Ranking Median Regression: Learning to Order through Local Consensus

Anna Korba* Stephan Clémençon* Eric Sibony[†]

* Telecom ParisTech, † Shift Technology

Journée de la Chaire April 19 2018

Outline

- 1. Ranking Regression
- 2. Background and Results on Ranking Aggregation
- 3. Risk Minimization for Ranking (Median) Regression
- 4. Algorithms Local Median Methods

Outline

Ranking Regression

Background and Results on Ranking Aggregation

Risk Minimization for Ranking (Median) Regression

Algorithms - Local Median Methods

Ranking Regression

Consider:

- A set of *n* items: $[n] = \{1, ..., n\}$ (Ex: $\{1, 2, 3, 4\}$)
- A individual expresses her preferences as (full) ranking, i.e a strict order ≻ over n :

$$a_1 \succ a_2 \succ \cdots \succ a_n \quad (\mathsf{Ex:} \ 2 \succ 1 \succ 3 \succ 4)$$

• Also seen as the permutation σ that maps an item to its rank:

 $a_1 \succ \cdots \succ a_n \quad \Leftrightarrow \quad \sigma \in \mathfrak{S}_n \text{ such that } \sigma(a_i) = i$

Ex: $\sigma(2) = 1, \sigma(1) = 2, \dots, \mathfrak{S}_n$: set of permutations of [n], the symmetric group.

Ranking Regression

Consider:

- A set of *n* items: $[n] = \{1, ..., n\}$ (Ex: $\{1, 2, 3, 4\}$)
- A individual expresses her preferences as (full) ranking, i.e a strict order ≻ over n :

 $a_1 \succ a_2 \succ \cdots \succ a_n \quad (\mathsf{Ex:} \ 2 \succ 1 \succ 3 \succ 4)$

• Also seen as the permutation σ that maps an item to its rank:

 $a_1 \succ \cdots \succ a_n \quad \Leftrightarrow \quad \sigma \in \mathfrak{S}_n \text{ such that } \sigma(a_i) = i$

Ex: $\sigma(2) = 1, \sigma(1) = 2, \dots \mathfrak{S}_n$: set of permutations of [n], the symmetric group.

Problem: Given a vector X (e.g, the characteristics of an individual), the goal is to predict (her preferences) as a random permutation Σ in \mathfrak{S}_n .

Ranking Regression

Consider:

- A set of *n* items: $[n] = \{1, ..., n\}$ (Ex: $\{1, 2, 3, 4\}$)
- A individual expresses her preferences as (full) ranking, i.e a strict order ≻ over n :

 $a_1 \succ a_2 \succ \cdots \succ a_n \quad (\mathsf{Ex:} \ 2 \succ 1 \succ 3 \succ 4)$

• Also seen as the permutation σ that maps an item to its rank:

 $a_1 \succ \cdots \succ a_n \quad \Leftrightarrow \quad \sigma \in \mathfrak{S}_n \text{ such that } \sigma(a_i) = i$

Ex: $\sigma(2) = 1, \sigma(1) = 2, \dots, \mathfrak{S}_n$: set of permutations of [n], the symmetric group.

Problem: Given a vector X (e.g, the characteristics of an individual), the goal is to predict (her preferences) as a random permutation Σ in \mathfrak{S}_n .



Related Work

- ► Has been referred to as **label ranking** in the literature [Tsoumakas et al., 2009], [Vembu and Gärtner, 2010]
- Related to multiclass and multilabel classification
- A lot of applications, e.g : document categorization, meta-learning
 - rank a set of topics relevant for a given document
 - rank a set of algorithms according to their suitability for a new dataset, based on the characteristics of the dataset
- A lot of approaches rely on parametric modelling [Cheng and Hüllermeier, 2009], [Cheng et al., 2009], [Cheng et al., 2010]

Related Work

- ► Has been referred to as **label ranking** in the literature [Tsoumakas et al., 2009], [Vembu and Gärtner, 2010]
- Related to multiclass and multilabel classification
- A lot of applications, e.g : document categorization, meta-learning
 - rank a set of topics relevant for a given document
 - rank a set of algorithms according to their suitability for a new dataset, based on the characteristics of the dataset
- A lot of approaches rely on parametric modelling [Cheng and Hüllermeier, 2009], [Cheng et al., 2009], [Cheng et al., 2010]

⇒ We develop an approach free of any parametric assumptions (**local learning**) relying on results and framework developped in [Korba et al., 2017] for **ranking aggregation**.

