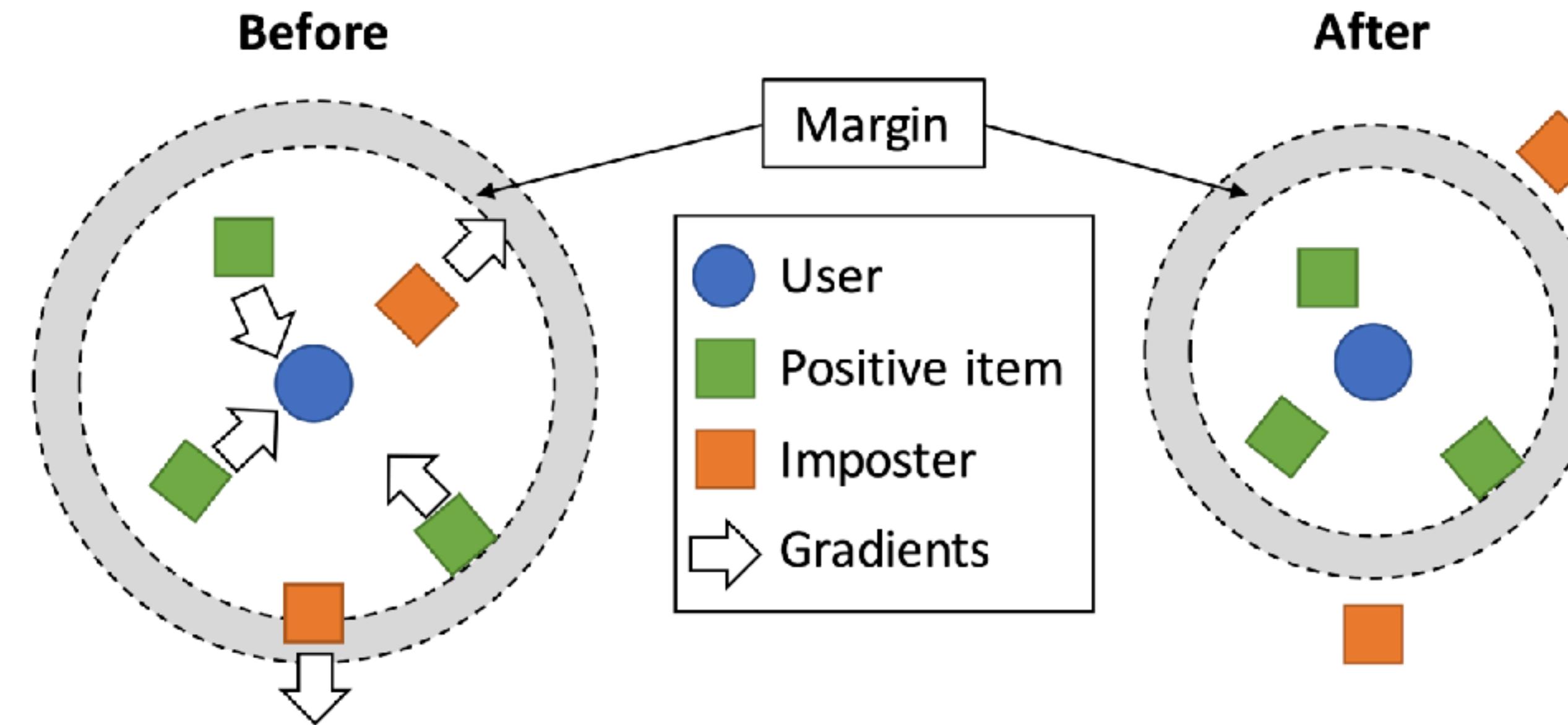
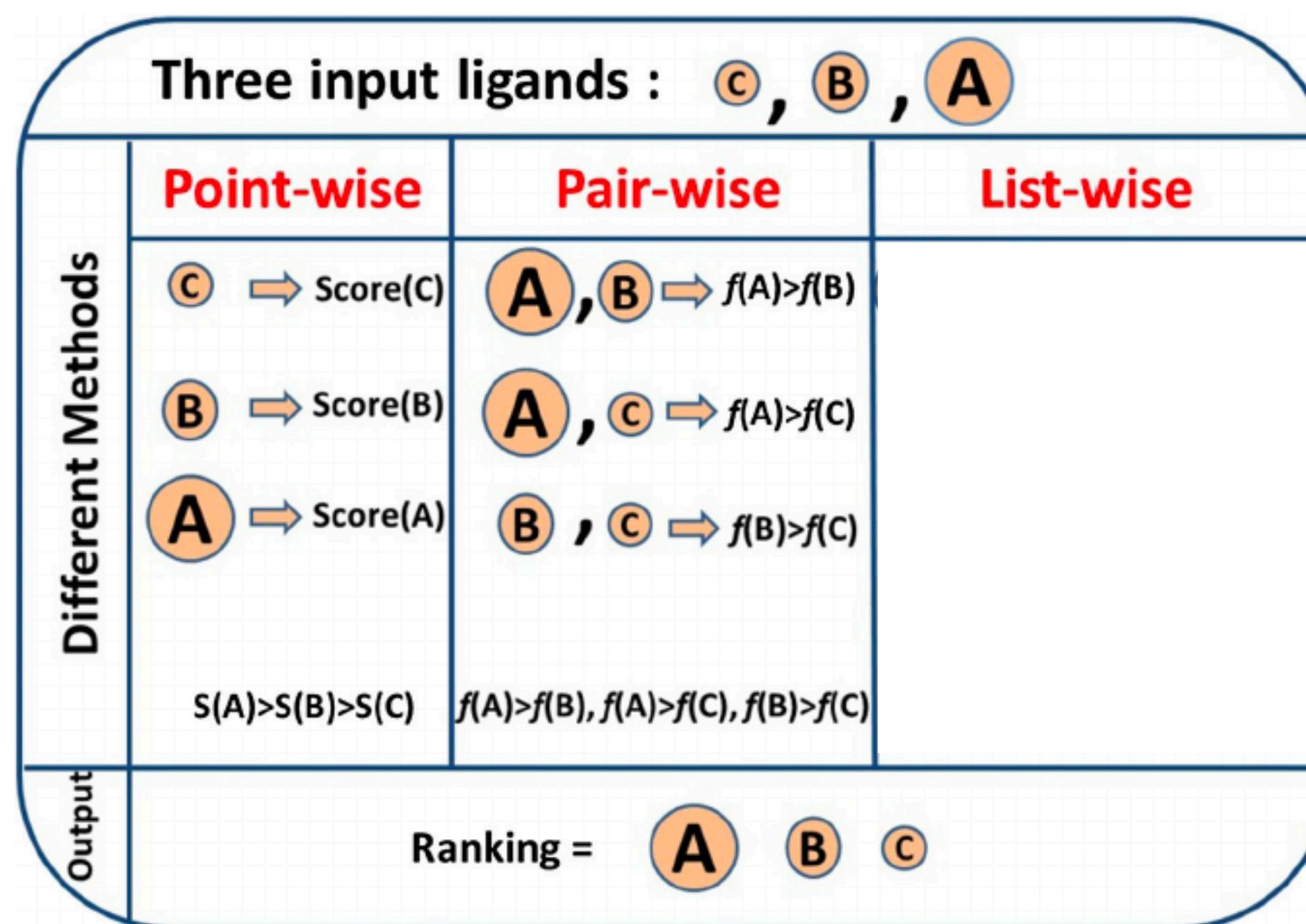


