# STAD Research in a Tree
## Learning from Preference Rankings

Antonio D'Ambrosio

# STAD research group

# STAD research in a tree



1994, 1997: Speeding-up growing trees
1997: Logistic classification trees
1998: Ternary trees
1999: Latent budget trees
2000: Multivariate trees
2006: Three-way trees
2016: Prediction trees for preference rankings
2017: Trees for functional predictors
2018: Trees for functional responses
2019: Recursive partitioning for ordinal responses

2015: Branch-and-Bound Solution Median Ranking Problem for Full Rankings;
2016: FAST Algorithms for Kemeny Problem for Full and Partial Rankings;
2017: Differential Evolution Algorithm for Rank Aggregation Problems;
2018, 2019, in progress: Avoiding Degenerate Solutions for Unfolding;
2018, 2019, progress: Projection Pursuit Clustering for Rank Data
2019: Non-Parametric Soft-Clustering Method for Preference Rankings;

2017: Visualization of Web-Sequence Rules

2004: New way to specify data editing
2009: Incremental imputation
2012: Boosted incremental imputation

2016: Parsimonious P-spline Based Clustering
2018, in progress: Fuzzy External Validation Criteria
2019, in progress: Depth-based medoid clustering
In progress: Co-clustering for directional data

Preference Rankings

Web mining

2017: Bibliometrix

Recursive Partitioning (TREES)

2016: Nested Stump Trees of Web Sequence Rules

Web Mining

Data Editing

2000: GAM-M
2002: GAM-MM
2015: V-curve for P-splines

Science mapping

2014: Concurvity in nonlinear and nonparametric regression models

Models diagnostic

Non-linear regression

2012: New Missing Data Imputation Paradigm within the Statistical Learning era
2018: Non-parametric spherical depths
2019: Median Constrained Bucket Order

Clustering

Key Issue

Key Issue

Data analysis

Finance;
Medicine;
Transportation;
Psychology;
Fraud detection;
Political sciences
.....

Methodological

Asset allocation;
Biology;
Structural engineering;
Psychometrics;
Social choice
...

Applied

Methodological

1990: ML approach to NSCA;
1992: Reduced rank models;
1994: Longitudinal NSCA;
1994: Latent Budget Analysis & Models;
2018: Unfolding;
In progress: Geometrical assessment evaluation criteria

Applied

Contribution

Supervised

Learning problem

Unsupervised

Contribution

The *optimal bucket order problem* (Gionis et al., 2006; Ukkonen et al., 2009; Kenkre et al., 2011; Aledo et al., 2017b) is a recent terminology for a old problem: dealing with rank aggregation by allowing tied ranking in the solution.

This problem was stated by Kemeny and Snell (1962) when defined the median ranking.

For long time the term 'preference rankings' has been a synonymous of permutations, tied rankings were interpreted as indifference declaration.

A bucket order is 'simply' a tied ranking.

# SDSS 2018

# Outline

# Preference rankings

# Preference rankings in a nutshell

Preference data are generally expressed by either *ratings* data or *rank* (or rankings, or preference rankings) data.

Both are data expressing individual's preference over a set of available alternatives.

Ratings data: *please assign a score in a range from 1 to 10 to the objects (sentences) A, B, C and D. The score 10 means "I completely agree". The score 1 means "I completely disagree".*

Rank data: *Please place the objects A, B, C and D in order in such a way that the resulting ordering reflects your preferences among these objects.*

## Type of rankings

When the subject assigns the integer values from 1 to $n$ to all the $n$ items we have a complete (or full) ranking.

| Item | A | B | C | D | E | F | G | H | I | L |
|------|---|---|---|---|---|---|---|---|----|---|
| Rank | 4 | 9 | 7 | 1 | 2 | 5 | 3 | 10 | 8 | 6 |

When a judge 'fails' to distinguish between two or more items and assigns to them the same integer number, we deal with tied (or weak) rankings

| Item | A | B | C | D | E | F | G | H | I | L |
|------|---|---|---|---|---|---|---|---|---|---|
| Rank | **3** | 7 | **5** | 1 | 2 | 4 | **3** | 7 | 6 | **5** |

We have a partial ranking (or incomplete rankings) when judges are asked to rank a subset of the entire set of objects (*pick k out of n*), or when there are some missingness in the ranked items

| Item | A | B | C | D | E | F | G | H | I | L |
|------|---|---|---|---|---|---|---|---|---|---|
| Rank | 1 |   |   | 4 | 2 |   |   |   | 3 |   |

# Geometry of rankings (1)

It is widely accepted that the geometrical space of preference rankings is the permutation polytope, which is the convex hull of a finite set of points in $\mathbb{R}^n$, in which the preference rankings are represented on its vertices (Thompson, 1993; Marden, 1996; Heiser, 2004; Heiser and D'Ambrosio, 2013; Alvo and Yu, 2014, ...).

