

Ranking Big Data Sets using Rank Aggregation Techniques



Sarah Cohen-Boulakia



Laboratoire de Recherche en Informatique

CNRS UMR 8623, Université Paris-Sud

Université Paris-Saclay



General context

- Large set of data items can be obtained
 - as answer to a query, as produced by a tool...
- **Ranking** data items is crucial
- Ranking data may be difficult: **how** to rank?
- Alternative **ranking criteria** may be considered
 - (Quality-like) Reliability, Completeness...
- Alternative **methods** may rank the same set of data items
 - Google vs Yahoo vs ...
 - Classifiers...

Several rankings taken as input to produce one **consensus ranking** minimizing their disagreements

$$\pi_1 = [A, D, C, B]$$

$$\pi_2 = [B, A, D, C]$$

$$\pi_3 = [D, A, B, C]$$

$$\pi^* = [A, D, B, C]$$

→ Rank Aggregation Techniques

More formally...

- Rank aggregation is based on a *distance* to produce the closest ranking from a set of input rankings

► **The Kendall- τ $D(\pi, \sigma)$ distance**

#pairs of elements **inversed between** two rankings

$$D(\pi, \sigma) = \left| \left\{ (i, j) : i < j \wedge \left(\begin{array}{l} \pi[i] < \pi[j] \wedge \sigma[i] > \sigma[j] \\ \vee \pi[i] > \pi[j] \wedge \sigma[i] < \sigma[j] \end{array} \right) \right\} \right|$$

$$\begin{aligned} \pi_1 &:= [A, D, C, B] \\ \pi_2 &:= [B, A, D, C] \\ D(\pi_1, \pi_2) &= 1_{A>B} \\ &\quad + 1_{B>D} \\ &\quad + 1_{B>C} \\ &= 3 \end{aligned}$$

More formally...

- Rank aggregation is based on a **distance** to produce the closest ranking from a set of input rankings

- ▶ **The Kendall- τ $D(\pi, \sigma)$ distance**

#pairs of elements **inversed between** two rankings

$$D(\pi, \sigma) = \left| \left\{ (i, j) : i < j \wedge \left(\begin{array}{l} \pi[i] < \pi[j] \wedge \sigma[i] > \sigma[j] \\ \vee \pi[i] > \pi[j] \wedge \sigma[i] < \sigma[j] \end{array} \right) \right\} \right|$$

- ▶ **Kemeny Score**

$$S(\pi, \mathcal{P}) = \sum_{\sigma \in \mathcal{P}} D(\pi, \sigma)$$

- ▶ **Optimal Consensus (median)**

$$\forall \pi \in \mathcal{S}_n : S(\pi^*, \mathcal{P}) \leq S(\pi, \mathcal{P})$$

Complexity [Dwork et al 2001, Biedl et al. 2009, Bachmeier et al. 2017] **NP-Difficult** \rightarrow **Numerous heuristics**

$$\pi_1 := [A, D, C, B]$$

$$\pi_2 := [B, A, D, C]$$

$$D(\pi_1, \pi_2) = 1_{A>B} + 1_{B>D} + 1_{B>C} = 3$$

$$\mathcal{P} \begin{cases} \pi_1 = [A, D, C, B] \\ \pi_2 = [B, A, D, C] \\ \pi_3 = [D, A, B, C] \end{cases}$$

$$\pi^* = [A, D, B, C]$$

$$S(\pi^*, \mathcal{P}) = 1_{A>B@ \pi_2} + 1_{A>D@ \pi_3} + 1_{B>C@ \pi_1} + 1_{B>D@ \pi_2} = 4$$

Our expertise @LRI (Data Science)

- Comparison of 15+ algorithms (exact, approx, heuristics) able to provide consensus rankings
 - Rank-and-ties platform
- Design (or tuning) of efficient rank aggregation algorithms
- Considering alternative distances to compute consensus rankings taking into account user needs, context
 - Candidates to an election → one single position per candidate

$$\begin{aligned}
 \pi_1 &= [A, D, C, B] \\
 \pi_2 &= [B, A, D, C] \\
 \pi_3 &= [D, A, B, C]
 \end{aligned}
 \quad
 \pi^* = [A, D, B, C]$$

- Participant to sport competitions → *ex-aequo* should be allowed $\pi_1 = [A, \{D, C\}, B]$...
- Movies people like → not all the movies are ranked by anyone
 - $\pi_1 = [A, D, B], \pi_2 = [A, C, B]$...

Our expertise @LRI (Bioinformatics)

- Using rank aggregation techniques to automatically rank results obtained to a query and equivalent reformulations

<http://conqur-bio.lri.fr>

Genes associated with *Breast cancer*?

[G1, G2, G3, G4]

Genes associated with *Mamamilian*

Carcinoma? [G1, G3, G4, G12]

...

ConQuR-Bio
Consensus ranking with Query Reformulation for biological data from NCBI

Your query:

Species: Search deeper Show rank changes

Results

Rank	Name	Id	Official Full Name
1	BRCA2	(ID:675)	BRCA2, DNA repair associated
2	BRCA1	(ID:672)	BRCA1, DNA repair associated
3	TERT	(ID:7015)	telomerase reverse transcriptase
4	ESR1	(ID:2099)	estrogen receptor 1
5	CDKN2A	(ID:1029)	cyclin dependent kinase inhibitor 2A
6	CCND1	(ID:595)	cyclin D1
7	CHEK2	(ID:11200)	checkpoint kinase 2

Details

- ✓ Finding reformulations
- ✓ Running queries 14/14
- ✓ Computing a consensus ranking

Open with GeneValorization
[All these results](#), or only [the top 20](#).
[NCBI's results](#), or only [the top 20](#).

Conclusion

- Rank aggregation techniques are interesting
 - When it is **difficult to decide on which criteria to rank**
 - When you have **several rankings** and want to highlight their **common points**
- We have expertise on
 - **Design** new rank aggregation algorithms
 - **Help you choose and tune** the right rank aggregation algorithm
- Current collab. with **APHP P. Brousse** on Conquer-Bio (Leukemia)
 - 6 months of M2 by CDS 2.0 for **Pierre Andrieu** (now PhD stud., MNRT)
 - A very efficient and promising heuristic has been designed!
- We are **open to new collaborations**
 - On the domain science side (**new use cases**)
 - On the data science side (**evaluation of consensus rankings**)

THANKS!



Alain Denise
LRI



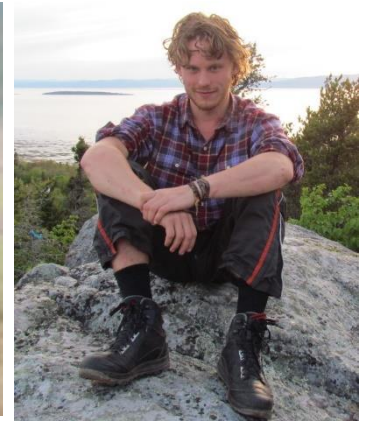
Bryan Brancotte
(now @Pasteur)



Pierre Andrieu
LRI



Adeline Pierrot
LRI



Robin Milosz
(Univ. Montréal &
LRI for 6 months)



Bastien Rance
APHP G. Pompidou



Ivan Sloma
APHP P. Brousse



Christophe Desterke
Inserm P. Brousse



Sylvie Hamel
(Univ. Montréal)