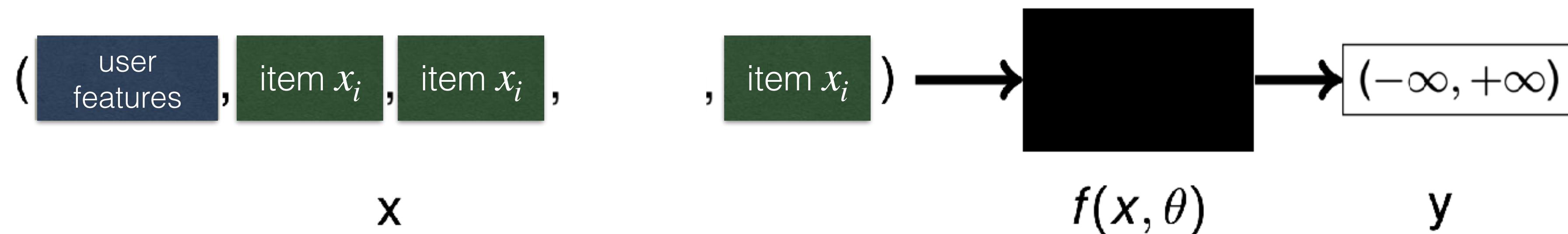
Figure 1: An illustration of collaborative metric learning. The hinge loss defined in Eq. 1 creates a gradient that *pulls* positive items closer to the user and *pushes* the intruding impostor items (i.e., items that the user did not like) away until they are beyond the safety margin.

Learning to rank

- Using Machine Learning, the goal is to **construct a ranking model** from the training data.
- The problem can be treat as a standard classification problem



Listwise approach



- predict for ranked list of documents based on feature vector x
- **effectiveness** of ranked list \mathbf{y} (e.g., MAP or nDCG)
- training determines the **parameter** θ based on a **loss function**

Advantage:

positional information visible to loss function

Main disadvantage:

- high training complexity

Lambda Mart

- LambdaMART combines LambdaRank and MART (Multiple Additive Regression Trees)
- MART uses gradient boosted decision trees for prediction tasks, LambdaMART uses gradient boosted decision trees using a cost function derived from LambdaRank for solving a ranking task

The screenshot shows a Microsoft Research publication page. At the top, there's a navigation bar with links for Microsoft, Research, Research areas, Researcher tools, Programs & Events, Careers, People, Blogs & Learning, and Labs & Locations. Below the navigation is the title "From RankNet to LambdaRank to LambdaMART: An Overview" by Chris J.C. Burges. It includes a link to "Download BibTeX" and a "View Publication" button. A summary text states: "LambdaMART is the boosted tree version of LambdaRank, which is based on RankNet. RankNet, LambdaRank, and LambdaMART have proven to be very successful algorithms for solving real world ranking problems: for example an ensemble of LambdaMART rankers won Track 1 of the 2010 Yahoo! Learning To Rank Challenge. The details of these algorithms are spread across several papers and reports, and so here we give a self-contained, detailed and complete description of them." Below the summary are "Research Areas" (Systems and networking) and a "View Publication" button. The URL at the bottom is <https://www.microsoft.com/en-us/research/publication/from-ranknet-to-lambdarank-to-lambdamart-an-overview/>.