# Just full rankings?

What about tied rankings? Just indifference declaration? Positive statement of agreement?

Nowadays dealing with tied rankings is the rule rather than an exception (ranking of Italian Universities, ranking of European Universities, ranking of World Universities, ranking of the Netflix series, ranking of the Amazon items,.....).

Times have changed, data have changed, sometimes the universe of the permutations is not enough.

Working with just *full rankings* can be a limitation in dealing with a lot of real problems ('*we consider the corresponding (tied) ranking positions as missing*',Jacques & Biernacki, C. (2014).

## Universe of rankings

The universe of rankings with $n$ items is equal to the ordered Bell number of $n$ elements

$$\mathcal{Z}^n = \sum_{b=1}^{n} b! \begin{Bmatrix} n \\ b \end{Bmatrix},$$

where $\begin{Bmatrix} n \\ b \end{Bmatrix} = \frac{1}{b!} \sum_{i=0}^{b} (-1)^i \binom{b}{i} (b-i)^n$ indicates the Stirling number of the second kind (the number of ways to partition a set of $n$ objects into $b$ non-empty subsets). These $b$ non-empty subsets are sometimes called *buckets*, so tied rankings are also known (in the computer science community) as bucket orders.
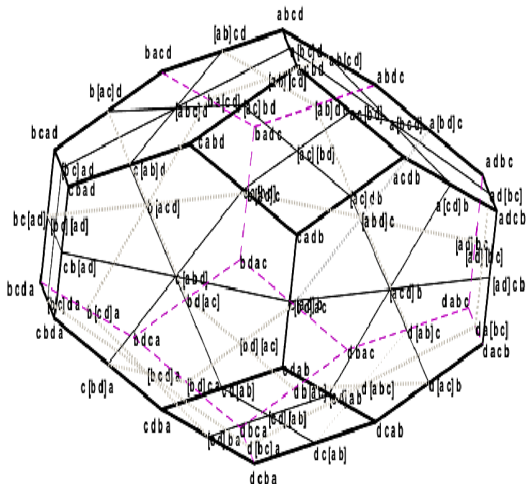
# Universe of rankings (2)

Cardinality of the universe of rankings containing ties for $n = 1, \ldots, 10$. The columns indicating the buckets (b) show the cardinality of the rankings of $n$ items constrained into $b$ buckets. Last column shows the universe of rankings with $n$ items

| n \ b | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | ... | $\mathcal{Z}^n$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | - | - | - | - | - | - | - | - | - | - | 1 |
| 2 | 1 | **2** | - | - | - | - | - | - | - | - | - | 3 |
| 3 | 1 | 6 | **6** | - | - | - | - | - | - | - | - | 13 |
| 4 | 1 | 14 | 36 | **24** | - | - | - | - | - | - | - | 75 |
| 5 | 1 | 30 | 150 | 240 | **120** | - | - | - | - | - | - | 541 |
| 6 | 1 | 62 | 540 | 1,560 | 1,800 | **720** | - | - | - | - | - | 4,683 |
| 7 | 1 | 126 | 1,806 | 8,400 | 16,800 | 15,120 | **5,040** | - | - | - | - | 47,293 |
| 8 | 1 | 254 | 5,796 | 40,824 | 126,000 | 191,520 | 141,120 | **40,320** | - | - | - | 545,835 |
| 9 | 1 | 510 | 18,150 | 186,480 | 834,120 | 1,905,120 | 2,328,480 | 1,451,520 | **362,880** | - | - | 7,087,261 |
| 10 | 1 | 1,022 | 955,980 | 818,520 | 5,103,000 | 16,435,440 | 29,635,200 | 30,240,000 | 16,329,600 | **3,628,800** | - | 102,247,563 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |

# Geometry of rankings (2)

Starting from the study of the permutation structure of partial (tied) rankings (with a pre-specified pattern of ties) made by Thompson (1993), Heiser and D'Ambrosio (2013) defined the following integrated graph of all full and partial (tied) rankings.

