# Nonparametric Regression with Comparisons: Escaping the Curse of Dimensionality with Ordinal Information

Yichong Xu<sup>1</sup> Hariank Muthakana<sup>1</sup> Sivaraman Balakrishnan<sup>2</sup> Artur Dubrawski<sup>3</sup> Aarti Singh<sup>1</sup>

### **Abstract**

In supervised learning, we leverage a labeled dataset to design methods for function estimation. In many practical situations, we are able to obtain alternative feedback, possibly at a low cost. A broad goal is to understand the usefulness of, and to design algorithms to exploit, this alternative feedback. We focus on a semi-supervised setting where we obtain additional ordinal (or comparison) information for potentially unlabeled samples. We consider ordinal feedback of varying qualities where we have either a perfect ordering of the samples, a noisy ordering of the samples or noisy pairwise comparisons between the samples. We provide a precise quantification of the usefulness of these types of ordinal feedback in nonparametric regression, showing that in many cases it is possible to accurately estimate an underlying function with a very small labeled set, effectively escaping the curse of dimensionality. We develop an algorithm called Ranking-Regression (R<sup>2</sup>) and analyze its accuracy as a function of size of the labeled and unlabeled datasets and various noise parameters. We also present lower bounds, that establish fundamental limits for the task and show that R<sup>2</sup> is optimal in a variety of settings. Finally, we present experiments that show the efficacy of R<sup>2</sup> and investigate its robustness to various sources of noise and model-misspecification.

#### 1. Introduction

Classical nonparametric regression is centered around the development and analysis of methods that use labeled observations,  $\{(X_1, y_1), \dots, (X_n, y_n)\}$ , where  $(X_i, y_i) \in$ 

Proceedings of the 35<sup>th</sup> International Conference on Machine Learning, Stockholm, Sweden, PMLR 80, 2018. Copyright 2018 by the author(s).

 $\mathbb{R}^d \times \mathbb{R}$ , in various tasks of estimation and inference. Nonparametric methods are appealing in practice owing to their flexibility, and the relatively weak a-priori structural assumptions that they impose on the unknown regression function. However, the price we pay is that nonparametric methods typically require a large amount of labeled data, scaling exponentially with the dimension, to estimate complex target functions - the so-called curse of dimensionality. This has motivated research on structural constraints – for instance, sparsity or manifold constraints - as well as research on active learning and semi-supervised learning where labeled samples are used judiciously. We consider a complementary approach, motivated by applications in material science, crowdsourcing, and healthcare, where we are able to supplement a small labeled dataset with a potentially larger dataset of ordinal information. Such ordinal information is obtained either in the form of a (noisy) ranking of unlabeled points or in the form of (noisy) pairwise comparisons between function values at unlabeled points.

In crowdsourcing we rely on human labeling effort, and in many cases humans are able to provide more accurate ordinal feedback with substantially less effort (see for instance (Tsukida & Gupta, 2011; Shah et al., 2015)). We investigate a task of this flavor in Section 6. In material synthesis the broad goal is to design complex new materials and machine learning approaches are gaining popularity (Xue et al., 2016; Faber et al., 2016). Typically given a setting of input parameters (temperature, pressure etc.) we are able to perform a synthesis experiment and measure the quality of resulting synthesized material. Understanding this quality landscape is essentially a task of high-dimensional function estimation. Synthesis experiments can be costly and material scientists when presented with pairs of input parameters are often able to cheaply provide noisy comparative assessments of synthesis quality. Similarly, in clinical settings, precise assessment of an individual patient's health readings can be difficult, expensive and/or risky, but comparing the relative status of two patients may be relatively easy and accurate. In each of these settings, it is important to develop methods for function estimation that combine standard supervision with (potentially) cheaper and abundant ordinal or comparative supervision.