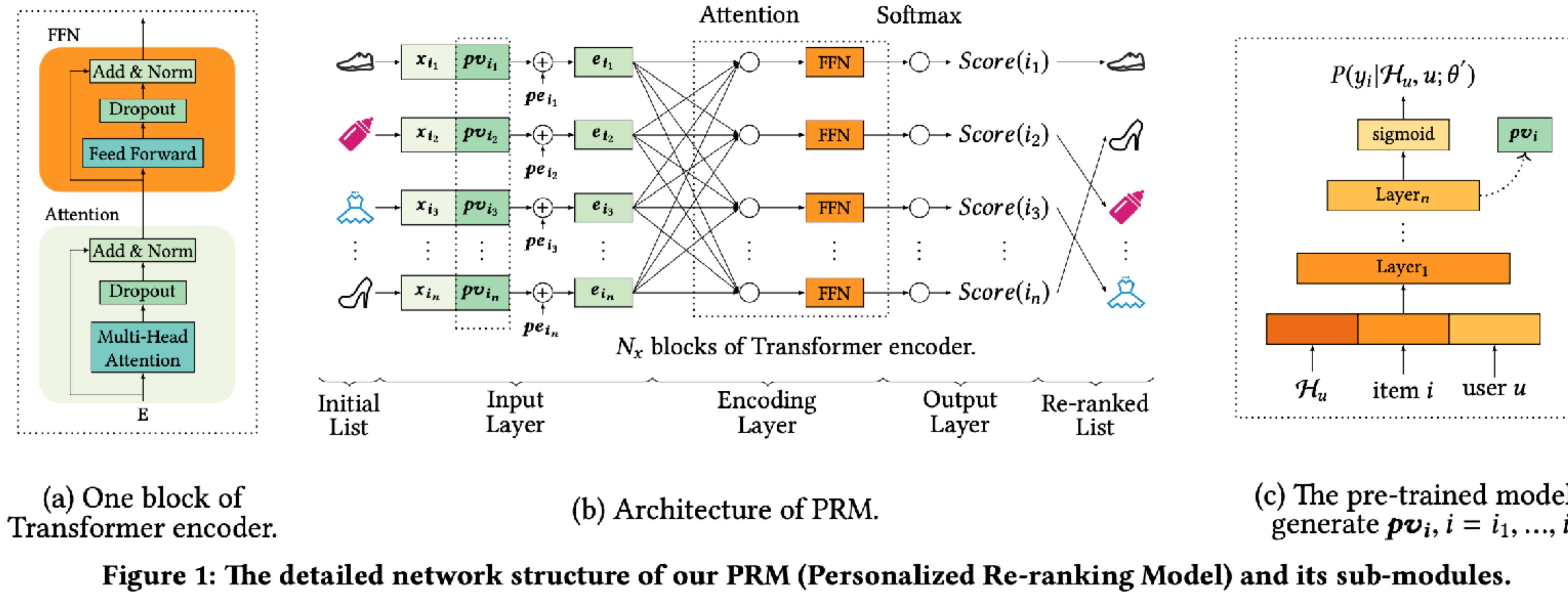
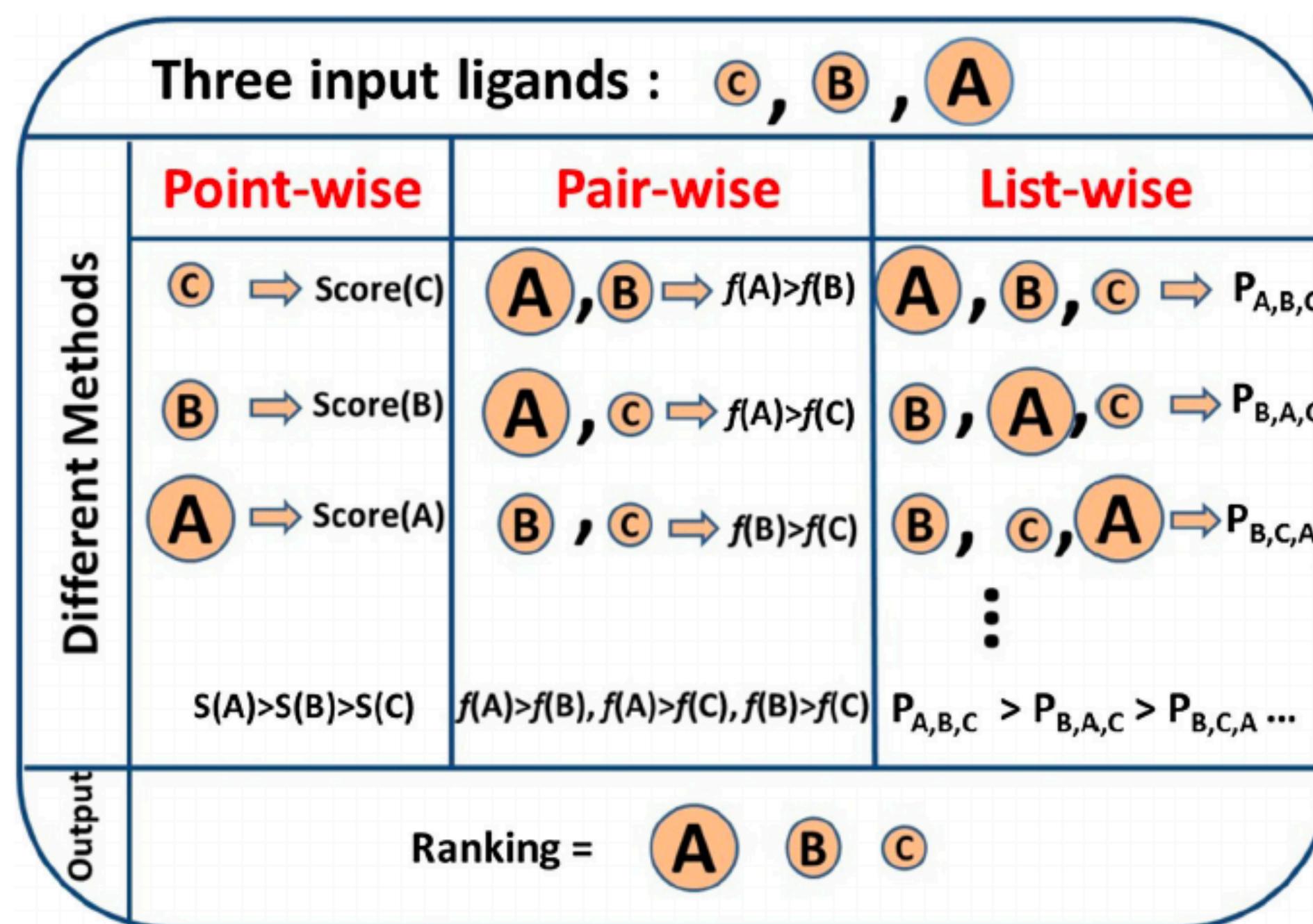