# Overview of statistical methods and models

Statistical methods and models for the analysis of preference rankings can be distinguished in (Marden, 1996):

- methods devoted to find the central ranking (de Borda, 1781; Condorcet, 1785; Mallows, 1957, ,...);

- methods based on badness-of-fit functions describing the multidimensional structure of rank data (Multidimensional Scaling, Unfolding, Vector model, Preference mapping,... Carroll 1972; Heiser and De Leeuw 1981; Meulman et al. 2004; Coombs 1950, 1964; Busing et al. 2005, 2010);

- methods based on probabilistic models, modeling either the ranking process or the population of rankers (distance-based models, multistage models,... Thurstone 1927; Bradley and Terry 1952; Mallows 1957; Fligner and Verducci 1986, 1988; Critchlow et al. 1991);

- methods that model the population of rankers assume heterogeneity among the judges with the goal to identify homogeneous sub-populations (mixtures of distance-based models, sorting insertion rank models, $K$-median cluster component analysis,... Croon 1989; Murphy and Martin 2003; Gormley and Murphy 2008a; Heiser and D'Ambrosio 2014; D'Ambrosio and Heiser 2018).

# Overview of statistical methods and models (with covariates)

Among the proposals that include covariates, the majority of them is based on:

- generalized linear models (Chapman and Staelin, 1982; Dittrich et al., 1998, 2000; Böckenholt, 2001; Gormley and Murphy, 2008b);
- recursive partitioning methods (D'Ambrosio, 2008; Cheng et al., 2009; Strobl et al., 2009; Lee and Yu, 2010; D'Ambrosio and Heiser, 2016; Plaia and Sciandra, 2017).

## Consensus Ranking

What is the common thread that combines **all** the methods and models dealing with preference rankings?

The detection of the so-called consensus ranking.

Given a series of judgments about a set of $n$ objects by a group of $m$ judges, what is the ranking that best represents the consensus opinion?

# Consensus ranking: a bit of history

It is:

- a very old problem (de Borda, 1781; Condorcet, 1785) ;

- that became a classical problem (Coombs, 1950; Black, 1958; Arrow, 1951; Goodman and Markowitz, 1952; Coombs, 1964; Davis et al., 1972; Bogart, 1973; Cook and Saipe, 1976; Cook and Seiford, 1978; Barthelemy and Monjardet, 1981; Beck and Lin, 1983; Barthélemy et al., 1989) ;

- remaining an actual problem (Emond and Mason, 2002; Meila et al., 2012; Cook et al., 2007; Biernacki and Jacques, 2013; D'Ambrosio et al., 2015; Amodio et al., 2016; Aledo et al., 2017a; D'Ambrosio et al., 2017; Yu and Xu, 2018) .

# Synonymous of consensus ranking

It has:

a lot of different names (Social choice problem, Consensus ranking problem, Rank aggregation problem, Kemeny problem, Median ranking problem, Kemeny aggregation problem, Preference learning problem.....),

also depending on the scientific field (Social sciences, Economics, Computer science, Statistics,...),

and the reference framework (ad hoc, distance-based, axiomatic, ...).

It is a NP-hard problem.

## Some distances for rankings: short list

In the framework of preference rankings, distance-based models and methods are largely used.

Several distance or dissimilarity measures have been defined: Spearman footrule, Spearman $\rho$, Kendall, Hulam, Hamming, Cayley, Kemeny,...

Each distance has some nice property, but which distance one should use? Is there some important desiderata? Is there some reference geometrical space?

