

## Label Ranking

Label Ranking (LR) corresponds to the problem of learning a hypothesis that maps instances to rankings over a finite set of labels.

Applications: pattern recognition, web advertisement and document categorization.

The importance of LR has spurred the development of applied CS approaches. However, formal theoretical guarantees for this problem are missing.

### Label Ranking through Nonparametric Regression

In LR, the learner observes labeled examples  $(x, \sigma) \in \mathbb{X} \times \mathbb{S}_k$  and the goal is to learn a function  $h : \mathbb{X} \rightarrow \mathbb{S}_k$  from instances to rankings of labels with low error on future predictions.

According to practical applications, a natural way to represent preferences is to map alternatives through a real-valued score function.

**Distribution-free Nonparametric LR:** Let  $\mathbb{X} \subseteq \mathbb{R}^d$ ,  $[k]$  be a set of labels,  $\mathcal{C}$  be a class of functions from  $\mathbb{X}$  to  $[0, 1]^k$  and  $\mathcal{D}_x$  be an arbitrary distribution over labels  $\mathcal{E}$  over  $\mathbb{R}^k$ . Let  $m$  be an unknown target function in  $\mathcal{C}$ .

An example oracle  $\text{Ex}(m, \mathcal{E})$  with complete rankings, works as follows:

- $x \sim \mathcal{D}_x$  and  $\xi \sim \mathcal{E}$  independently,
- $\sigma = \text{argsort}(m(x) + \xi)$  and
- it returns a labeled example  $(x, \sigma) \in \mathbb{X} \times \mathbb{S}_k$ .

In the noiseless case ( $\xi = 0$  almost surely), we simply write  $\text{Ex}(m)$ .

# Label Ranking through Nonpar

Dimitris Fotakis Alkis Karayannidis

NT

## Main Results

that maps feature vectors

ategorization.

Approaches for tackling this  
problem:

## Regression

The goal is to learn a hypothesis  
from seen and unseen examples.

The goal is to evaluate individual

labels,  $\mathcal{C}$  be a class of  
functions  $\mathbb{X} \rightarrow \mathbb{R}$ . Consider a noise

We give the first theoretical guarantees for the  
robustness of shallow decision trees and random forests in the noiseless setting.

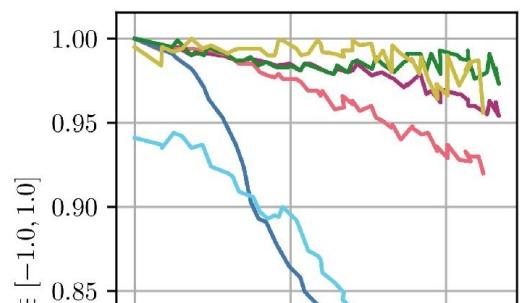
### Noiseless Setting

Under sparsity and approximate submodularity assumptions, we show that there exists a set of Decision Trees via Level-Splits that draws  $\tilde{O}(\log d)$  trees in time  $\text{Ex}(m)$  and, in  $\text{poly}_{C,r}(d, k, 1/\epsilon)$  time, computes a function  $\{0, 1\}^d \rightarrow \mathbb{S}_k$  which, with probability 99%, satisfies the following properties:

We further obtain results for random forests and gradient boosting.

**Noisy Setting** The noise distribution  $\mathcal{E}$  satisfies

- (i) the  $\alpha$ -inconsistency property if  $\mathbf{E}_{x \sim \mathcal{D}_x} [\Pr_{\xi \sim \mathcal{E}}]$
- (ii) the  $\beta$ - $k_\tau$  gap property if  $\mathbf{E}_{x \sim \mathcal{D}_x} \mathbf{E}_{\xi \sim \mathcal{E}} [k_\tau(h(x)) - k_\tau(h^*(x))]$