Figure 1: The detailed network structure of our PRM (Personalized Re-ranking Model) and its sub-modules.

Learning to rank

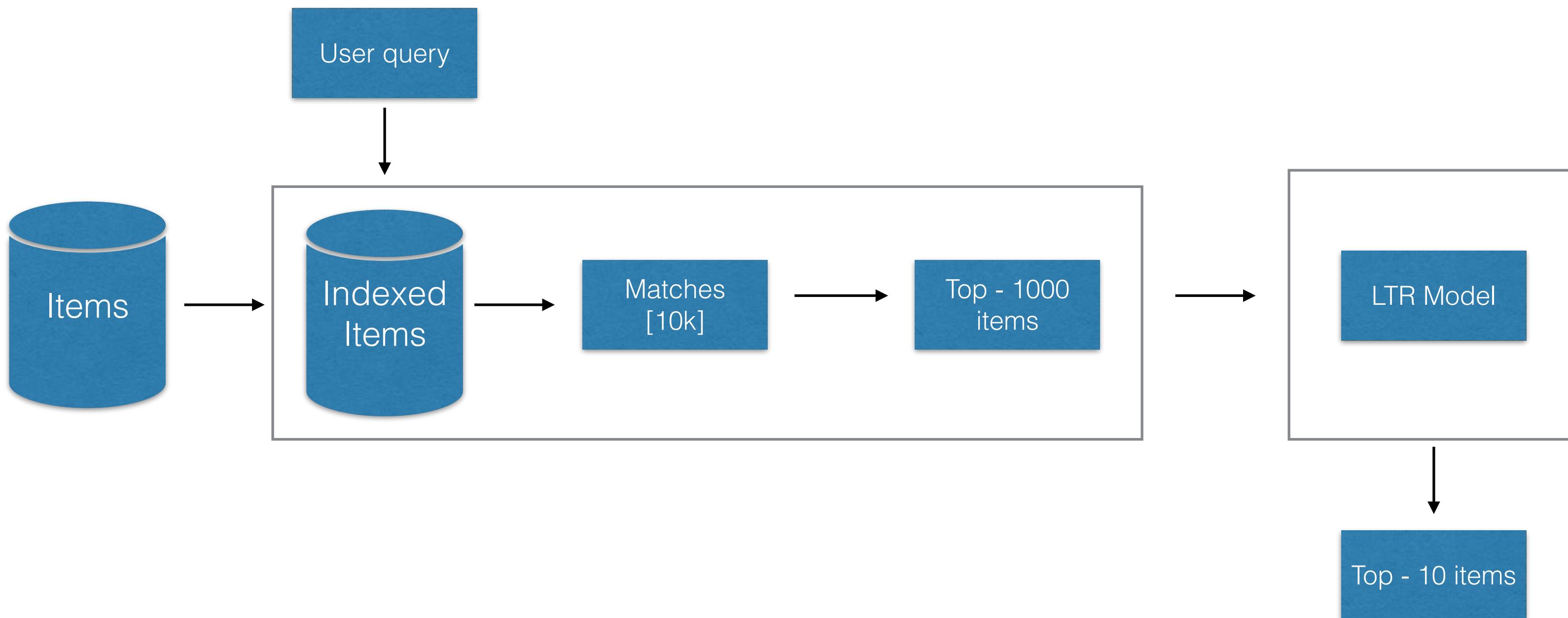
- Using Machine Learning, the goal is to **construct a ranking model** from the training data.
- The problem can be treat as a standard classification problem