Our Problem

Suppose we observe $(X_1, \Sigma_1), \ldots, (X_N, \Sigma_N)$ i.i.d. copies of the pair (X, Σ) , where

- $X \sim \mu$, where μ is a distribution on some feature space \mathcal{X}
- $\Sigma \sim P_X$, where P_X is the conditional probability distribution (on \mathfrak{S}_n): $P_X(\sigma) = \mathbb{P}[\Sigma = \sigma | X]$

Ex: Users *i* with characteristics X_i order items by preference resulting in Σ_i .

Goal: Learn a predictive ranking rule :

 $s : \mathcal{X} \to \mathfrak{S}_n$ $x \mapsto s(x)$ which given a random vector X, predicts the permutation Σ on the n items.

Our approach: build *piecewise constant* ranking rules, i.e: Ranking rules that are constant on each cell of a partition of \mathcal{X} built from the training data $(X_1, \Sigma_1), \ldots, (X_N, \Sigma_N)$.

Let $\mathcal{P} = \{\mathcal{C}_1, \ldots, \mathcal{C}_K\}$ be a partition of the feature space \mathcal{X} . Any $s \in S_{\mathcal{P}}$ (ranking rules that are constant on each cell of \mathcal{P}) can be written as:

$$s_{\mathcal{P},\bar{\sigma}}(x) = \sum_{k=1}^{K} \sigma_k \cdot \mathbb{I}\{x \in \mathcal{C}_k\}$$
 where $\bar{\sigma} = (\sigma_1, \dots, \sigma_K)$

Let $\mathcal{P} = \{\mathcal{C}_1, \ldots, \mathcal{C}_K\}$ be a partition of the feature space \mathcal{X} .

Any $s\in\mathcal{S}_{\mathcal{P}}$ (ranking rules that are constant on each cell of $\mathcal{P})$ can be written as:

$$s_{\mathcal{P},\bar{\sigma}}(x) = \sum_{k=1}^{K} \sigma_k \cdot \mathbb{I}\{x \in \mathcal{C}_k\}$$
 where $\bar{\sigma} = (\sigma_1, \dots, \sigma_K)$

Two methods are investigated:

k-nearest neighbor (Voronoi partitioning)



decision tree (Recursive partitioning)



Compute Local Labels/Medians

For classification, the label of a cell (ex: a leaf) is the **majority** label among the training data which fall in this cell.



Compute Local Labels/Medians

For classification, the label of a cell (ex: a leaf) is the **majority** label among the training data which fall in this cell.



Problem: Our labels are *permutations* σ :

For a cell C_k , if $\sigma_1, \ldots, \sigma_N \in C_k$, how do we aggregate them into a final label σ^* ?

 \implies Ranking aggregation problem.

Outline

Ranking Regression

Background and Results on Ranking Aggregation

Risk Minimization for Ranking (Median) Regression

Algorithms - Local Median Methods

Ranking Aggregation

Suppose we have a dataset of rankings/permutations

 $\mathcal{D}_N = (\sigma_1, \dots, \sigma_N) \in \mathfrak{S}_n^N.$ We want to find a global order

("consensus") σ^* on the n items that best represents the dataset.

Ranking Aggregation

Suppose we have a dataset of rankings/permutations $\mathcal{D}_N = (\sigma_1, \dots, \sigma_N) \in \mathfrak{S}_n^N$. We want to find a global order ("consensus") σ^* on the *n* items that best represents the dataset.