| id | $\pi$ | $d(\pi_i, \pi_1)$ | | | | | | |
|----|-------|----------|----------|---------|--------|---------|------|--------|
| $\pi$ | | Footrule | Spearman | Kendall | Cayley | Hamming | Ulam | Kemeny |
| 1 | a b c d | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | a b d c | 2 | 2 | 1 | 1 | 2 | 1 | 2 |
| 3 | a c b d | 2 | 2 | 1 | 1 | 2 | 1 | 2 |
| 4 | a d b c | 4 | 6 | 2 | 2 | 3 | 1 | 4 |
| 5 | a c d b | 4 | 6 | 2 | 2 | 3 | 1 | 4 |
| 6 | a d c b | 4 | 8 | 3 | 1 | 2 | 2 | 6 |
| 7 | b a c d | 2 | 2 | 1 | 1 | 2 | 1 | 2 |
| 8 | b a d c | 4 | 4 | 2 | 2 | **4** | 2 | 4 |
| 9 | c a b d | 4 | 6 | 2 | 2 | 3 | 1 | 4 |
| 10 | d a b c | 6 | 12 | 3 | **3** | **4** | 1 | 6 |
| 11 | c a d b | 6 | 10 | 3 | **3** | **4** | 2 | 6 |
| 12 | d a c b | 6 | 14 | 4 | 2 | 3 | 2 | 8 |
| 13 | b c a d | 4 | 6 | 2 | 2 | 3 | 1 | 4 |
| 14 | b d a c | 6 | 10 | 3 | **3** | **4** | 2 | 6 |
| 15 | c b a d | 4 | 8 | 3 | 1 | 2 | 2 | 6 |
| 16 | d b a c | 6 | 14 | 4 | 2 | 3 | 2 | 8 |
| 17 | c d a b | **8** | 16 | 4 | 2 | **4** | 2 | 8 |
| 18 | d c a b | **8** | 18 | 5 | **3** | **4** | 2 | 10 |
| 19 | b c d a | 6 | 12 | 3 | **3** | **4** | 1 | 6 |
| 20 | b d c a | 6 | 14 | 4 | 2 | 3 | 2 | 8 |
| 21 | c b d a | 6 | 14 | 4 | 2 | 3 | 2 | 8 |
| 22 | d b c a | 6 | 18 | 5 | 1 | 2 | 2 | 10 |
| 23 | c d b a | **8** | 18 | 5 | **3** | **4** | 2 | 10 |
| 24 | d c b a | **8** | **20** | **6** | 2 | **4** | **3** | **12** |

# Distances and geometrical space

- Kendall and Kemeny are well defined in the permutation polytope.
- Spearman 'enter' in the polytope: only adjacent points are consistent with the polytope provided that the length of each edge equals $\sqrt{2}$.
- Cayley, Hamming and Ulam are not properly defined in the permutation polytope.
- Cayley does not reach the maximum distance between a ranking and its *reverse*.
- Hamming gets a lot of maximum distances.
- Kendall and Kemeny are equivalent for full rankings
- If ties are allowed, Kendall fails the triangular inequality and Spearman is sensitive to the *irrelevant alternatives*
- Kemeny is the unique distance defined in the generalized permutation polytope
- The Kemeny distance can be used in *Mallows* model only for full rankings. For tied rankings it is not possible: its exact distribution is not known (yet!)

# Kemeny's axiomatic framework

Let $A$ and $B$ be two rankings and let $d(A, B)$ be a distance between them (Kemeny, 1959; Kemeny and Snell, 1962):

- Axiom 1: $d(A, B)$ must be a metric;

- Axiom 2 : invariance:
  $d(A, B) = d(A', B')$, where $A'$ and $B'$ result from $A$ and $B$ respectively by the same permutation of the alternatives.

- Axiom 3: consistency in measurement:
  If two rankings $A$ and $B$ agree except for a set $S$ of $k$ elements, which is a segment of both, then $d(A, B)$ may be computed as if these $k$ objects were the only objects being ranked.

- Axiom 4: scaling: The minimum positive distance is 1.

# Kemeny distance

Suppose we have $n$ objects to be ranked. In defining his distance, Kemeny (1959) made use of the same matrix representation of rankings as was used earlier by Kendall (1948).

Let $a_{ij}$ ($b_{ij}$) be the generic element of the $n \times n$ squared preference matrix $A$ ($B$) called score matrix, with $i, j \in 1, \cdots, n$.

$a_{ij} = 1$ if the $i$th object is preferred to the $j$th object;
$a_{ij} = -1$ if the $j$th object is preferred to the $i$th object;
$a_{ij} = 0$ if the objects are tied.

The distance is defined as

$$d(A, B) = \frac{1}{2} \sum_{i}^{n} \sum_{j}^{n} |a_{ij} - b_{ij}|.$$

## Kemeny distance: properties

The Kemeny distance is the **unique** measure satisfying these axioms, working with **any** kind of ranking (full, partial, incomplete, tied), and naturally defined on the extended permutation polytope (Heiser and D'Ambrosio, 2013).

Except for the Kendall distance, any other (widely used) distance (e.g., Spearman's Footrule, Spearman $\rho$, Hamming, Cayley, Ulam) either is not defined in the polytope (do not preserve its geodesic nature) or assumes *strange* behaviors in dealing with tied rankings.