**LTR are great but
task consuming....**

LTR framework

- LTR methods are computationally expensive.
- Usually used as re-rankers



If you have implicit data

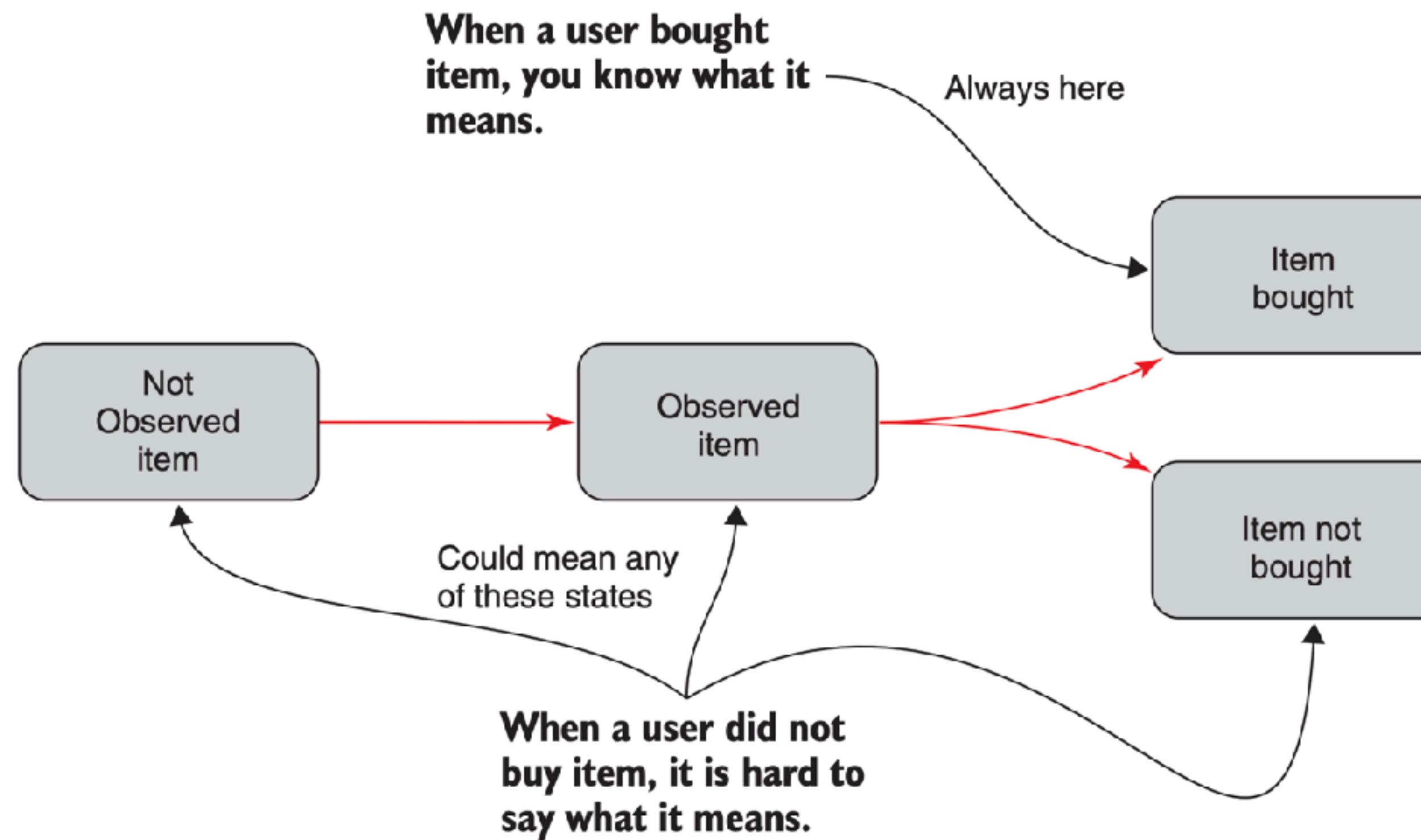


Figure 13.6 Different states of a user-item relationship. You know that when a user buys an item it's bought, but if a user doesn't buy an item, what does that mean?