Kemeny's rule (1959) Find the solution of :

$$\sigma^* = \operatorname*{argmin}_{\sigma \in \mathfrak{S}_n} \sum_{k=1}^{N} d(\sigma, \sigma_k)$$

where d is the Kendall's tau distance:

$$d_{\tau}(\sigma, \sigma') = \sum_{i < j} \mathbb{I}\{(\sigma(i) - \sigma(j))(\sigma'(i) - \sigma'(j)) < 0\},\$$

Ex: σ = 1234, σ' = 2413 $\Rightarrow d_{\tau}(\sigma, \sigma') = 3$ (disagree on (12),(14),(34)).

Ranking Aggregation

Suppose we have a dataset of rankings/permutations $\mathcal{D}_N = (\sigma_1, \dots, \sigma_N) \in \mathfrak{S}_n^N$. We want to find a global order ("consensus") σ^* on the *n* items that best represents the dataset.

Kemeny's rule (1959) Find the solution of :

$$\sigma^* = \operatorname*{argmin}_{\sigma \in \mathfrak{S}_n} \sum_{k=1}^{N} d(\sigma, \sigma_k)$$

where d is the Kendall's tau distance:

$$d_{\tau}(\sigma, \sigma') = \sum_{i < j} \mathbb{I}\{(\sigma(i) - \sigma(j))(\sigma'(i) - \sigma'(j)) < 0\},\$$

Ex: σ = 1234, σ '= 2413 \Rightarrow $d_{\tau}(\sigma, \sigma') = 3$ (disagree on (12),(14),(34)).

Problem: Solving (1) is NP-hard.

Other methods

Idea: Compute a score for each object $i \in \{1, ..., n\}$, and sort the objets according to these scores.



• Copeland score of *i* (nb or pairwise victories): $s_C(i) = \frac{1}{N} \sum_{t=1}^N \sum_{j \neq i} \mathbb{I}[\sigma_t(i) < \sigma_t(j)]$



Other methods

Idea: Compute a score for each object $i \in \{1, ..., n\}$, and sort the objets according to these scores.



• Copeland score of *i* (nb or pairwise victories): $s_C(i) = \frac{1}{N} \sum_{t=1}^N \sum_{j \neq i} \mathbb{I}[\sigma_t(i) < \sigma_t(j)]$



Problem: Do not verify as many properties as Kemeny's rule.

Statistical Ranking Aggregation [Korba et al., 2017]

Probabilistic Modeling

 $\mathcal{D}_N = (\Sigma_1, \dots, \Sigma_N)$ with $\Sigma_k \sim P$

where $P \sim \mathfrak{S}_n$.

Statistical Ranking Aggregation [Korba et al., 2017]

Probabilistic Modeling

 $\mathcal{D}_N = (\Sigma_1, \dots, \Sigma_N)$ with $\Sigma_k \sim P$

where $P \sim \mathfrak{S}_n$.

Definition A **Kemeny median** of *P* is solution of:

 $\sigma_{P}^{*} = \operatorname*{argmin}_{\sigma \in \mathfrak{S}_{n}} L_{P}(\sigma), \tag{1}$

where $L_{\mathbf{P}}(\sigma) = \mathbb{E}_{\Sigma \sim \mathbf{P}}[d(\sigma, \Sigma)]$ is **the risk** of σ .

Question: Can we exhibit some conditions on P so that solving (1) is tractable?

Let $p_{i,j} = \mathbb{P}[\Sigma(i) < \Sigma(j)]$ (probability that item $i \succ j$).

Let $p_{i,j} = \mathbb{P}[\Sigma(i) < \Sigma(j)]$ (probability that item $i \succ j$).

Strict Stochastic Transitivity (SST):

$$p_{i,j} > 1/2 \text{ and } p_{j,k} > 1/2 \Rightarrow p_{i,k} > 1/2.$$

Let $p_{i,j} = \mathbb{P}[\Sigma(i) < \Sigma(j)]$ (probability that item $i \succ j$).