## Median ranking

Let $X_1, \ldots, X_m$ be a set of $m$ rankings of $n$ objects.
Kemeny and Snell (1962) defined the median ranking as that
ranking (or those rankings)

$$\hat{Y} = \arg\min_{Y \in \mathcal{Z}^n} \sum_{k=1}^{m} d(X_k, Y).$$

## $\tau_X$ rank correlation coefficient

Emond and Mason (2002) defined a new rank correlation coefficient, named *tau extension*, in this way:

$$\tau_X(A, B) = \frac{\sum_{i,j=1}^{n} a_{ij} b_{ij}}{n(n-1)},$$

where $a_{ij}$ and $b_{ij}$, $i, j = 1, \ldots, n$, are the elements of the score matrices of the rankings $A$ and $B$ slightly modified with respect to the original Kendall's formulation.

($a_{ij} = 1$ if the $i$th object is preferred to or is in a tie with the $j$th object).

## Median ranking: Emond and Mason's reformulation

They proved that

$$\tau_X(A, B) = 1 - 2\frac{d(A, B)}{n(n - 1)}.$$

The original Kemeny problem has been reformulated in this way:

$$\hat{Y} = \underset{Y \in \mathcal{Z}^n}{\arg\max} \frac{\sum_{k=1}^m w_k(\sum_{i,j=1}^n x_{ij}^{(k)} y_{ij})}{n(n - 1) \sum_{k=1}^m w_k} = \underset{Y \in \mathcal{Z}^n}{\arg\max} \sum_{i,j=1}^n c_{ij} y_{ij}, \text{ where}$$

- $w_k$ is a weight associated to the $k$-th ranking,
- $x_{ij}^{(1)}, \ldots, x_{ij}^{(m)}$ is the set of $m$ modified score matrices associated to $m$ rankings,
- $c_{ij} = \sum_{k=1}^m w_k x_{ij}^{(k)}$,
- $y_{ij}$ represents the elements of the *modified* score matrix associated to the ranking $Y$.

# Rank aggregation problem: STAD contribution 1

### Two algorithms for finding optimal solutions of the Kemeny rank aggregation problem for full rankings

Antonio D'Ambrosio*[a], Sonia Amodio[b], and Carmela Iorio[a]

[a]Department of Economics and Statistics, University of Naples Federico II, Naples, Italy
[b]Department of Medical Statistics, Leiden University Medical Centre, Leiden, The Netherlands

The analysis of ranking data has recently received increasing attention in many fields (i.e. political sciences, computer sciences, social sciences, medical sciences, etc.). Typically when dealing with preference rankings one of the main issue is to find a ranking that best represents the set of input rankings. Among several measures of agreement proposed in the literature, the Kendall distance is probably the most known. We propose a branch-and-bound algorithm to find the solution(s) even when we take into account a relatively large number of objects to be ranked. We also propose a heuristic variant of the branch-and-bound algorithm useful when the number of objects to rank is particularly high. We show how the solution(s) achieved by the algorithm can be employed in different analysis of rank data such as Mallow's $\phi$ model, mixtures of distance-based models, cluster analysis and so on.

keywords: Consensus ranking, Branch and bound, Mallows-$\phi$ model, exponential models.

- Branch-and-bound algorithm for full rankings
- Connection with Mallows Model
- One-to-one correspondence $\tau_a$ rcc with spread parameter $\lambda$

# Rank aggregation problem: STAD contribution 2

Decision Support

### Accurate algorithms for identifying the median ranking when dealing with weak and partial rankings under the Kemeny axiomatic approach

S. Amodio[a,1], A. D'Ambrosio[b,*], R. Siciliano[b]

[a] Department of Economics and Statistics, University of Naples Federico II, Italy
[b] Department of Industrial Engineering, University of Naples Federico II, Italy

ABSTRACT

Preference rankings virtually appear in all fields of science (political sciences, behavioral sciences, machine learning, decision making and so on). The well-known social choice problem consists in trying to find a reasonable procedure to use the aggregate preferences or rankings expressed by subjects to reach a collective decision. This turns out to be equivalent to estimate the consensus (central) ranking from data and it is known to be a NP-hard problem. A useful solution has been proposed by Emond and Mason in 2002 through the Branch-and-Bound algorithm (BB) within the Kemeny and Snell axiomatic framework. As a matter of fact, BB is a time demanding procedure when the complexity of the problem becomes untractable, i.e. a large number of objects, with weak and partial rankings, in presence of a low degree of consensus. As an alternative, we propose an accurate heuristic algorithm called FAST that finds at least one of the consensus ranking solutions found by BB saving a lot of computational time. In addition, we show that the building block of FAST is an algorithm called QUICK that finds already one of the BB solutions so that it can be fruitfully considered to speed up even more the overall searching procedure of the median ranking when the number of objects is low. Simulation studies and applications on real data allows to show the accuracy and the computational efficiency of our proposal.