Presentations

TASK 2

Paper presentations. 8 -10 minutes

| Team | Date | Students | Paper |
|-------------|-------------|-----------------|--------------|
| 1 | March 16th | | |
| 2 | March 16th | | |
| 3 | March 16th | | |
| 4 | March 16th | | |
| 5 | March 26th | | |
| 6 | March 26th | | |
| 7 | March 26th | | |
| 8 | March 26th | | |

Paper Suggestions I

Local Item-Item Models for Top-N Recommendation

Evangelia Christakopoulou and George Karypis
Computer Science & Engineering
University of Minnesota, Minneapolis, MN
{evangel,karypis}@cs.umn.edu

ABSTRACT

Item-based approaches based on SLIM (Sparse LInear Methods) have demonstrated very good performance for top- N recommendation; however they only estimate a single model for all the users. This work is based on the intuition that not all users behave in the same way – instead there exist subsets of like-minded users. By using different item-item models for these user subsets, we can capture differences in their prefer-

based methods, which include item k-NN [8] and Sparse LInear Methods (SLIM) [16] have been shown to outperform the user-based schemes for the top- N recommendation task.

However, item-based methods have the drawback of estimating only a single model for all users. In many cases, there are differences in users' behavior, which cannot be captured by a single model. For example, there could be a pair of items that are extremely similar for a specific user subset, while they have low similarity for another user subset. D...

Paper Suggestions I

Paper Suggestions I

Restricted Boltzmann Machines for Collaborative Filtering

Ruslan Salakhutdinov

Andriy Mnih

Geoffrey Hinton

University of Toronto, 6 King's College Rd., Toronto, Ontario M5S 3G4, Canada

RSALAKHU@CS.TORONTO.EDU

AMNIH@CS.TORONTO.EDU

HINTON@CS.TORONTO.EDU

Abstract

Most of the existing approaches to collaborative filtering cannot handle very large data sets. In this paper we show how a class of two-layer undirected graphical models, called Restricted Boltzmann Machines (RBM's), can be used to model tabular data, such as user's ratings of movies. We present efficient learning and inference procedures for this class of models and demonstrate that RBM's can be successfully applied to the Netflix data set, containing over 100 million user/movie ratings. We also show that RBM's slightly outperform carefully-tuned SVD models. When the predictions of multiple RBM models and multiple SVD models are linearly combined, we achieve an error rate that is well over 6% better than the score of Netflix's own system.

Low-rank approximations based on minimizing the sum-squared distance can be found using Singular Value Decomposition (SVD). In the collaborative filtering domain, however, most of the data sets are sparse, and as shown by Srebro and Jaakkola (2003), this creates a difficult non-convex problem, so a naive solution is not going work.¹

In this paper we describe a class of two-layer undirected graphical models that generalize Restricted Boltzmann Machines to modeling tabular or count data (Welling et al., 2005). Maximum likelihood learning is intractable in these models, but we show that learning can be performed efficiently by following an approximation to the gradient of a different objective function called "Contrastive Divergence" (Hinton, 2002).

2. Restricted Boltzmann Machines (RBM's)

Paper Suggestions I

Restricted Boltzmann Machines for Collaborative Filtering

Ruslan Salakhutdinov

Andriy Mnih

Geoffrey Hinton

University of Toronto, 6 King's College Rd., Toronto, Ontario M5S 3G4, Canada

RSALAKHU@CS.TORONTO.EDU

AMNIH@CS.TORONTO.EDU

HINTON@CS.TORONTO.EDU

Abstract

Most of the existing approaches to collaborative filtering cannot handle very large data sets. In this paper we show how a class of two-layer undirected graphical models, called Restricted Boltzmann Machines (RBM's), can be used to model tabular data, such as user's ratings of movies. We present efficient learning and inference procedures for this class of models and demonstrate that RBM's can be successfully applied to the Netflix data set, containing over 100 million user/movie ratings. We also show that RBM's slightly outperform carefully-tuned SVD models. When the predictions of multiple RBM models and multiple SVD models are linearly combined, we achieve an error rate that is well over 6% better than the score of Netflix's own system.

Low-rank approximations based on minimizing the sum-squared distance can be found using Singular Value Decomposition (SVD). In the collaborative filtering domain, however, most of the data sets are sparse, and as shown by Srebro and Jaakkola (2003), this creates a difficult non-convex problem, so a naive solution is not going work.¹

In this paper we describe a class of two-layer undirected graphical models that generalize Restricted Boltzmann Machines to modeling tabular or count data (Welling et al., 2005). Maximum likelihood learning is intractable in these models, but we show that learning can be performed efficiently by following an approximation to the gradient of a different objective function called "Contrastive Divergence" (Hinton, 2002).

2. Restricted Boltzmann Machines (RBM's)

Paper Suggestions II



Recommending music on Spotify with deep learning

AUGUST 05, 2014

This summer, I'm interning at Spotify in New York City, where I'm working on content-based music recommendation using convolutional neural networks. In this post, I'll explain my approach and show some preliminary results.

Overview

This is going to be a long post, so here's an overview of the different sections. If you want to skip ahead, just click the section title to go there.