Strict Stochastic Transitivity (SST):

$$p_{i,j} > 1/2 \text{ and } p_{j,k} > 1/2 \Rightarrow p_{i,k} > 1/2.$$

Low-Noise condition NA(h) for some h > 0:

$$\min_{i< j} |p_{i,j} - 1/2| \ge h.$$

Let $p_{i,j} = \mathbb{P}[\Sigma(i) < \Sigma(j)]$ (probability that item $i \succ j$).

Strict Stochastic Transitivity (SST):

$$p_{i,j} > 1/2 \text{ and } p_{j,k} > 1/2 \Rightarrow p_{i,k} > 1/2.$$

Low-Noise condition NA(h) for some h > 0:

$$\min_{i< j} |p_{i,j} - 1/2| \ge h.$$

Our result Suppose P satisfies **SST and NA**(h) for a given h > 0. Then with overwhelming probability $1 - \frac{n(n-1)}{4}e^{-\alpha_h N}$:

Let $p_{i,j} = \mathbb{P}[\Sigma(i) < \Sigma(j)]$ (probability that item $i \succ j$).

Strict Stochastic Transitivity (SST):

$$p_{i,j} > 1/2 \text{ and } p_{j,k} > 1/2 \Rightarrow p_{i,k} > 1/2.$$

Low-Noise condition **NA**(h) for some h > 0:

$$\min_{i< j} |p_{i,j} - 1/2| \ge h.$$

Our result Suppose P satisfies **SST and NA**(h) for a given h > 0. Then with overwhelming probability $1 - \frac{n(n-1)}{4}e^{-\alpha_h N}$: \hat{P} also verifies **SST**...

Let $p_{i,j} = \mathbb{P}[\Sigma(i) < \Sigma(j)]$ (probability that item $i \succ j$).

Strict Stochastic Transitivity (SST):

$$p_{i,j} > 1/2 \text{ and } p_{j,k} > 1/2 \Rightarrow p_{i,k} > 1/2.$$

Low-Noise condition **NA**(h) for some h > 0:

$$\min_{i< j} |p_{i,j} - 1/2| \ge h.$$

Our result

Suppose P satisfies **SST and NA**(h) for a given h > 0. Then with overwhelming probability $1 - \frac{n(n-1)}{4}e^{-\alpha_h N}$:

 \widehat{P} also verifies **SST**...and the Kemeny median of P is given by the empirical Copeland ranking:

$$\sigma^*_P(i) = 1 + \sum_{j \neq i} \mathbb{I}\{\widehat{p}_{i,j} < \frac{1}{2}\} \quad \text{for } 1 \le i \le n$$

Graph of pairwise probabilities



 \Rightarrow sort the *i*'s by increasing input degree

In practice: Pseudo-empirical Kemeny Medians • If \hat{P} is SST, compute $\sigma_{\hat{P}}^*$ with Copeland method based on $\hat{p}_{i,j}$

In practice: Pseudo-empirical Kemeny Medians

- ▶ If \widehat{P} is SST, compute $\sigma_{\widehat{P}}^*$ with Copeland method based on $\widehat{p}_{i,j}$
- Else, compute $\tilde{\sigma}_{\hat{p}}^*$ with empirical Borda count ([Jiang et al., 2011])

$$\widetilde{\sigma}^*_{\widehat{P}}(i) = \frac{1}{N} \sum_{k=1}^N \Sigma_k(i) \quad \text{ for } 1 \leq i \leq n$$



Locally consistent (curl-free)

FIGURE 2. Hodge/Helmholtz decomposition of pairwise rankings

Outline

Ranking Regression

Background and Results on Ranking Aggregation

Risk Minimization for Ranking (Median) Regression

Algorithms - Local Median Methods

Our Problem

Suppose we observe $(X_1, \Sigma_1), \ldots, (X_N, \Sigma_N)$ i.i.d. copies of the pair (X, Σ) , where

- $X \sim \mu$, where μ is a distribution on some feature space \mathcal{X}
- $\Sigma \sim P_X$, where P_X is the conditional probability distribution (on \mathfrak{S}_n): $P_X(\sigma) = \mathbb{P}[\Sigma = \sigma | X]$

Ex: Users *i* with characteristics X_i order items by preference resulting in Σ_i .

