

Unsupervised Rank Aggregation with Distance-Based Models

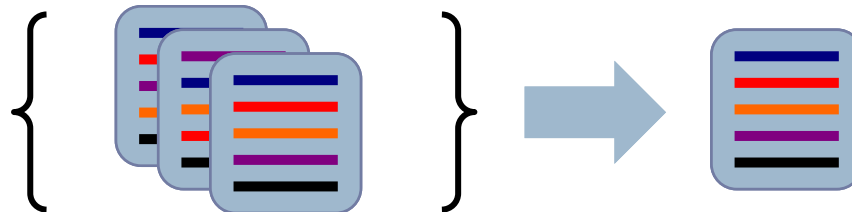
Alexandre Klementiev, Dan Roth, and Kevin Small
University of Illinois at Urbana-Champaign

Motivation

- ▶ Consider a panel of judges
 - ▶ Each (independently) generates preferences over items, i.e. (partial) rankings
- ▶ The need to meaningfully aggregate their rankings is a fundamental problem
 - ▶ Applications are plentiful in Information Retrieval and Natural Language Processing

Meta-search

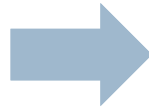
- ▶ *Meta-search*: combine results of multiple search engines into a single ranking
 - ▶ Sets of ranked pages are different across rankers
 - ▶ *Supervision is difficult to get*, often collected indirectly (e.g. clickthrough data)



Multilingual Named Entity Discovery

- ▶ *Named Entity Discovery [Klementiev& Roth, ACL 06]:* given a bilingual corpus one side of which is annotated with Named Entities, find their counterparts in the other

guimaraes



Candidate	r1	r2	r3	r4
гуимареша	1	1	3	3
муаммар	3	2 7	4	1 4
гимараешу	2	3	1	1
футбол	5 9	2	3 1	1 2
гамма	7	1 7	5 1	3 2
...				

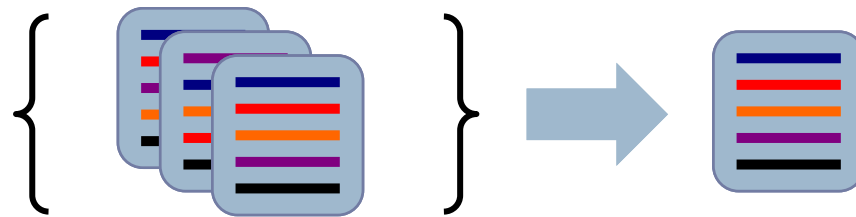
- ▶ NEs are often transliterated: rare ... *ion model*
- ▶ NEs tend to co-occur across languages: rank according to *temporal*
- ▶ NEs tend to co-occur in similar contexts: rank according to *contextual*
- ▶ NEs tend to co-occur in similar topics: rank according to *topic similarity*
- ▶ etc.

Outline

- ▶ Motivation
 - ▶ Problem Statement
 - ▶ Overview of our approach
 - ▶ Background
 - ▶ Mallows models
 - ▶ Extended Mallows models
 - ▶ Unsupervised Learning and Inference
 - ▶ Instantiations of the framework
 - ▶ Combining permutations
 - ▶ Combining top-k lists
 - ▶ Experiments
 - ▶ Conclusions and Current/Future work
- Introduction and background
- Our contribution

Problem

How can we combine (partial) object preferences from multiple judges into a joint ranking?



- ▶ In *IR*, many approaches (data fusion) aggregate rankings heuristically
 - ▶ Linear score/rank aggregation is frequently used
 - ▶ Assume domain knowledge is available
- ▶ *Supervised machine learning* techniques require labeled training data

Overview of Our Approach

- ▶ We propose a formal framework for *unsupervised* rank aggregation
 - ▶ Judges independently generate a (partial) ranking attempting to reproduce the true underlying ranking based on their level of expertise
 - ▶ We derive an EM-based algorithm treating the votes of individual judges and the true ranking as the observed and unobserved data, respectively
- ▶ We instantiate the framework for the cases of combining permutations and combining *top-k*lists