<sup>&</sup>lt;sup>1</sup>Machine Learning Department, Carnegie Mellon University, Pittsburgh, USA <sup>2</sup>Department of Statistics and Data Science, Carnegie Mellon University, Pittsburgh, USA <sup>3</sup>Auton Lab, Carnegie Mellon University, Pittsburgh, USA. Correspondence to: Yichong Xu <yichongx@cs.cmu.edu>.

**Related Work:** There is considerable work in supervised and unsupervised learning on incorporating additional types of feedback beyond labels. For instance, the papers (Zou et al., 2015) and (Poulis & Dasgupta, 2017) study the benefits of different types "feature feedback" in clustering and supervised learning respectively. There is also a vast literature on models and methods for analyzing pairwise comparison data, like the classical Bradley-Terry (Bradley & Terry, 1952) and Thurstone (Thurstone, 1927) models. In this literature, the typical focus is on ranking or quality estimation for a fixed set of objects. In contrast, we focus on function estimation and the resulting models and methods are quite different. We build on work on "noisy sorting" (Braverman & Mossel, 2009) to extract a consensus ranking from noisy pairwise comparisons. Most close in spirit to our own work are the two recent papers (Kane et al., 2017; Xu et al., 2017), which consider binary classification with ordinal information. These works differ from ours in their focus on classification, emphasis on active querying strategies and use of quite different ordinal feedback models. Finally, given ordinal information of sufficient fidelity, the problem of nonparametric regression is related to the problem of regression with shape constraints, or more specifically isotonic regression (Barlow, 1972; Zhang, 2002). Accordingly, we leverage algorithms from this literature in our work and we comment further on the connections in Section 3. Some salient differences between this literature and our work are that we design methods that work in a semisupervised setting, and further that our target is an unknown d-dimensional (smooth) regression function as opposed to a univariate shape-constrained function.

Our Contributions: We develop the Ranking-Regression  $(R^2)$  algorithm for nonparametric regression that can leverage ordinal information, in addition to direct labels. Theoretical analysis and practical experiments show the strength of our algorithm.

- To establish the usefulness of ordinal information in nonparametric regression, in Section 3 we consider the idealized setting where we obtain a perfect ordering of the unlabeled set. We show that the Mean Squared Error (MSE) of  $\mathbb{R}^2$  can be bounded by  $\widetilde{O}(m^{-2/3}+n^{-2/d})^1$ , where m denotes the number of labeled samples and n the number of ranked samples. To achieve an MSE of  $\varepsilon$ , the number of labeled samples required by  $\mathbb{R}^2$  is *independent* of dimension. This result establishes that sufficient ordinal information of high quality can allow us to effectively circumvent the curse of dimensionality.
- In Section 4 we analyze R<sup>2</sup> when using a noisy ranking. We show that the MSE is bounded as  $\widetilde{O}(m^{-2/3} + \sqrt{\nu} +$

- $n^{-2/d}$ ), where  $\nu$  is the Kendall-Tau distance between the true and noisy ranking.
- As a corollary, we develop results for  $R^2$  using pairwise comparisons. If the comparison noise is bounded, the  $R^2$  algorithm can be combined with algorithms for ranking from pairwise comparisons (Braverman & Mossel, 2009) to obtain an MSE of  $\widetilde{O}(m^{-2/3} + n^{-2/d})$  when  $d \geq 4$ .
- We give information-theoretic lower bounds to characterize the fundamental limits of combining ordinal and standard supervision. These lower bounds show that our algorithms are almost optimal. In particular, the R<sup>2</sup> algorithm under perfect ranking, as well as under bounded noise comparisons, is optimal up to log factors.
- In our experiments, we test R<sup>2</sup> on simulated data, on UCI datasets and on various age-estimation tasks. Our experimental results show the advantage of R<sup>2</sup> over algorithms that only use labeled data when this labeled data is scarce. Our experiments with the age-estimation data also show the practicality of R<sup>2</sup>.