Goal: Learn a predictive ranking rule :

 $s : \mathcal{X} \to \mathfrak{S}_n$ $x \mapsto s(x)$ which given a random vector X, predicts the permutation Σ on the n items.

Performance: Measured by the risk:

 $\mathcal{R}(s) = \mathbb{E}_{X \sim \mu, \Sigma \sim P_X} \left[d_\tau \left(s(X), \Sigma \right) \right]$

Ranking Median Regression Approach

 $\mathcal{R}(s) = \mathbb{E}_{X \sim \mu} \left[\mathbb{E}_{\Sigma \sim \mathbf{P}_{\mathbf{X}}} \left[d_{\tau} \left(s(X), \Sigma \right) \right] \right] = \mathbb{E}_{X \sim \mu} \left[L_{\mathbf{P}_{\mathbf{X}}}(s(X)) \right]$ (2)

Ranking Median Regression Approach

$$\mathcal{R}(s) = \mathbb{E}_{X \sim \mu} \left[\mathbb{E}_{\Sigma \sim \mathbf{P}_{\mathbf{X}}} \left[d_{\tau} \left(s(X), \Sigma \right) \right] \right] = \mathbb{E}_{X \sim \mu} \left[L_{\mathbf{P}_{\mathbf{X}}}(s(X)) \right] \quad (2)$$

Assumption

 $\text{For }X\in\mathcal{X}, \textbf{\textit{P}}_{\textbf{\textit{X}}}\text{ is }\textbf{\textbf{SST}} \Rightarrow \sigma^*_{\textbf{\textit{P}}_{\textbf{\textit{X}}}} = \text{argmin}_{\sigma\in\mathfrak{S}_n}L_{\textbf{\textit{P}}_{\textbf{\textit{X}}}}(\sigma) \text{ is unique}.$

Optimal elements

The predictors s^* minimizing (2) are the ones that maps any point $X \in \mathcal{X}$ to the **conditional** Kemeny median of P_X :

$$s^* = \operatorname*{argmin}_{s \in \mathcal{S}} \mathcal{R}(s) \iff s^*(X) = \sigma^*_{P_X}$$

Ranking Median Regression Approach

$$\mathcal{R}(s) = \mathbb{E}_{X \sim \mu} \left[\mathbb{E}_{\Sigma \sim \mathbf{P}_{\mathbf{X}}} \left[d_{\tau} \left(s(X), \Sigma \right) \right] \right] = \mathbb{E}_{X \sim \mu} \left[L_{\mathbf{P}_{\mathbf{X}}}(s(X)) \right]$$
(2)

Assumption

 $\text{For }X\in\mathcal{X}, \textbf{\textit{P}}_{\textbf{\textit{X}}}\text{ is }\textbf{\textbf{SST}} \Rightarrow \sigma^*_{\textbf{\textit{P}}_{\textbf{\textit{X}}}} = \text{argmin}_{\sigma\in\mathfrak{S}_n}L_{\textbf{\textit{P}}_{\textbf{\textit{X}}}}(\sigma) \text{ is unique}.$

Optimal elements

The predictors s^* minimizing (2) are the ones that maps any point $X \in \mathcal{X}$ to the **conditional** Kemeny median of P_X :

$$s^* = \mathop{\mathrm{argmin}}_{s \in \mathcal{S}} \mathcal{R}(s) \ \Leftrightarrow \ s^*(X) = \sigma^*_{{\pmb{P_X}}}$$

To minimize (2) approximately:

$$\sigma^*_{P_X}$$
 for any $X \Longrightarrow \sigma^*_{P_C}$ for any $X \in \mathcal{C}$

 \Rightarrow We develop Local consensus methods.