Concepts and Notation

- ▶ Permutation π over n objects $x_1 \dots x_n$
 - ▶ $\pi(i)$ is the rank assigned to object x_i
 - ▶ $\pi^{-1}(j)$ is the index of the object assigned to rank j
 - ▶ $e = \pi\pi^{-1} = \pi^{-1}\pi$ is the identity permutation
- ▶ Set S_n of all $n!$ permutations
- ▶ Distance $d : S_n \times S_n \rightarrow R_+$ between permutations
 - ▶ E.g. *Kendall's tau distance*: minimum number of adjacent transpositions needed to turn π into σ
- ▶ d is assumed to satisfy the *right invariance* property: does not depend on arbitrary re-labeling of the n objects
 - ▶ $d(\pi, \sigma) = d(\pi\pi^{-1}, \sigma\pi^{-1}) = d(e, \nu) = D(\nu)$. If ν is a r.v., so is $D = D(\nu)$

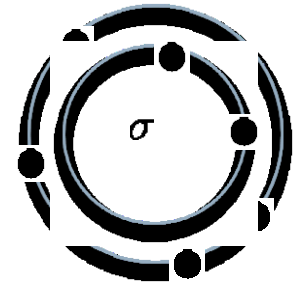
Background: Mallows Models

Uniform when
 $\theta = 0$

“Peaky” when
 $|\theta|$ is large

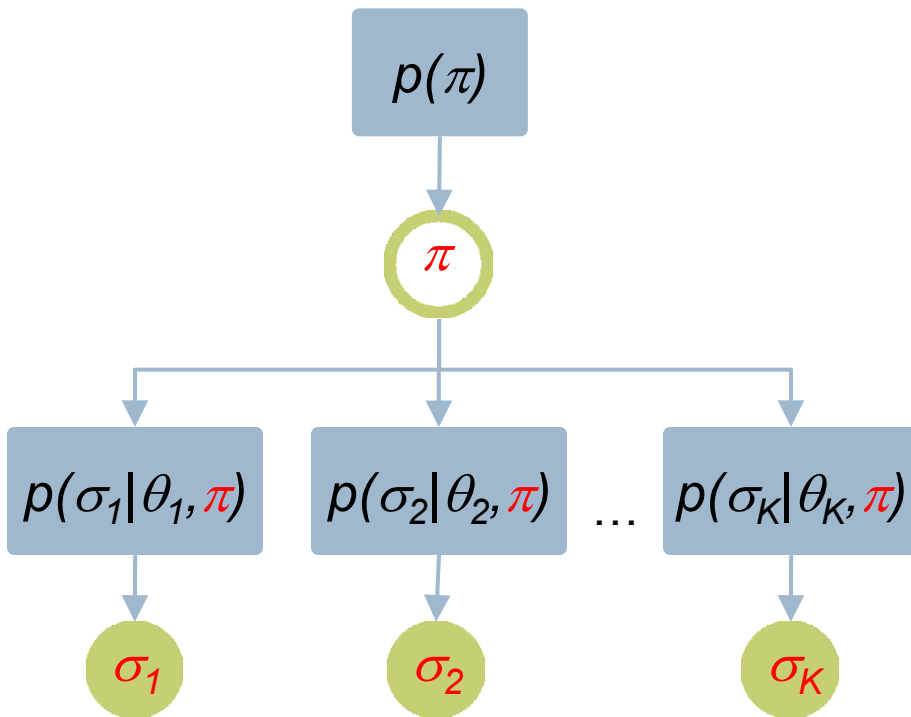
Expensive to
compute

where



- ▶ $\theta \in \mathbb{R}$, $\theta \leq 0$ is the *dispersion* parameter
- ▶ $\sigma \in \mathcal{S}_n$ *location* parameter
- ▶ $d(\cdot, \cdot)$ right-invariant, $sZ(\theta, \sigma)$ not depend on σ
- ▶ If D can be decompose $\epsilon D(\pi) = \sum_{i=1}^m V_i(\pi)$ where V_i are indep. r.v. $E_{\theta}(D)$ may be efficient to compute [Fligner and Verducci '86]