# ametric Regression

ICML 20:

lavasis Eleni Psaroudaki

UA

## Results

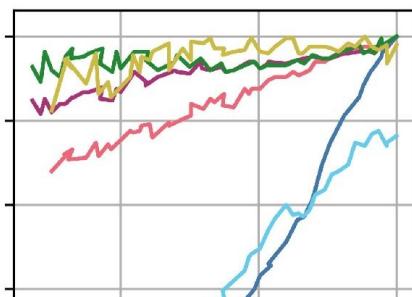
Label Ranking problem for algorithms based on less setting and extensive experimental evidence random forests in the noisy case.

conditions, there exists an algorithm based on  $\mathcal{O}(d) \cdot \text{poly}_{C,r}(k/\epsilon)$  independent samples from a set of splits  $S_n$  and an estimate  $h^{(n)}(\cdot; S_n)$ : ities  $\mathbf{E}_{x \sim \mathcal{D}_x} [\Delta_{\text{Spearman}}(h(x), h^{(n)}(x; S_n))] \leq \epsilon$ .

for the Breiman's criterion.

for some  $\alpha \in [0, 1]$  and  $\beta \in [-1, 1]$ :

$\mathbb{P}[h(x) \neq \sigma] \leq \alpha$ , and  
 $\mathbb{P}[h(x) = \sigma)] = \beta$ .



## Statistical Learning in ra

$$h^* = a$$

The learner is given i.i.d hypothesis  $\hat{h} : \mathbb{X} \rightarrow \mathcal{S}_k$  fr  $\Pr_{x \sim \mathcal{D}_x} [\hat{h}(x) \neq h^*(x)] \leq \epsilon$

**Related Work** What is Clemenccon et al. (2018)

C

**The model** Consider ar over features. Let  $y : \mathbb{R}^d \rightarrow [0, 1]^k$ . Each sample  $(x, c)$

- Draw  $x \in \mathbb{X}$  from  $\mathcal{D}_x$
- Draw  $q(x)$ -biased coin
- If  $c_i > 0$ , set  $y_i = m_i(x)$
- Compute  $\sigma = \text{argsort}($

**Assumptions** We wor  $x \in \mathbb{R}^d$ . For any  $1 \leq i <$

1. **(Strict Stochastic Tran**  
 $(n_{i+1}(x) > 1/2 \wedge n_{i+1}(x) < 1)$

## Problem Formulation

**Ranking metric  $\Delta$  and incomplete rankings** Consider the median problem  
 $\operatorname{argmin}_h \mathbf{E}_{(x,\sigma)} [\Delta(h(x), \sigma)]$ , where  $(x, \sigma) \sim \text{Ex}(m, \mathcal{E})$ .

I. samples from the oracle  $\text{Ex}(m, \mathcal{E}, \mathcal{M})$  and its goal is to output a hypothesis from some hypothesis class  $\mathcal{H}$  such that, with high probability, the error against the median  $h^*$  is small.

↳ the optimal solution  $h^*$ ? When does it exist? See Korba et al. (2017); Clemenccon and Korba (2018); Clemenccon and Vogel (2020).

## Conditions for Theoretical Guarantees

Given an underlying score hypothesis  $m : \mathbb{X} \rightarrow [0, 1]^k$  and let  $\mathcal{D}_x$  be a distribution  $[\star] \rightarrow [0, 1] \cup \{\star\}$  and consider the survival probabilities vector  $q : \mathbb{X} \rightarrow \mathbb{R}^k$ .  $\tau) \sim \text{Ex}(m, \mathcal{E}, q)$  is generated as follows:

and  $\xi \in [-\frac{1}{4}, \frac{1}{4}]^k$  from  $\mathcal{E}$ .

is  $c \in \{-1, +1\}^k$ .

$c_i) + \xi_i$ , else  $y_i = \star$ .

$y$ ), ignoring the  $\star$  symbol.

Let  $k$  with  $\Delta = \Delta_{KT}$ . Let  $p_{ij}(x) = \Pr_{\xi \sim \mathcal{E}}[m_i(x) + \xi_i > m_j(x) + \xi_j | x]$  for  $i, j \leq k$ , we assume that the following hold.