Statistical Framework- ERM

Optimize a statistical version of the theoretical risk based on the training data (X_k, Σ_k) 's:

$$\min_{s \in \mathcal{S}} \widehat{\mathcal{R}}_{N}(s) = \frac{1}{N} \sum_{k=1}^{N} d_{\tau}(s(\boldsymbol{X}_{k}), \boldsymbol{\Sigma}_{k})$$

where $\ensuremath{\mathcal{S}}$ is the set of measurable mappings.

Statistical Framework- ERM

Optimize a statistical version of the theoretical risk based on the training data (X_k, Σ_k) 's:

$$\min_{s \in \mathcal{S}} \widehat{\mathcal{R}}_{N}(s) = \frac{1}{N} \sum_{k=1}^{N} d_{\tau}(s(\boldsymbol{X}_{k}), \boldsymbol{\Sigma}_{k})$$

where ${\cal S}$ is the set of measurable mappings.

- \Rightarrow We consider a subset $\mathcal{S}_{\mathcal{P}} \subset \mathcal{S}$:
 - ▶ rich enough so that $\inf_{s \in S_P} \mathcal{R}(s) \inf_{s \in S} \mathcal{R}(s)$ is "small"
 - ideally appropriate for greedy optimization.

 \Rightarrow $\mathcal{S}_{\mathcal{P}}\text{=}$ space of piecewise constant ranking rules ("local consensus methods")

Outline

Ranking Regression

Background and Results on Ranking Aggregation

Risk Minimization for Ranking (Median) Regression

Algorithms - Local Median Methods

Let $\mathcal{P} = {\mathcal{C}_1, \ldots, \mathcal{C}_K}$ be a partition of the feature space \mathcal{X} . Any $s \in S_{\mathcal{P}}$ (ranking rules that are constant on each cell of \mathcal{P}) can be written as:

$$s_{\mathcal{P},\bar{\sigma}}(x) = \sum_{k=1}^{K} \sigma_k \cdot \mathbb{I}\{x \in \mathcal{C}_k\}$$
 where $\bar{\sigma} = (\sigma_1, \dots, \sigma_K)$

Let $\mathcal{P} = \{\mathcal{C}_1, \ldots, \mathcal{C}_K\}$ be a partition of the feature space \mathcal{X} . Any $s \in S_{\mathcal{P}}$ (ranking rules that are constant on each cell of \mathcal{P}) can be written as:

$$s_{\mathcal{P},\bar{\sigma}}(x) = \sum_{k=1}^{K} \sigma_k \cdot \mathbb{I}\{x \in \mathcal{C}_k\}$$
 where $\bar{\sigma} = (\sigma_1, \dots, \sigma_K)$

Local Learning

Let $P_{\mathcal{C}_k}$ the cond. distr. of Σ given $X \in \mathcal{C}_k$: $P_{\mathcal{C}_k}(\sigma) = \mathbb{P}[\Sigma = \sigma | X \in \mathcal{C}_k]$

Let $\mathcal{P} = {\mathcal{C}_1, \ldots, \mathcal{C}_K}$ be a partition of the feature space \mathcal{X} . Any $s \in S_{\mathcal{P}}$ (ranking rules that are constant on each cell of \mathcal{P}) can be written as:

$$s_{\mathcal{P},\bar{\sigma}}(x) = \sum_{k=1}^{K} \sigma_k \cdot \mathbb{I}\{x \in \mathcal{C}_k\}$$
 where $\bar{\sigma} = (\sigma_1, \dots, \sigma_K)$

Local Learning

Let $P_{\mathcal{C}_k}$ the cond. distr. of Σ given $X \in \mathcal{C}_k$: $P_{\mathcal{C}_k}(\sigma) = \mathbb{P}[\Sigma = \sigma | X \in \mathcal{C}_k]$

Recall: P_X is SST for any $X \in \mathcal{X}$.