Generative Story for Aggregation



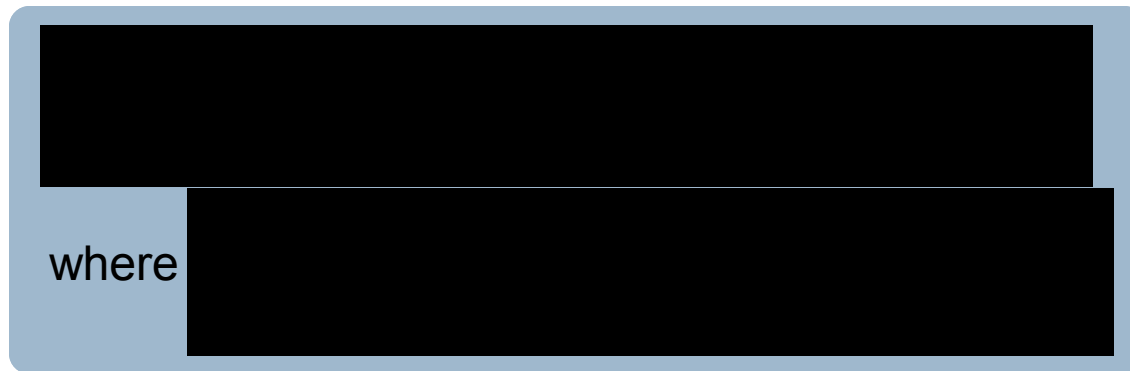
Generate the *true* π according to prior $p(\pi)$

Draw $\sigma_1 \dots \sigma_K$ independently from K Mallows models $p(\sigma_i | \theta_i, \pi)$, with the *same location parameter* π

$$p(\pi, \boldsymbol{\sigma} | \boldsymbol{\theta}) = p(\pi) \prod_{i=1}^K p(\sigma_i | \theta_i, \pi)$$

Background: Extended Mallows Models

The associated conditional model (when votes of K judges $\sigma \in \mathcal{S}_n^K$ are available) proposed in [Lebanon and Lafferty '02]:



Free parameters $\theta \in \mathbb{R}^K$, $\theta \leq \mathbf{0}$ represent the degree of expertise of individual judges.

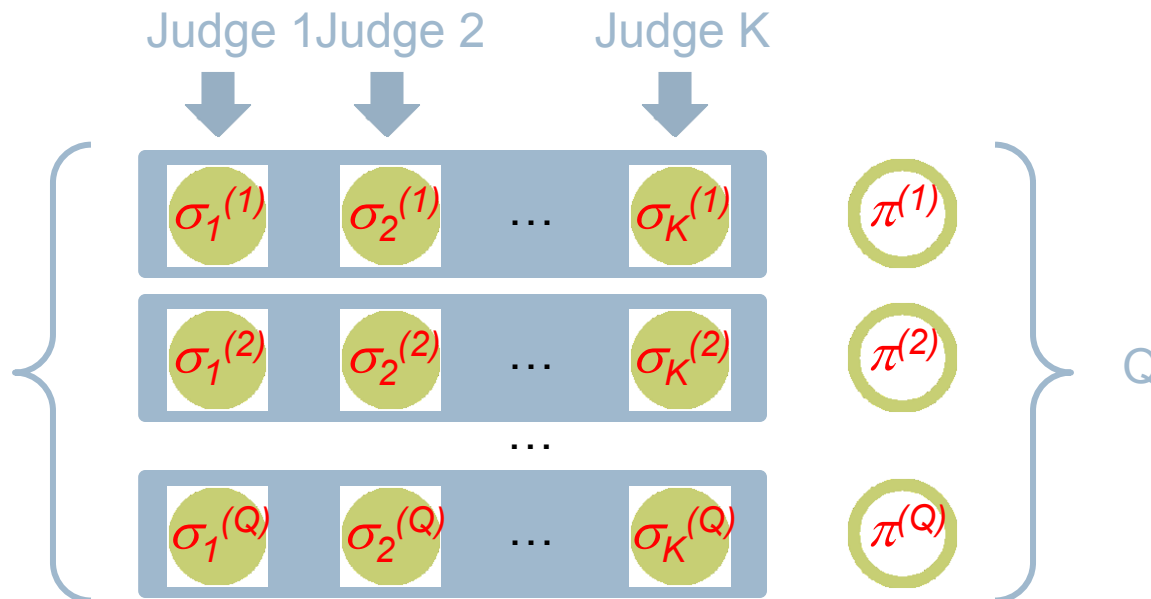
It is straightforward to generalize both models to partial rankings by constructing *appropriate distance functions*