### 2. Background and Problem Setup

We consider a non-parametric regression model with random design, i.e. we suppose first that we are given access to an unlabeled set  $\mathcal{U} = \{X_1, \dots, X_n\}$ , where  $X_i \in \mathcal{X} \subset [0,1]^d$ , and  $X_i$  are drawn i.i.d. from a distribution  $\mathbb{P}_{\mathcal{X}}$ . We assume that  $\mathbb{P}_{\mathcal{X}}$  has a density p(x) which is upper and lower bounded as  $0 < p_{\min} \leq p(x) \leq p_{\max}$  for  $x \in [0,1]^d$ . Our goal is to estimate a function  $f: \mathcal{X} \mapsto \mathbb{R}$ , where following classical work (Györfi et al., 2006; Tsybakov, 2009) we assume that f is bounded in [-M, M] and belongs to a Hölder ball  $\mathcal{F}_{s,L}$ , with  $0 < s \leq 1$  where:

$$\mathcal{F}_{s,L} = \{ f : |f(x) - f(y)| \le L ||x - y||_2^s, \forall x, y \in \mathcal{X} \}.$$

For s=1 this is the class of Lipschitz functions. We discuss the estimation of smoother functions (i.e. the case when s>1) in Section 7. We obtain two forms of supervision:

1. Classical supervision: For a (uniformly) randomly chosen subset  $\mathcal{L} \subseteq \mathcal{U}$  of size m (we assume throughout that  $m \leq n$  and focus on settings where  $m \ll n$ ) we make noisy observations of the form:

$$y_i = f(X_i) + \epsilon_i, i \in \mathcal{L},$$

where  $\epsilon_i$  are i.i.d.  $\mathbb{E}[\epsilon_i] = 0$ ,  $\mathrm{Var}[\epsilon_i] = \sigma^2$ . We denote the indices of the labeled samples as  $\{t_1, \ldots, t_m\} \subset \{1, \ldots, n\}$ .

2. **Ordinal supervision:** For the given dataset  $\{X_1,\ldots,X_n\}$  we let  $\pi$  denote the *true ordering*, i.e.  $\pi$  is a permutation of  $\{1,\ldots,n\}$  such that for  $i,j\in\{1,\ldots,n\}$ , with  $\pi(i)\leq\pi(j)$  we have that  $f(X_i)\leq f(X_j)$ . We assume access to one of the following types of ordinal supervision:

 $<sup>^{1}</sup>$ We use the standard big-O notation throughout this paper, and use  $\widetilde{O}$  when we suppress log-factors.

(1) We are given access to a noisy ranking  $\widehat{\pi}$ , i.e. for a parameter  $\nu \in [0,1]$  we assume that the Kendall-Tau distance between  $\widehat{\pi}$  and the true-ordering is upper-bounded as:

$$\sum_{i,j\in[n]} \mathbb{I}[(\pi(i) - \pi(j))(\widehat{\pi}(i) - \widehat{\pi}(j)) < 0] \le \nu n^2.$$
(1)

(2) For each pair of samples  $(X_i, X_j)$ , with i < j we obtain a comparison  $Z_{ij}$  where for some constant  $\lambda > 0$ :

$$\mathbb{P}(Z_{ij} = \mathbb{I}(f(X_i) > f(X_j))) \ge \frac{1}{2} + \lambda. \tag{2}$$

As we discuss in Section 5 it is straightforward to extend our results to a setting where only a randomly chosen subset of all pairwise comparisons are observed.

Although classical supervised learning estimates a regression function with labels only and without ordinal supervision, we note that we cannot consistently estimate the underlying function with only ordinal supervision and without direct observations. In the case when no direct measurements are available the underlying function is only identifiable up to certain monotonic transformations.