Idea: P_{C_k} is still SST and $\sigma_{P_C}^* = \sigma_{P_X}^*$ provided the C_k 's are small enough.

Theorem Suppose that:

There exists $M < \infty$ such that:

 $orall (x,x') \in \mathcal{X}^2, \ \ \sum_{i < j} |p_{i,j}(x) - p_{i,j}(x')| \le M \cdot ||x - x'||.$ Then:

$$\mathcal{R}(s_{\mathcal{P}}^*) - \mathcal{R}(s^*) \le M.\delta_{\mathcal{P}}$$

where $\delta_{\mathcal{P}}$ is the max. diameter of \mathcal{P} 's cells.

Theorem Suppose that:

There exists $M < \infty$ such that: $\forall (x, x') \in \mathcal{X}^2, \quad \sum_{i < j} |p_{i,j}(x) - p_{i,j}(x')| \le M \cdot ||x - x'||.$ Then:

$$\mathcal{R}(s_{\mathcal{P}}^*) - \mathcal{R}(s^*) \le M.\delta_{\mathcal{P}}$$

where $\delta_{\mathcal{P}}$ is the max. diameter of \mathcal{P} 's cells.

Suppose in addition that: For all $x \in \mathcal{X}$, $P_x \in \mathcal{T}$ and $H = \inf_{x \in \mathcal{X}} \min_{i < j} |p_{i,j}(x) - 1/2| > 0$. and that $P_{\mathcal{C}} \in \mathcal{T}$ for all $\mathcal{C} \in \mathcal{P}$. Then,

$$\mathbb{E}\left[d_{\tau}\left(\sigma_{P_{X}}^{*}, s_{\mathcal{P}}^{*}(X)\right)\right] \leq \sup_{x \in \mathcal{X}} d_{\tau}\left(\sigma_{P_{x}}^{*}, s_{\mathcal{P}}^{*}(x)\right) \leq (M/H) \cdot \delta_{\mathcal{P}}$$

Partitioning Methods

Goal: Generate partitions \mathcal{P}_N from the training data $(X_1, \Sigma_1), \ldots, (X_N, \Sigma_N)$.

For $\mathcal{C} \in \mathcal{P}_N$, consider its empirical distribution:

$$\widehat{P}_{\mathcal{C}} = \frac{1}{N_{\mathcal{C}}} \sum_{k: X_k \in \mathcal{C}} \delta_{\Sigma_k}$$

and compute locally its Pseudo-Empirical Kemeny median $\widetilde{\sigma}_{\widehat{P}_{\mathcal{C}}}^{*}$.

Partitioning Methods

Goal: Generate partitions \mathcal{P}_N from the training data $(X_1, \Sigma_1), \ldots, (X_N, \Sigma_N)$.

For $\mathcal{C} \in \mathcal{P}_N$, consider its empirical distribution:

$$\widehat{P}_{\mathcal{C}} = \frac{1}{N_{\mathcal{C}}} \sum_{k: X_k \in \mathcal{C}} \delta_{\Sigma_k}$$

and compute locally its Pseudo-Empirical Kemeny median $\widetilde{\sigma}_{\widehat{p}_{\alpha}}^{*}$.

Two methods are investigated:

k-nearest neighbor (Voronoi partitioning)



decision tree (Recursive partitioning)