<http://benanne.github.io/2014/08/05/spotify-cnns.html>

Paper Suggestions III

Deep content-based music recommendation

Aäron van den Oord, Sander Dieleman, Benjamin Schrauwen

Electronics and Information Systems department (ELIS), Ghent University

{aaron.vandenoord, sander.dieleman, benjamin.schrauwen}@ugent.be

Abstract

Automatic music recommendation has become an increasingly relevant problem in recent years, since a lot of music is now sold and consumed digitally. Most recommender systems rely on collaborative filtering. However, this approach suffers from the cold start problem: it fails when no usage data is available, so it is not effective for recommending new and unpopular songs. In this paper, we propose to use a latent factor model for recommendation, and predict the latent factors from music audio when they cannot be obtained from usage data. We compare a traditional approach using a bag-of-words representation of the audio signals with deep convolutional neural networks, and evaluate the predictions quantitatively and qualitatively on the Million Song Dataset. We show that using predicted latent factors produces sensible recommendations, despite the fact that there is a large semantic gap between the characteristics of a song that affect user preference and the corresponding audio signal. We also show that recent advances in deep learning translate very well to the music recommendation setting, with deep convolutional neural networks significantly outperforming the traditional approach.

Paper Suggestions IV

Deep Neural Networks for YouTube Recommendations

Paul Covington, Jay Adams, Emre Sargin
Google
Mountain View, CA
{pcovington, jka, msargin}@google.com

ABSTRACT

YouTube represents one of the largest scale and most sophisticated industrial recommendation systems in existence. In this paper, we describe the system at a high level and focus on the dramatic performance improvements brought by deep learning. The paper is split according to the classic two-stage information retrieval dichotomy: first, we detail a deep candidate generation model and then describe a separate deep ranking model. We also provide practical lessons and insights derived from designing, iterating and maintaining a massive recommendation system with enormous user-facing impact.

Keywords

recommender system; deep learning; scalability

1. INTRODUCTION

YouTube is the world's largest platform for creating, sharing and discovering video content. YouTube recommendations are responsible for helping more than a billion users discover personalized content from an ever-growing corpus of videos. In this paper we will focus on the immense impact deep learning has recently had on the YouTube video recommendations system. Figure 1 illustrates the recommendations on the YouTube mobile app home.

Recommending YouTube videos is extremely challenging

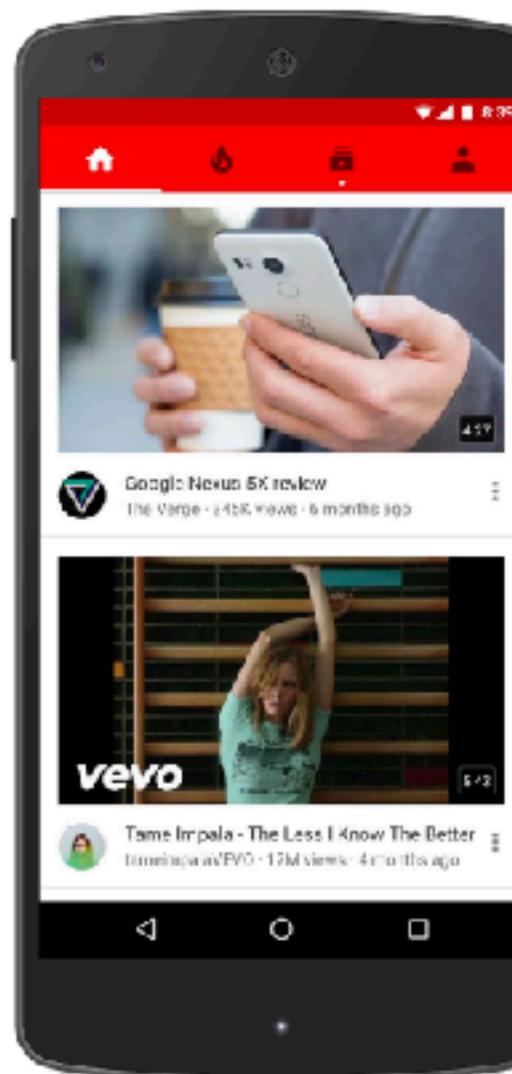


Figure 1: Recommendations displayed on YouTube mobile app home.

with well-established videos can be understood from an exploration/exploitation perspective.

Paper Suggestions V

Wide & Deep Learning for Recommender Systems

Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra,
Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, Rohan Anil,
Zakaria Haque, Lichan Hong, Vihan Jain, Xiaobing Liu, Hemal Shah

*
Google Inc.

Paper Suggestions VI

Collaborative Knowledge Base Embedding for Recommender Systems

Fuzheng Zhang[†], Nicholas Jing Yuan[†], Defu Lian[‡], Xing Xie[†], Wei-Ying Ma[†]

[†]Microsoft Research

[‡]Big Data Research Center, University of Electronic Science and Technology of China