Our goal is to estimate f, and the quality of an estimate  $\widehat{f}$  is assessed using the mean squared error  $\mathbb{E}(\widehat{f}(X) - f(X))^2$ , where the expectation is taken over the labeled and unlabeled training samples, as well as the new test point X. We also study the fundamental information-theoretic limits of estimation with classical and ordinal supervision by establishing lower (and upper) bounds on the minimax risk. Letting  $\eta$  denote various problem dependent parameters (the Hölder parameters s, L and various noise parameters), the minimax risk:

$$\mathfrak{M}(m, n; \eta) = \inf_{\widehat{f}} \sup_{f \in \mathcal{F}_{s, L}} \mathbb{E}(\widehat{f}(X) - f(X))^2, \quad (3)$$

provides an information-theoretic benchmark to assess the performance of an estimator. We conclude this section recalling a well-known fact: given access to only classical supervision the minimax risk  $\mathfrak{M}(m;\eta) = \Theta(m^{-\frac{2s}{2s+d}})$ , suffers from an exponential curse of dimensionality.

# 3. Nonparametric Regression with Perfect Ranking

To establish the value of ordinal information we first consider an idealized setting, where we are given a perfect ranking  $\pi$  of the unlabeled samples in  $\mathcal{U}$ . We present our Ranking-Regression (R<sup>2</sup>) algorithm with performance guarantees in Section 3.1, and a lower bound in Section 3.2 which shows that R<sup>2</sup> is optimal up to log factors.

# **Algorithm 1** R<sup>2</sup>: Ranking-Regression

**Input:** Unlabeled data  $\mathcal{U} = \{X_1, ..., X_n\}$ , a labeled set of size m and corresponding labels, i.e. samples  $\{(X_{t_1}, y_{t_1}), \ldots, (X_{t_m}, y_{t_m})\}$ , and a ranking  $\widehat{\pi}$ .

- 1: Order elements in  $\mathcal{U}$  as  $(X_{\widehat{\pi}(1)},...,X_{\widehat{\pi}(n)})$ .
- 2: Run isotonic regression (see (4)) on  $\{y_{t_1}, \dots, y_{t_m}\}$ . Denote the estimated values by  $\{\hat{y}_{t_1}, \dots, \hat{y}_{t_m}\}$ .
- 3: For i=1,2,...,n, let  $\widetilde{i}=t_k$ , where  $\widehat{\pi}(t_k)$  is the largest value such that  $\widehat{\pi}(t_k) \leq \widehat{\pi}(i), k=0,1,...,m$ , and  $\widetilde{i}=\star$  if no such  $t_k$  exists. Set

$$\widehat{y}_i = \begin{cases} \widehat{y}_{\widetilde{i}} & \text{if } \widetilde{i} \neq \star \\ 0 & \text{otherwise.} \end{cases}$$

**Output:** Function  $\widehat{f} = \text{NearestNeighbor}(\{(X_i, \widehat{y}_i)\}_{i=1}^n)$ .

### 3.1. Upper bounds for the R<sup>2</sup> Algorithm

Our non-parametric regression estimator is described in Algorithm 1 and Figure 1. We first rank all the samples in  $\mathcal U$  according to the (given or estimated) permutation  $\widehat \pi$ . We then run isotonic regression (Barlow, 1972) on the labeled samples in  $\mathcal L$  to de-noise them and borrow statistical strength. In more detail, we solve the following program to de-noise the labeled samples:

$$\min_{\substack{\{\widehat{y}_{\widehat{\pi}(t_1)}, \dots, \widehat{y}_{\widehat{\pi}(t_m)}\}}} \sum_{k=1}^{m} (\widehat{y}_{\widehat{\pi}(t_k)} - y_{\widehat{\pi}(t_k)})^2$$
s.t.  $\widehat{y}_{t_k} \leq \widehat{y}_{t_l} \ \forall \ (k, l) \ \text{such that} \ \widehat{\pi}(t_k) < \widehat{\pi}(t_l)$ 

$$-M \leq \{y_{\widehat{\pi}(t_1)}, \dots, y_{\widehat{\pi}(t_m)}\} \leq M.$$

$$(4)$$

We introduce the bounds  $\{M,-M\}$  in the above program to ease our analysis. In our experiments, we simply set M to be a large positive value so that it has no influence on our estimator. We then leverage the ordinal information in  $\widehat{\pi}$  to impute regression estimates for the unlabeled samples in  $\mathcal{U}$ , by assigning each unlabeled sample the value of the nearest (de-noised) labeled sample which has a smaller function value according to  $\widehat{\pi}$ . Finally, for a new test point, we use the imputed (or estimated) function value of the nearest neighbor in  $\mathcal{U}$ .