K-Nearest Neigbors Algorithm

- 1. Fix $k \in \{1, \ldots, N\}$ and a query point $x \in \mathcal{X}$
- 2. Sort $(X_1, \Sigma_1), \ldots, (X_N, \Sigma_N)$ by increasing order of the distance to $x : ||X_{(1,N)} x|| \le \ldots \le ||X_{(N,N)} x||$
- 3. Consider next the empirical distribution calculated using the k training points closest to \boldsymbol{x}

$$\widehat{P}(x) = \frac{1}{k} \sum_{l=1}^{k} \delta_{\Sigma_{(l,N)}}$$

and compute the pseudo-empirical Kemeny median, yielding the *k*-NN prediction at *x*:

$$s_{k,N}(x) \stackrel{def}{=} \widetilde{\sigma}^*_{\widehat{P}(x)}.$$

 \Rightarrow We recover the classical bound $\mathcal{R}(s_{k,N}) - \mathcal{R}^* = \mathcal{O}(\frac{1}{\sqrt{k}} + \frac{k}{N})$

Decision Tree

Split recursively the feature space $\ensuremath{\mathcal{X}}$ to minimize some impurity criterion.

Analog to Gini criterion in classification: m classes, f_i proportion of class $i \rightarrow I_G(f) = \sum_{i=1}^m f_i(1-f_i)$

Decision Tree

Split recursively the feature space $\ensuremath{\mathcal{X}}$ to minimize some impurity criterion.

Analog to Gini criterion in classification: m classes, f_i proportion of class $i \rightarrow I_G(f) = \sum_{i=1}^m f_i(1-f_i)$

Here, for a cell $C \in \mathcal{P}_N$:

Impurity [Alvo and Philip, 2014]:

$$\gamma_{\widehat{P}_{\mathcal{C}}} = \frac{1}{2} \sum_{i < j} \widehat{p}_{i,j}(\mathcal{C}) \left(1 - \widehat{p}_{i,j}(\mathcal{C})\right)$$

which is tractable and satisfies the double inequality

$$\widehat{\gamma}_{\widehat{P}_{\mathcal{C}}} \leq \min_{\sigma \in \mathfrak{S}_n} L_{\widehat{P}_{\mathcal{C}}}(\sigma) \leq 2\widehat{\gamma}_{\widehat{P}_{\mathcal{C}}}.$$

Terminal value : Compute the pseudo-empirical median of a cell C:

$$s_{\mathcal{C}}(x) \stackrel{def}{=} \widetilde{\sigma}^*_{\widehat{P}_{\mathcal{C}}}.$$

Conclusion

Interesting challenges:

- Most of the maths from euclidean spaces cannot be applied
- But our insights still hold
- Based on our results on ranking aggregation, we develop a novel approach to ranking regression/label ranking
- Theoretical guarantees (approximation error, rates of convergence)
- We propose two practical algorithms

Openings: How to extend to incomplete rankings (+with ties)?

Alvo, M. and Philip, L. (2014).
 Decision tree models for ranking data.
 In *Statistical Methods for Ranking Data*, pages 199–222.
 Springer.

- Cheng, W., Dembczyński, K., and Hüllermeier, E. (2010).
 Label ranking methods based on the Plackett-Luce model.
 In Proceedings of the 27th International Conference on Machine Learning (ICML-10), pages 215–222.
- Cheng, W., Hühn, J., and Hüllermeier, E. (2009).
 Decision tree and instance-based learning for label ranking.
 In Proceedings of the 26th International Conference on Machine Learning (ICML-09), pages 161–168.
- Cheng, W. and Hüllermeier, E. (2009).
 A new instance-based label ranking approach using the mallows model.

Advances in Neural Networks–ISNN 2009, pages 707–716.

Jiang, X., Lim, L.-H., Yao, Y., and Ye, Y. (2011). Statistical ranking and combinatorial hodge theory. *Mathematical Programming*, 127(1):203–244.

- Korba, A., Clémençon, S., and Sibony, E. (2017). A learning theory of ranking aggregation. In *Proceeding of AISTATS 2017*.
- Tsoumakas, G., Katakis, I., and Vlahavas, I. (2009).
 Mining multi-label data.
 In Data mining and knowledge discovery handbook, pages 667–685. Springer.
- Vembu, S. and Gärtner, T. (2010).
 Label ranking algorithms: A survey.
 In *Preference learning*, pages 45–64. Springer.