- QUICK accurate algorithm for median ranking problem
- FAST solution for problems with large number of objects

# Rank aggregation problem: STAD contribution 3

### A differential evolution algorithm for finding the median ranking under the Kemeny axiomatic approach

Antonio D'Ambrosio[a,\*], Giulio Mazzeo[b], Carmela Iorio[c], Roberta Siciliano[c]

[a] Department of Economics and Statistics, University of Naples Federico II, Italy
[b] Hewlett Packard Enterprise, HPE Tech Partners India S.r.l., Italy
[c] Department of Industrial Engineering, University of Naples Federico II, Italy

ABSTRACT

In recent years the analysis of preference rankings has become an increasingly important topic. One of the most important tasks in dealing with preference rankings is the identification of the median ranking, namely that ranking that best represents the preferences of a population of judges. This task is known with several alternative names, such as rank aggregation problem, consensus ranking problem, social choice problem. In this paper we propose a Differential Evolution algorithm for the Consensus Ranking detection (DECoR) within the Kemeny's axiomatic framework. The algorithm works with full, partial and incomplete rankings. A simulation study shows that our proposal is particularly feasible when working with a very large number of objects to be ranked, because it is accurate and also faster than other proposals. Some applications on real data sets show the practical utility of our proposal in helping the users in taking decisions.

- Differential evolution proposal for discrete optimization problem
- Accurate solution for 'intractable" problems in a *reasonable* computing time

# STAD contribution to supervised learning for preference learning

Psychometric Society   CrossMark

A RECURSIVE PARTITIONING METHOD FOR THE PREDICTION OF PREFERENCE
RANKINGS BASED UPON KEMENY DISTANCES

ANTONIO D'AMBROSIO

UNIVERSITY OF NAPLES FEDERICO II

WILLEM J. HEISER

LEIDEN UNIVERSITY

Preference rankings usually depend on the characteristics of both the individuals judging a set of objects and the objects being judged. This topic has been handled in the literature with log-linear representations of the generalized Bradley-Terry model and, recently, with distance-based tree models for rankings. A limitation of these approaches is that they only work with full rankings or with a pre-specified pattern governing the presence of ties, and/or they are based on quite strict distributional assumptions. To overcome these limitations, we propose a new prediction tree method for ranking data that is totally distribution-free. It combines Kemeny's axiomatic approach to define a unique distance between rankings with the CART approach to find a stable prediction tree. Furthermore, our method is not limited by any particular design of the pattern of ties. The method is evaluated in an extensive full-factorial Monte Carlo study with a new simulation design.

Key words: prediction trees, kemeny distance, preference rankings, consensus ranking.

- Prediction trees for *any* kind of rankings
- New general simulation settings for *any* kind of tree-based methods
- It works with several sampling distributions
- Better, it works with real data

# STAD contribution to unsupervised learning for preference learning

CrossMark

## A distribution-free soft-clustering method for preference rankings

Antonio D'Ambrosio[1] · Willem J. Heiser[2]

**Abstract**

Typically, ranking data consist of a set of individuals, or judges, who have ordered a set of items—or objects—according to their overall preference or some pre-specified criterion. When each judge has expressed his or her preferences according to his own best judgment, such data are characterized by systematic individual differences. In the literature, several approaches have been proposed to decompose heterogeneous populations of judges into a defined number of homogeneous groups. Often, these approaches work by assuming that the ranking process is governed by some distance-based probability models. We use the flexible class of methods proposed by Ben-Israel and Iyigun, which consists in a probabilistic distance clustering approach, and define the disparity between a ranking and the center of a cluster as the Kemeny distance. This class of methods allows for probabilistic allocation of cases to classes, thus being a form of soft or fuzzy, clustering. The allocation probability is unequivocally related to the chosen distance measure.

**Keywords** Preference rankings · Soft clustering · Kemeny distance

- Probabilistic clustering for preference data
- Distribution free
- It works with several sampling distributions
- Better, it works with real data

# Optimal bucket order problem

The so-called *optimal bucket order problem* (OBOP) (Gionis et al., 2006; Ukkonen et al., 2009; Kenkre et al., 2011; Aledo et al., 2017b), namely dealing with rank aggregation while allowing ties in the solution, is in fact a recent terminology for the problem already stated by Kemeny and Snell (1962) when defined the median ranking.