Outline

- ▶ Motivation
 - ▶ Problem Statement
 - ▶ Overview of our approach
 - ▶ Background
 - ▶ Mallows models
 - ▶ Extended Mallows models
 - ▶ **Unsupervised Learning and Inference**
 - ▶ **Instantiations of the framework**
 - ▶ Combining permutations
 - ▶ Combining top-k lists
 - ▶ **Experiments**
 - ▶ **Conclusions and Current/Future work**
- Introduction and background
- Our contribution

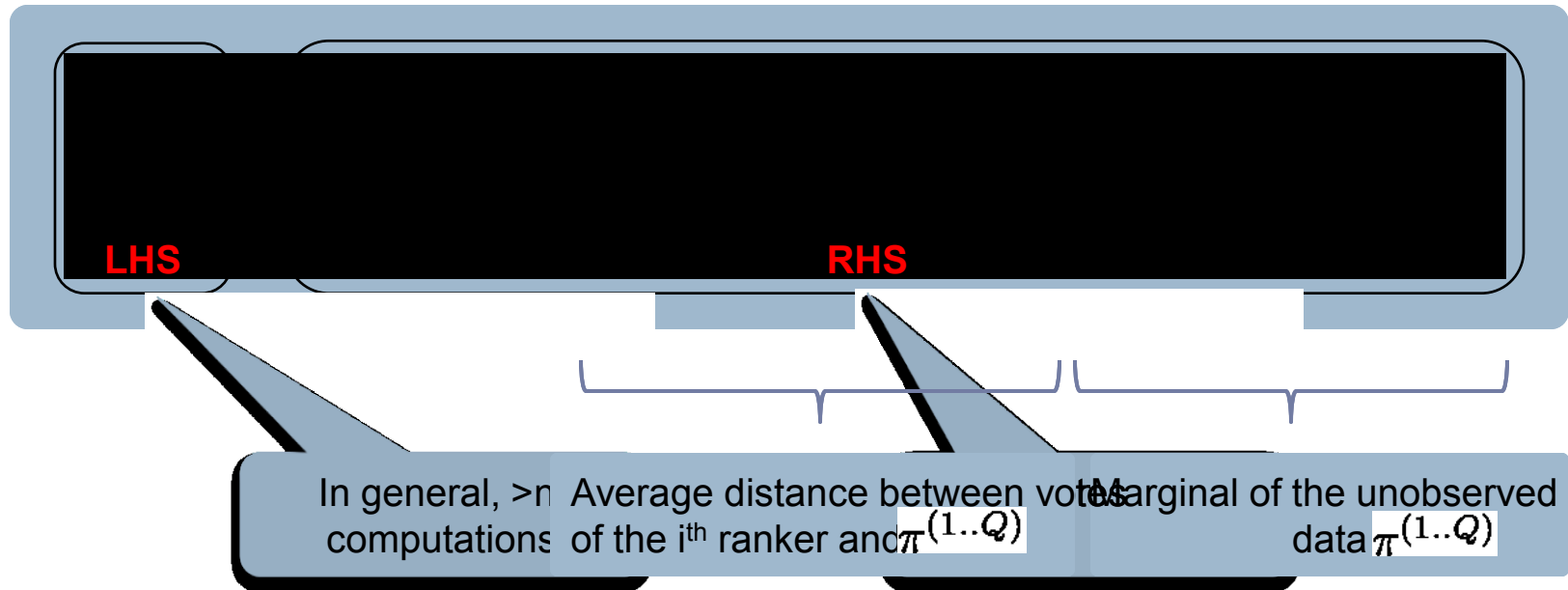
Our Approach

- ▶ We propose a formal framework for unsupervised rank aggregation based on the *extended Mallows model* formalism
- ▶ We derive an EM-based algorithm to estimate model parameters θ
 - ▶ Observed data: votes of individual judges
 - ▶ Unobserved data: true ranking



Learning

Denoting θ' to be the value of parameters from the previous iteration, the M step for the i^{th} ranker is:



Learning and Inference

LHS

RHS

Learning (estimating θ)

- ▶ For K constituent rankers, repeat:
 - ▶ Estimate the **RHS** given current parameter values θ
 - ▶ Sample with Metropolis-Hastings
 - ▶ Or use heuristics
 - ▶ Solve the **LHS** to update θ
 - ▶ Efficient estimation can be done using various types of distance functions