In the setting where we use a perfect ranking the following theorem characterizes the performance of R<sup>2</sup>:

**Theorem 1.** For constants  $C_1, C_2 > 0$  the MSE of  $\hat{f}$  is bounded by

$$\mathbb{E}(\widehat{f}(X) - f(X))^2 \le C_1 m^{-2/3} \log^2 n \log m + C_2 n^{-2s/d}.$$

Before we turn our attention to the proof of this result, we examine some consequences.

**Remarks:** (1) Theorem 1 shows a surprising dependency on the sizes of the labeled and unlabeled sets (m and n).

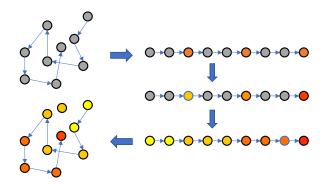


Figure 1. Top Left: A group of unlabeled points are ranked according to function values using ordinal information only. Top Right: We obtain function values of m randomly chosen samples. Middle Right: The values are adjusted using isotonic regression. Bottom Right: Function values of other unlabeled points are inferred. Bottom Left: For a new point, the estimated value is given by the nearest neighbor in  $\mathcal{U}$ .

The MSE of nonparametric regression using only the labeled samples is  $\Theta(m^{-\frac{2s}{2s+d}})$  which is exponential in d and makes non-parametric regression impractical in high-dimensions. Focusing on the dependence on m, Theorem 1 improves the rate to  $m^{-2/3} \operatorname{polylog}(m,n)$ , which is no longer exponential in d. By using enough ordinal information we can avoid the curse of dimensionality.

- (2) On the other hand, the dependence on n (which dictates the amount of ordinal information needed) is still exponential. This illustrates that ordinal information is most beneficial when it is copious. We show in Section 3.2 that this is unimprovable in an information-theoretic sense.
- (3) Somewhat surprisingly, we also observe that the dependence on n is faster than the  $n^{-\frac{2s}{2s+d}}$  rate that would be obtained if all the samples were labeled.
- (4) In the case where all points are labeled (i.e., m=n), the MSE is of order  $n^{-2/3}+n^{-2s/d}$ , again improving slightly on the rate when no ordinal information is available. The improvement is largest when  $m \ll n$ .
- (5) Finally, we also note in passing that the above theorem provides an upper bound on the minimax risk in (3).

*Proof Sketch.* We provide a brief outline and defer technical details to the Supplementary Material. For a randomly drawn point  $X \in \mathcal{X}$ , we denote by  $X_{\alpha}$  the nearest neighbor of X in  $\mathcal{U}$ . We decompose the MSE as

$$\mathbb{E}\left[(\widehat{f}(X) - f(X))^2\right] \le 2\mathbb{E}\left[(\widehat{f}(X) - f(X_\alpha))^2\right] + 2\mathbb{E}\left[(f(X_\alpha) - f(X))^2\right]. \quad (5)$$

The second term corresponds roughly to the finite-sample bias induced by the discrepancy between the function value at X and the closest labeled sample. We use standard sample-spacing arguments (see (Györfi et al., 2006)) to

bound this term. This term contributes the  $n^{-2s/d}$  rate to the final result. For the first term, we show a technical result in the Appendix (Lemma 9). Without loss of generality suppose  $f(X_{t_1}) \leq \cdots f(X_{t_m})$ . By conditioning on a probable configuration of the points and enumerating over choices of the nearest neighbor we find that roughly (see Lemma 9 for a precise statement):

$$\mathbb{E}\left[\left(\widehat{f}(X) - f(X_{\alpha})\right)^{2}\right] \leq \left(\frac{\log^{2} n \log m}{m}\right) \times \\ \mathbb{E}\left(\sum_{k=1}^{m} \left(\left(\widehat{f}(X_{t_{k}}) - f(X_{t_{k}})\right)^{2} + \left(f\left(X_{t_{k+1}}\right) - f\left(X_{t_{k}}\right)\right)^{2}\right)\right).$$

$$(6)$$

Intuitively, these terms are related to the estimation error arising in isotonic regression (first term) and a term that captures the variance of the function values (second term). When the function f is bounded, we show that the dominant term is the isotonic estimation error which is on the order of  $m^{-2/3}$ . Putting these pieces together we obtain the theorem.