**Positivity**) For any  $x \in \mathbb{R}^d$  and any  $u \in [k]$ , we have that  $p_{uj}(x) \neq 1/2$  and  $p_{uj}(x) > 1/2 \Rightarrow n_{uz}(x) > 1/2$ .

## Problem Formulation

**Computational Learning in ranking metric  $\Delta$**  The learner is given i.i.d.  $\text{Ex}(m, \mathcal{E})$  and its goal is to efficiently output a hypothesis  $\hat{h} : \mathbb{X} \rightarrow \mathbb{S}_k$  such that, the error  $\mathbb{E}_{x \sim \mathcal{D}_x} [\Delta(\hat{h}(x), h(x))]$  is small.

**In practice** This problem is mainly solved using decision trees and random forests. Theoretical guarantees were not known for this task.

## Conditions for Theoretical Guarantee

Set  $\mathbb{X} = \{0, 1\}^d$  and the regression function  $m : \{0, 1\}^d \rightarrow [0, 1]^k$  with  $m$  be the distribution over features. We assume that the following hold for

1. **(Sparsity)**  $m_j : \{0, 1\}^d \rightarrow [0, 1]$  is  $r$ -sparse, i.e., it depends on  $r$  out of  $d$
2. **(Approximate Submodularity)** The mean squared error  $\tilde{L}_j$  of  $m_j$  is  $C$ -approximate-submodular, i.e., for any  $S \subseteq T \subseteq [d]$ ,  $i \in [d]$ , it holds that

$$\tilde{L}_j(T) - \tilde{L}_j(T \cup \{i\}) \leq C \cdot (\tilde{L}_j(S) - \tilde{L}_j(S \cup \{i\}))$$

The **Mean Squared Error** (MSE) of a function  $f : \{0, 1\}^d \rightarrow [0, 1]$  is equal

$$\tilde{L}(f, S) = \mathbb{E}_{x \sim \mathcal{D}_x} \left[ (f(x) - \mathbb{E}_{w \sim \mathcal{D}_x} [f(w) | w_S = x_S])^2 \right]$$

where  $x_S$  is the sub-vector of  $x$ , where we observe only the coordinate in  $S$  and  $x_S \in \{0, 1\}^{|S|}$ .

samples from the oracle such that, with high probability,

random forests; however,

**S**

$= (m_1, \dots, m_k)$ . Let  $\mathcal{D}_x$  be any  $j \in [k]$ .

coordinates.

hat

) .

to

$^2]$ ,

s with indices in  $S \subseteq [d]$

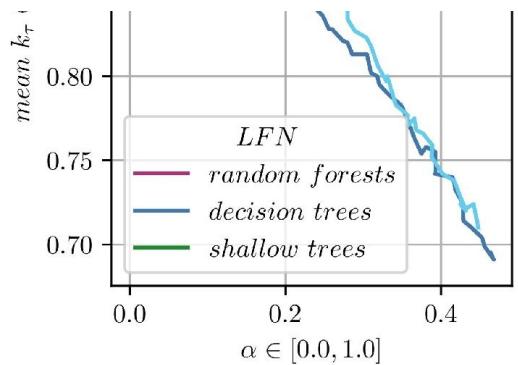


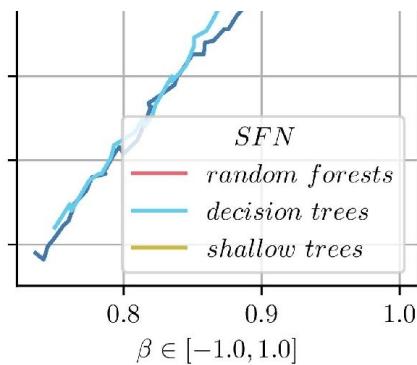
Figure 1. Experimental results in terms of mean  $k_{\tau}$