Depends on ranking type,
more about this later

Inference (computing the most likely ranking)

- ▶ Sample with Metropolis-Hastings or use heuristics as above

Instantiating the Framework

- ▶ We have *not committed* to a particular type of ranking
- ▶ In order to instantiate the framework:
 - ▶ Design a distance function appropriate for the setting
 - ▶ If a function is right invariant and decomposable [LHS] estimation can be done quickly (more about this later)
 - ▶ Design a sampling procedure for learning [RHS] and inference

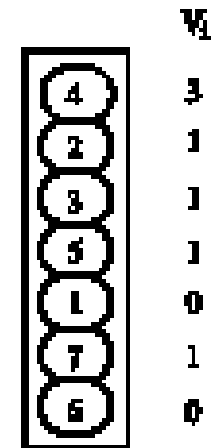
Case 1: Combining Permutations [LHS]

- ▶ Kendall tau distance D_K is the minimum number of adjacent transpositions needed to transform one permutation into another
- ▶ Can be decomposed into a sum of independent random variables:

$$D_K(\pi) = \sum_{i=1}^{n-1} V_i(\pi) \text{ where } V_i(\pi) = \sum_{j>i} I(\pi^{-1}(i) - \pi^{-1}(j))$$

- ▶ And the expected value can be shown to be:

$$E_{\theta}(D_K) = \frac{ne^{\theta}}{1 - e^{\theta}} - \sum_{j=1}^n \frac{je^{\theta j}}{1 - e^{\theta j}}$$



Monotonically decreasing,
can find θ with *line search*

Case 1: Combining Permutations [RHS]

Sampling from the base chain of random transpositions

- ▶ Start with a random permutation π

- ▶ If $\alpha = \frac{p(\pi'|\theta, \sigma)}{p(\pi|\theta, \sigma)} \geq 1$, use two objects π and π' forming a chain moves to π'
 - ▶ If $\alpha \geq 1$, chain moves to π' with probability $\frac{1}{\alpha}$
 - ▶ Else, chain moves to π with probability $\frac{1}{\alpha}$

- ▶ Note that we can compute distance incrementally, i.e. add the change due to a single transposition

- ▶ Convergence

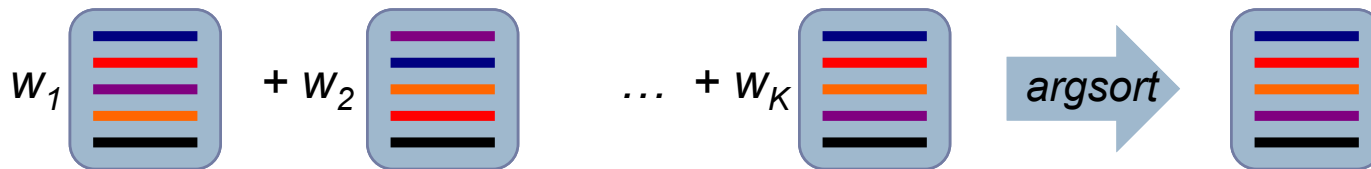
- ▶ $n \log(n)$ if d is Cayley's distance [Diaconis '98], likely similar for some others

▶ No convergence results for general case, but it works well in practice

Case 1: Combining Permutations [RHS]

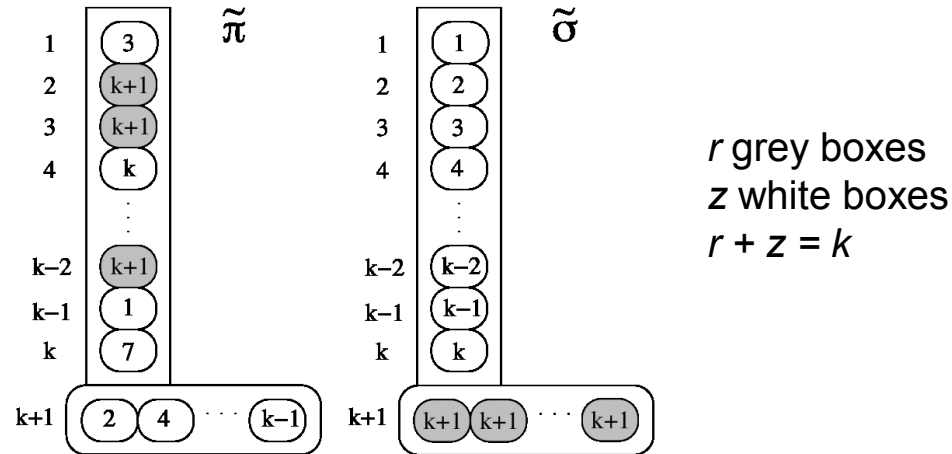
An alternative heuristic: weighted *Bordacount*, i.e.