#### 3.2. Lower bounds with Ordinal Data

To understand the fundamental limits on the usefulness of ordinal information, as well as to study the optimality of the  $R^2$  algorithm we now turn our attention to establishing lower bounds on the minimax risk. In our lower bounds we choose  $\mathbb{P}_{\mathcal{X}}$  to be uniform on  $[0,1]^d$ . Our estimators  $\widehat{f}$  are functions of the labeled samples:  $\{(X_{t_1},y_{t_1}),\ldots,(X_{t_m},y_{t_m})\}$ , the set  $\mathcal{U}=\{X_1,\ldots,X_n\}$  and the true ranking  $\pi$ . We have the following result:

**Theorem 2.** For any estimator  $\hat{f}$  we have that for a universal constant C > 0,

$$\inf_{\widehat{f}} \sup_{f \in \mathcal{F}_{s,L}} \mathbb{E}\left[ (f(X) - \widehat{f}(X))^2 \right] \ge C(m^{-2/3} + n^{-2s/d}).$$

Comparing with the result in Theorem 1 we conclude that the  $R^2$  algorithm is optimal up to log factors, when the ranking is noiseless.

*Proof Sketch.* We establish each term in the lower bound separately. Intuitively, for the  $n^{-2s/d}$  lower bound we consider the case when all the n points are labeled perfectly (in which case the ranking is redundant) and show that even in this setting the MSE of any estimator is at least  $n^{-2s/d}$  due to the finite resolution of the sample.

To prove the  $m^{-2/3}$  lower bound we construct a novel packing set of functions in the class  $\mathcal{F}_{s,L}$ , and use information-theoretic techniques (Fano's inequality) to establish the lower bound. The functions we construct are all increasing functions, and as a result the ranking  $\pi$  provides no additional information for these functions easing the analysis. Figure 2 contrasts the classical construction for lower

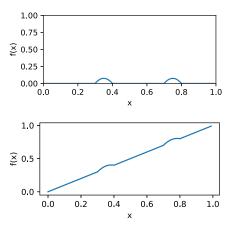


Figure 2. Original construction for nonparametric regression in 1-d (above), and our construction (below).

bounds in non-parametric regression (where tiny bumps are introduced to a reference function) with our construction where we additionally ensure the perturbed functions are all increasing. To complete the proof, we provide bounds on the cardinality of the packing set we create, as well as bounds on the Kullback-Leibler divergence between the induced distributions on the labeled samples. We provide the technical details in the Appendix.

# 4. Nonparametric Regression using Noisy Ranking

In this section, we study the setting where the ordinal information is noisy. We focus here on the setting where as in Equation (1) we obtain a ranking  $\widehat{\pi}$  whose Kendall-Tau distance from the true ranking  $\pi$  is at most  $\nu n^2$ . We show that the R<sup>2</sup> algorithm is quite robust to ranking errors and achieves an MSE of  $\widetilde{O}(m^{-2/3}+\sqrt{\nu}+n^{-2s/d})$ . We establish a complementary lower bound of  $\widetilde{O}(m^{-2/3}+\nu^2+n^{-2s/d})$  in Section 4.2.