'The optimal bucket order problem consists in obtaining a complete consensus ranking (ties are allowed) from a matrix of preferences...' (Aledo et al., 2018);

'...the problem is known as the Kemeny ranking problem (...) Both problems have in common that the solution is a permutation (i.e. a complete ranking without ties) defined over all the items' (Aledo et al., 2017b);

'We address the question of finding a bucket order for a set of items...' (Gionis et al., 2006);

...

## Optimal bucket orders (cont'd)

Within the Kemeny's axiomatic approach, both exact (Emond and Mason, 2002) and accurate heuristic algorithms (Amodio et al., 2016; D'Ambrosio et al., 2017) have been proposed. These algorithms, no matter about the nature of the rankings in input, search the best solution in $\mathcal{Z}^n$.

Other distance-based axiomatic frameworks allow for tied rankings as a *consensus ranking* solution (Cook et al., 1986, 1997)

# Median constrained bucket order

New *concept* (D'Ambrosio, 2017; D'Ambrosio et al, 2019):

let $X^{(1)}, \ldots, X^{(k)}$ be a set of rankings of $n$ items each of them bearing a weight $w_h$, with $\sum_{h=1}^{k} w_h = m$.

The median constrained bucket order is that ranking (or those rankings) $\hat{Y}$ for which

$$\hat{Y} = \arg\min_{Y \in \mathcal{Z}^{n \setminus b}} \sum_{h=1}^{k} w_h d(X^{(h)}, Y) = \arg\max_{Y \in \mathcal{Z}^{n \setminus b}} \frac{\sum_{i,j=1}^{n} c_{ij} y_{ij}}{m(n(n-1))},$$

where $\mathcal{Z}^{n \setminus b}$ is the subset of $\mathcal{Z}^n$ in which there are *exactly b* buckets.

# Rewind: Universe of rankings

Cardinality of the universe of rankings containing ties for $n = 1, \ldots, 10$. The columns indicating the buckets (b) show the cardinality of the rankings of $n$ items constrained into $b$ buckets. Last column shows the universe of rankings with $n$ items

| n \ b | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | ... | $\mathcal{Z}^n$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | - | - | - | - | - | - | - | - | - | - | 1 |
| 2 | 1 | 2 | - | - | - | - | - | - | - | - | - | 3 |
| 3 | 1 | 6 | 6 | - | - | - | - | - | - | - | - | 13 |
| 4 | 1 | 14 | 36 | 24 | - | - | - | - | - | - | - | 75 |
| 5 | 1 | 30 | 150 | 240 | 120 | - | - | - | - | - | - | 541 |
| 6 | 1 | 62 | 540 | 1,560 | 1,800 | 720 | - | - | - | - | - | 4,683 |
| 7 | 1 | 126 | 1,806 | 8,400 | 16,800 | 15,120 | 5,040 | - | - | - | - | 47,293 |
| 8 | 1 | 254 | 5,796 | 40,824 | 126,000 | 191,520 | 141,120 | 40,320 | - | - | - | 545,835 |
| 9 | 1 | 510 | 18,150 | 186,480 | 834,120 | 1,905,120 | 2,328,480 | 1,451,520 | 362,880 | - | - | 7,087,261 |
| 10 | 1 | 1,022 | 955,980 | 818,520 | 5,103,000 | 16,435,440 | 29,635,200 | 30,240,000 | 16,329,600 | 3,628,800 | - | 102,247,563 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |

# Median constrained bucket order

Why we are interested in such a constrained solution?

For example, according to the Bordeaux Official Wine Classification, wines are ranked in quality from first to fifth growths (Premier Cru, ..., Cinquieme Cru). In that wine tasting experiment the final solution is requested to be constrained into five buckets.

We were inspired by a possible *solution to a real problem*, by following the -too often forgotten- scheme according to which any **real problem should (must) be translated into a statistical problem, and the solution to the latter problem can help us to give a possible solution to the real one**.

# Triage prioritization example

An experiment was conducted in an Emergency Department (ED) of two popular Hospitals in Naples regarding the so-called triage, namely the admission phase to the ED.

A sample of 18 nurses for the Hospital named $\alpha$ and a sample of 35 nurses for the Hospital named $\beta$ had to place in order $n = 25$ cases according to their severity into $b = 4$ ordered categories: red ($R$), yellow ($Y$), green ($G$) and white ($W$). We assume that the cases can be ordered in terms of severity in this way: $R \succ Y \succ G \succ W$.

# Triage prioritization example (cont'd)

The 25 cases are the same for both Hospitals.

This experiment is equivalent to asking a set of $m$ judges to rank $n$ items allowing only $b$ different buckets, with $1 < b < n$.