- ▶ Linearly combine ranks of each object and argsort
- ▶ Model parameters θ represent relative expertise, so it makes sense to weigh ranks $e^{(-\theta_i)} V_i =$



Case 2: Combining Top-k [LHS]

- ▶ We extend Kendall tau to top-k



$$\tilde{D}_K(\tilde{\pi}) = \sum_{\substack{i=1 \\ \tilde{\pi}^{-1}(i) \notin Z}}^k \tilde{U}_i(\tilde{\pi}) + \frac{r(r+1)}{2} + \sum_{\substack{i=1 \\ \tilde{\pi}^{-1}(i) \in Z}}^k \tilde{V}_i(\tilde{\pi})$$

- Bring grey boxes to bottom
- Switch with objects in $(k+1)$
- Kendall's tau for the k elements

Case 2: Combining Top-k [LHS & RHS]

- ▶ R.v.'s \tilde{V}_i and \tilde{U}_i are independent, we can use the same trick to show that [LHS] is:

$$E_{\theta}(\tilde{D}_K) = \frac{ke^{\theta}}{1 - e^{\theta}} - \sum_{j=r+1}^k \frac{je^{j\theta}}{1 - e^{j\theta}} + \frac{r(r+1)}{2} - r(z+1) \frac{e^{\theta(z+1)}}{1 - e^{\theta(z+1)}}$$

- ▶ Also monotonically decreasing, can again use line search
- ▶ Both \tilde{D}_K and $E_{\theta}(\tilde{D}_K)$ reduce to Kendall tau results when same elements are ranked in both lists, i.e. $r = 0$
- ▶ Sampling / heuristics for [RHS] and inference are similar to the permutation case

Outline

- ▶ Motivation
 - ▶ Problem Statement
 - ▶ Overview of our approach
 - ▶ Background
 - ▶ Mallows models
 - ▶ Extended Mallows models
 - ▶ Unsupervised Learning and Inference
 - ▶ Instantiations of the framework
 - ▶ Combining permutations
 - ▶ Combining top-k lists
 - ▶ Experiments
 - ▶ Conclusions and Current/Future work
- Introduction and background
- Our contribution