## Statistical Label Rank

**Distribution-free Nonparametric Incomplete Labeling**: Let  $\mathbb{X}$  be a class of functions from  $\mathbb{X}$  to  $[0, 1]^k$  and  $\mathcal{D}_x$  be an arbitrary distribution over  $\mathbb{X}$ . Let  $m$  be an unknown mechanism that given a tuple  $(x, y) \in \mathbb{X} \times \mathbb{R}^k$  generates

An example oracle  $\text{Ex}(m, \mathcal{E}, \mathcal{M})$  with incomplete labeling

- $x \sim \mathcal{D}_x$  and  $\xi \sim \mathcal{E}$  independently,
- $y = m(x) + \xi$ ,
- $\sigma = \mathcal{M}(x, y)$  and
- it returns a labeled example  $(x, \sigma) \in \mathbb{X} \times \mathbb{S}_{\leq}$



KT coefficient for different noise distributions  $\mathcal{E}$ .

## Ranking with incomplete rankings

LR: Let  $\mathbb{X} \subseteq \mathbb{R}^d$ ,  $[k]$  be a set of labels,  $\mathcal{C}$  be an arbitrary distribution over  $\mathbb{X}$ . Consider a noise target function in  $\mathcal{C}$ . Let  $\mathcal{M}$  be a randomized algorithm that generates an incomplete ranking  $\mathcal{M}(x, y) \in \mathbb{S}_{\leq k}$ .

te rankings, works as follows:

$k$ .

2. **(Tsybakov's Noise Coefficient)** random feature  $x \sim \mathcal{I}$
3. **(Deletion Tolerance)** the survival probability

Let  $\epsilon, \delta \in (0, 1)$ . Consider

Under conditions (1)-(3)

samples from  $\text{Ex}(m, \mathcal{E})$  and  $h^\star(x)$ ] is, with probability

where  $h^\star$  is the optimal

$$h^\star(x; i)$$

and  $L_{i,j}^\star$  is the loss of constant depending on

$$L_{i,j}^\star(g)$$

**Condition**) There exists  $a \in [0, 1]$  and  $B > 0$  so that the probability that a  $x$  satisfies  $|p_{ij}(x) - 1/2| < 2t$ , is at most  $B \cdot t^{a/(1-a)}$  for all  $t \geq 0$ .

There exists  $\phi \in (0, 1]$  so that  $q_{i,j}(x) \geq \phi$  for any  $x \in \mathbb{R}^d$ , where  $q_{i,j}(x)$  is  $y$  of the pair  $i < j$  in  $x$ .

## Main Results

der a hypothesis class  $\mathcal{G}$  of binary classifiers with finite VC dimension. 3) with parameters  $a, B, \phi$ , there exists an algorithm that draws

$$n = \tilde{O} \left( \frac{k^{\frac{4(1-a)}{a}}}{\text{poly}_a(\phi \cdot \epsilon)} \cdot \max \left\{ \log \left( \frac{k}{\delta} \right), \text{VC}(\mathcal{G}) \right\} \right)$$

,  $q$ ) and computes an estimate  $\hat{h} : \mathbb{R}^d \rightarrow \mathbb{S}_k$  so that  $\Pr_{x \sim \mathcal{D}_x}[\hat{h}(x) \neq$  ity  $1 - \delta$ , at most

$$\frac{C_{a,B}}{\phi^2} \left( 2 \sum_{i < j} \left( \inf_{g \in \mathcal{G}} L_{i,j}(g) - L_{i,j}^* \right)^a \right) + \epsilon,$$

l predictor given by

$$) = 1 + \sum_{j \neq i} 1\{p_{ij}(x) < 1/2\} = 1 + \sum_{j \neq i} 1\{h_{ij}^*(x) = -1\}$$

the binary Bayes classifiers  $h_{i,j}^*$  for  $1 \leq i < j \leq k$ , where  $C_{a,B}$  is a  $a, B$ . In particular,

$$) = \mathbf{E}_{(x,\sigma) \sim \text{Ex}(m, \mathcal{E}, q)} [g(x) \neq \text{sgn}(\sigma(i) - \sigma(j)) | \sigma \ni i, j].$$