{fuzzhang,nicholas.yuan,xingx,wyma}@microsoft.com,

dove.ustc@gmail.com

ABSTRACT

Among different recommendation techniques, collaborative filtering usually suffer from limited performance due to the sparsity of user-item interactions. To address the issues, auxiliary information is usually used to boost the performance. Due to the rapid collection of information on the web, the knowledge base provides heterogeneous information including both structured and unstructured data with different semantics, which can be consumed by various applications. In this paper, we investigate how to leverage the heterogeneous information in a knowledge base to improve the quality of recommender systems. First, by exploiting the knowledge base, we design three components to extract items' semantic representations from structural content, textual content and visual content, respectively. To be specific, we adopt a heterogeneous network embedding method, termed as TransR, to extract items' structural representations by considering the heterogeneity of both nodes and relationships. We apply stacked denoising auto-encoders and stacked convolutional auto-encoders, which are two types of deep learning based embedding techniques, to extract items' textual representations and visual representations, respectively. Finally, we propose our final integrated framework, which is termed as Collaborative Knowledge Base Embedding (CKE), to jointly learn the latent representations in collaborative filtering as well as items' semantic representations from the knowledge base. To evaluate the performance of each embedding component as well as the whole system, we conduct extensive experiments with two real-

filtering (CF) based methods, which make use of historical interactions or preferences, have made significant success [23]. However, CF methods usually suffer from limited performance when user-item interactions are very sparse, which is very common for scenarios such as online shopping where the item set is extremely large. In addition, CF methods can not recommend new items since these items have never received any feedbacks from users in the past. To tackle these problems, hybrid recommender systems, which combine collaborative filtering and auxiliary information such as item content, can usually achieve better recommendation results and have gained increasing popularity in recent years [2].

Over the past years, more and more semantic data are published following the Linked Data principles¹, by connecting various information from different topic domains such as people, books, music, movies and geographical locations in a unified global data space. These heterogeneous data, interlinked with each other, forms a huge information resource repository called knowledge base. Several typical knowledge bases have been constructed, including academic projects such as YAGO², NELL³, DBpedia⁴, and DeepDive⁵, as well as commercial projects, such as Microsoft's Satori⁶ and Google's Knowledge Graph⁷. Using the heterogeneous connected information from the knowledge base can help to develop insights on problems which are difficult to uncover with data from a single domain [6]. To date, information retrieval [9], community detection [25], sentiment analysis [4] - to name a few - are the noteworthy applications that successfully leverage the knowledge base.

Paper Suggestions VII

SESSION-BASED RECOMMENDATIONS WITH RECURRENT NEURAL NETWORKS

Balázs Hidasi *
Gravity R&D Inc.
Budapest, Hungary
balazs.hidasi@gravityrd.com

Alexandros Karatzoglou
Telefonica Research
Barcelona, Spain
alexk@tid.es

Linas Baltrunas †
Netflix
Los Gatos, CA, USA
lbaltrunas@netflix.com

Domonkos Tikk
Gravity R&D Inc.
Budapest, Hungary
domonkos.tikk@gravityrd.com

ABSTRACT

We apply recurrent neural networks (RNN) on a new domain, namely recommender systems. Real-life recommender systems often face the problem of having to base recommendations only on short session-based data (e.g. a small sportswear website) instead of long user histories (as in the case of Netflix). In this situation the frequently praised matrix factorization approaches are not accurate. This problem is usually overcome in practice by resorting to item-to-item recommendations, i.e. recommending similar items. We argue that by modeling the whole session, more accurate recommendations can be provided. We therefore propose an RNN-based approach for session-based recommendations. Our approach also considers practical aspects of the task and introduces several modifications to classic RNNs such as a ranking loss function that make it more viable for this specific problem. Experimental results on two data-sets show marked improvements over widely used approaches.

Paper Suggestions VIII

Graph Convolutional Matrix Completion

Rianne van den Berg
University of Amsterdam

Thomas N. Kipf
University of Amsterdam

Max Welling
University of Amsterdam, CIFAR¹

Abstract

We consider matrix completion for recommender systems from the point of view of link prediction on graphs. Interaction data such as movie ratings can be represented by a bipartite user-item graph with labeled edges denoting observed ratings. Building on recent progress in deep learning on graph-structured data, we propose a graph auto-encoder framework based on differentiable message passing on the bipartite interaction graph. Our model shows competitive performance on standard

in the form of node features. Predicting ratings then reduces to predicting labeled links in the bipartite user-item graph.

We propose graph convolutional matrix completion (GC-MC): a graph-based auto-encoder framework for matrix completion, which builds on recent progress in deep learning on graphs [2, 6, 19, 5, 15, 30, 14]. The auto-encoder produces latent features of user and item nodes through a form of message passing on the bipartite interaction graph. These latent user and item representations are used to reconstruct the rating links through a bilinear decoder.

Paper Suggestions IX

Deep Models of Interactions Across Sets

Jason Hartford^{*1} Devon R Graham^{*1} Kevin Leyton-Brown¹ Siamak Ravanbakhsh¹

Abstract

We use deep learning to model interactions across two or more sets of objects, such as user–movie ratings or protein–drug bindings. The canonical representation of such interactions is a matrix (or tensor) with an exchangeability property: the encoding’s meaning is not changed by permuting rows or columns. We argue that

$\langle n, m, x \rangle \in \mathbb{X}$. Learning our function corresponds to *matrix completion*: using patterns in \mathbb{X} to predict values for the remaining elements of X .

X is what we will call an *exchangeable matrix*: any row- and column-wise permutation of X represents the same set of ratings and hence the same matrix completion problem. Exchangeability has a long history in machine learning and statistics. For example, the common iid assumption implies

huge list here:

**[https://github.com/
hongleizhang/RSPapers](https://github.com/hongleizhang/RSPapers)**