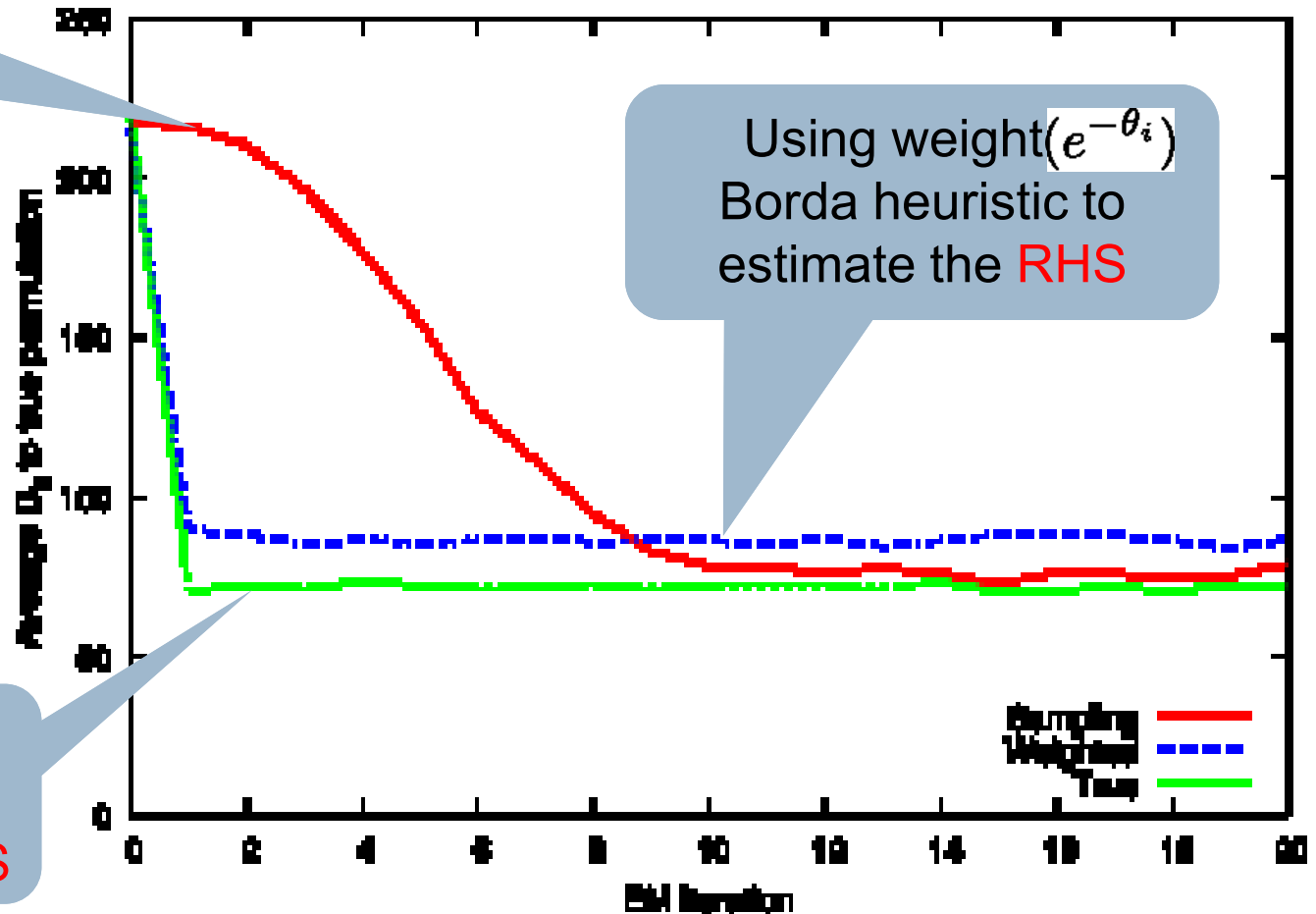
Exp. 1 Combining permutations

Using sampling to estimate the RHS

- Judges: $K = 10$
(Mallows models)
- Objects: $n = 30$
- $Q = 10$ sets of votes

Using weight ($e^{-\theta_i}$) Borda heuristic to estimate the RHS

Using true rankings to evaluate the RHS



Exp. 2 Meta-search dispersion parameters

- ▶ *Judges*: $K = 4$ search engines (S1, S2, S3, S4)
- ▶ *Documents*: Top $k = 100$
- ▶ *Queries*: $Q = 50$ queries

Define Mean Reciprocal Page Rank (MRPR): mean rank of the page containing the correct document