# Triage prioritization example (cont'd)

The median constrained bucket order for Hospital $\alpha$ ($\tau_X = 0.6865$) is

[3 24] [1 5 6 7 10 15 16 20] [8 9 11 12 14 17 19 21 22 25] [2 4 13 18 23].

The median constrained bucket order for Hospital $\beta$ ($\tau_X = 0.6903$) is

[3 24] [1 5 7 10 16 21] [2 6 8 9 11 12 14 15 17 19 20 22 25] [4 13 18 23].

The buckets correspond to the coding $R$, $Y$, $G$ and $W$ respectively. The numbers correspond to the ID of each single patient.

## Triage prioritization example (cont'd)

After the experiment, a supervisor revealed the 'true' coding for each case, which is:

[3 24] [1 5 6 7 10 12 15 16 20] [8 9 11 14 17 19 21 22 25] [2 4 13 18 23].

The agreement between the true bucket order and the median constrained bucket orders is clear for Hospital $\alpha$ ($\tau_X = 0.917$), showing a good decision process of the nurses.

The same measure for the Hospital $\beta$ is equal to 0.697, showing a less good global decision process.

## Triage prioritization example (cont'd)

We can statistically check the equality of the median constrained bucket order by using the $R^2$ statistic as described in Marden (1996, Chapter 4, pag. 102)

$$R^2 = 1 - \frac{\sum_{l=1}^{L} \sum_{i=1}^{m^{(l)}} d(X^{(li)} \hat{Y}^{(l)})}{\sum_{l=1}^{L} \sum_{i=1}^{m^{(l)}} d(X^{(li)} \hat{Y})},$$

where $L$ and $m^{(l)}$ are the groups and the sample size within each group, $X^{(li)}$ is the $i$-th ranking in the $l$-th group, $\hat{Y}^{(l)}$ and $\hat{Y}$ are the median constrained bucket order for the $l$-th group and for the entire sample respectively.

If the bucket orders in the two samples are equal then $R^2 = 0$, which constitutes the null hypothesis of the test.

## Triage prioritization example (cont'd)

In our case $R^2 = 0.0477$ (even if the theoretical maximum value of $R^2$ is equal to one, practically it often achieves values close to zero. Marden, 1996).
The test has been performed by computing a randomized p-value with 1,000 replications (Feigin and Cohen, 1978; Marden, 1996), which resulted to be less than 0.001.

Nurses in Hospital $\beta$ need a more 'general' training phase than the ones working in Hospital $\alpha$.

This example shows the usefulness of the novel concept of constraining the median ranking to be expressed with a pre-specified number of buckets.

# Algorithmic details

- Branch-and-bound: *cut branches that generate rankings with more than b buckets. Cut branches whose penalty is larger than the incremental penalty if there are less than b buckets.*
- QUICK: *store rankings that have exactly b buckets. Discharge rankings with penalty larger than incremental penalty.*
- DECoR: *restrict the searching space and use the bounded-closest-integer approach instead of hierarchical approach.*

## Concluding remarks I

The median constrained bucket order problem is a **new concept**.

It *can* be tackled under several axiomatic frameworks, but

distance-based approaches to rank aggregation problems *must* take in account to deal with (a lot) of ties

# Concluding remarks II

It can be used *only* when there is a good reason for searching the solution in a restricted space (see triage prioritization data set, there are other -not shown- cases, such as the study of priorities for students with disabilities)

## Concluding remarks III

We propose both branch-and-bound and differential evolution
solutions, modifying the algorithms proposed by Emond and Mason
(2002), Amodio et al. (2016) and D'Ambrosio et al. (2017).

Any other proposal dealing with tied rankings can be 'adjusted' to
return a median constrained bucket order.

Algorithms can change, the idea remains.

Thank you

and thanks for being
still awake!!

**ORIGINAL PAPER**

CrossMark

## Median constrained bucket order rank aggregation

Antonio D'Ambrosio[1] · Carmela Iorio[2] · Michele Staiano[2] ·
Roberta Siciliano[2]

**Abstract**
The rank aggregation problem can be summarized as the problem of aggregating individual preferences expressed by a set of judges to obtain a ranking that represents the best synthesis of their choices. Several approaches for handling this problem have been proposed and are generally linked with either axiomatic frameworks or alternative strategies. In this paper, we present a new definition of median ranking and frame it within the Kemeny's axiomatic framework. Moreover, we show the usefulness of our approach in a practical case about triage prioritization.

**Keywords** Tied rankings · Median ranking · Kemeny distance · Triage prioritization