- ▶ Our model gets **0.92**

	$S1$	$S2$	$S3$	$S4$
θ	-0.065	0.0	-0.066	-0.049
$MRPR$	0.86	0.43	0.82	0.78

Model parameters correspond to ranker quality

Exp. 3 Top-k rankings: robustness to noise

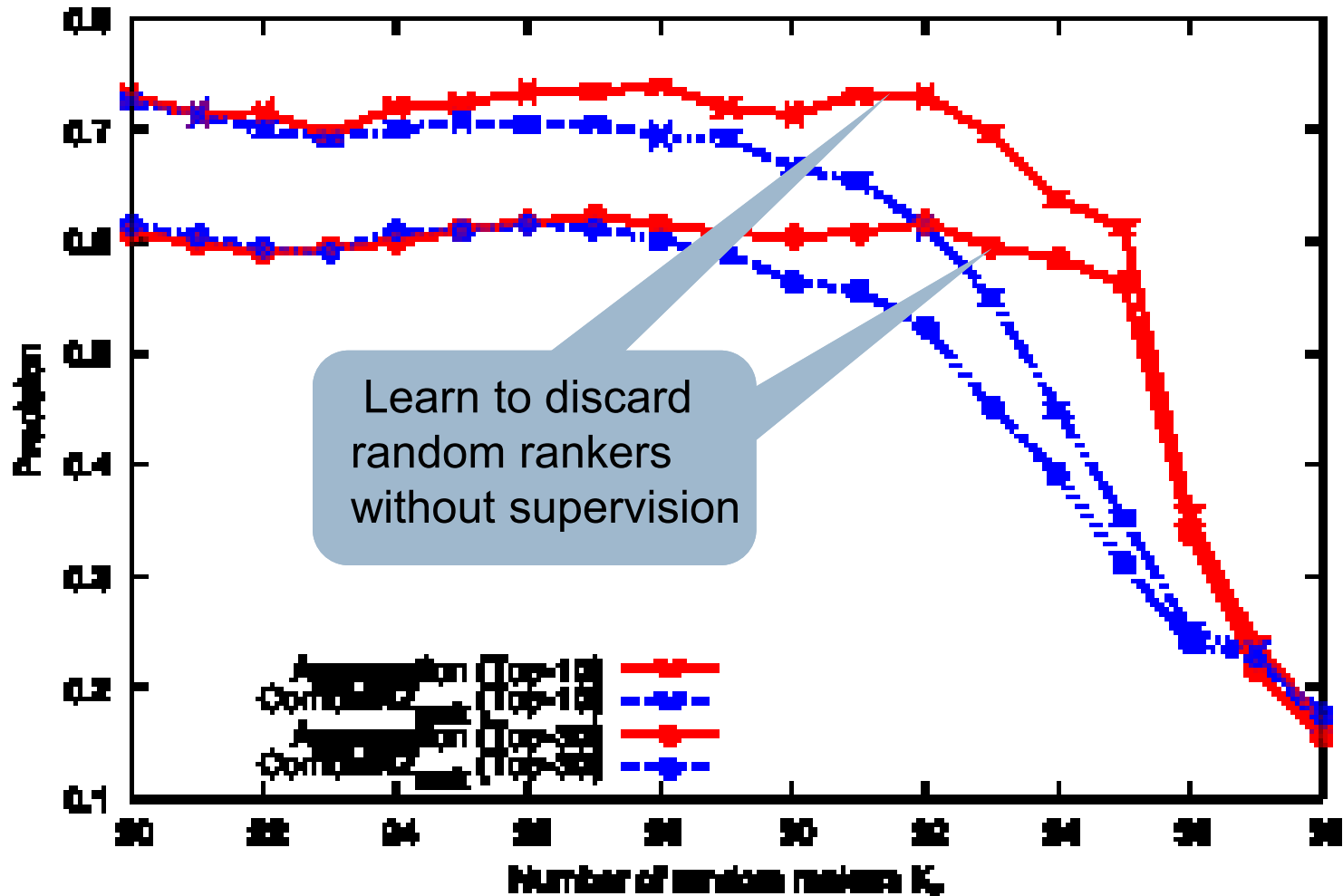
- ▶ *Judges*: $K = 38$ TREC-3 ad-hoc retrieval shared task participants
- ▶ *Documents*: Top $k = 100$ documents
- ▶ *Queries*: $Q = 50$ queries

Replaced $K_r \in [0, K]$ randomly chosen participants with random rankers. Baseline: rank objects according to score:

$$\text{CombMNZ}_{rank} = N_x \times \sum_{i=1}^K (k - r_i(x, q))$$

where $r_i(x, q)$ is the rank of x returned by i for query q , εN_x is the number of participants with x in top- k

Exp. 3 Top-k rankings: robustness to noise



Conclusions

- ▶ Propose a formal mathematical and algorithmic framework for aggregating (partial) rankings *without* supervision
 - ▶ Show that learning can be made efficient for decomposable distance functions
- ▶ Instantiate the framework for combining *permutations* and combining *top-k*lists
 - ▶ Introduce a novel distance function for *top-k*

Future work / work in progress

- ▶ Instantiate to other types of partial rankings
 - ▶ E.g. MT system aggregation: combine alignments
- ▶ Query / document type dependence: experts quality may depend on types of queries or objects being ranked
- ▶ Position dependence:
 - ▶ Right-invariant d , which is *position-dependent* (e.g. favors agreement at the top), need to be able to simplify the LHS.
 - ▶ [Flinger, Verducci '86] propose a multistage extension where objects are chosen going from top down position dependent θ . *Can we combine the ideas?*
- ▶ Domain Adaptation