### 4.1. Upper Bounds for the R<sup>2</sup> Algorithm

We characterize the robustness of  $\mathbb{R}^2$  to ranking errors, i.e. when  $\widehat{\pi}$  satisfies the condition in (1), in the following theorem:

**Theorem 3.** For constants  $C_1, C_2 > 0$ , the MSE of the  $R^2$  estimate  $\hat{f}$  is bounded by

$$\mathbb{E}[(\widehat{f}(X) - f(X))^{2}] \le C_{1} \left( \log^{2} n \log m \left( m^{-2/3} + \sqrt{\nu} \right) \right) + C_{2} n^{-2s/d}.$$

**Remarks:** (1) Once again we observe that in the regime where sufficient ordinal information is available, i.e. n is large, the rate no longer has an exponential dependence on the dimension d.

(2) This result also shows that the  $R^2$  algorithm is inherently robust to noise in the ranking, and the mean squared error degrades gracefully as a function of the noise parameter  $\nu$ . We investigate the optimality of the  $\sqrt{\nu}$ -dependence in the next section

(3) Finally, in settings where  $\nu$  is large  $R^2$  can be led astray by the ordinal information, and a standard non-parametric regressor can achieve the (possibly faster)  $O\left(m^{-\frac{2s}{2s+d}}\right)$  rate by ignoring the ordinal information. As we show in Appendix E a simple cross-validation procedure can combine the benefits of the two estimators to achieve a rate of  $O\left(m^{-2/3} + \min\{\sqrt{\nu}, m^{-\frac{2s}{2s+d}}\} + n^{-2s/d}\right)$ . This rate can converge to 0 if we have sufficiently many labels, even if the comparisons are very noisy. The cross validation process is standard and computationally efficient: we estimate the regression function twice, once using  $R^2$  and once using k-nearest neighbors, and choose the regression function that performs better on a held-out validation set.

We now turn our attention to the proof of this result.

*Proof Sketch.* When using an estimated permutation  $\widehat{\pi}$  the true function of interest f is no longer an increasing (isotonic) function with respect to  $\widehat{\pi}$ , and this results in a model-misspecification *bias*. The core technical novelty of our proof is in relating the upper bound on the error in  $\widehat{\pi}$  to an upper bound on this bias. Concretely, in the Appendix we show the following lemma:

**Lemma 4.** For any permutation  $\widehat{\pi}$  satisfying the condition in (1)

$$\sum_{i=1}^{n} (f(X_{\pi^{-1}(i)}) - f(X_{\widehat{\pi}^{-1}(i)}))^2 \le 8M^2 \sqrt{2\nu} n.$$

Using this result we bound the minimal error of approximating an increasing sequence according to  $\pi$  by an increasing sequence according to the estimated ranking  $\widehat{\pi}$ . We denote this error by  $\Delta$ , and using Lemma 4 we show that in expectation (over the random choice of the labeled set)

$$\mathbb{E}[\Delta] \le 8M^2 \sqrt{2\nu} m.$$

With this technical result in place we follow the same decomposition and subsequent steps before we arrive at the expression in Equation (6). In this case, the first term for some constant C>0 is bounded as:

$$\mathbb{E}\left(\sum_{k=1}^{m} \left(\widehat{f}(X_{t_k}) - f(X_{t_k})\right)^2\right) \le 2\mathbb{E}[\Delta] + Cm^{1/3},$$

where the first term corresponds to the modelmisspecification bias and the second corresponds to the usual isotonic regression rate. Putting these terms together in the decomposition in Equation (6) we obtain the theorem.

#### 4.2. Lower bounds with Noisy Ordinal Data

In this section we turn our attention to lower bounds in the setting with noisy ordinal information. In particular, we construct a permutation  $\widehat{\pi}$  such that for a pair  $(X_i, X_j)$  of points randomly chosen from  $\mathbb{P}_{\mathcal{X}}$ :

$$\mathbb{P}[(\pi(i) - \pi(j))(\widehat{\pi}(i) - \widehat{\pi}(j)) < 0] \le \nu.$$

We analyze the minimax risk of an estimator which has access to this noisy permutation  $\hat{\pi}$ , in addition to the labeled and unlabeled sets (as in Section 3.2).

**Theorem 5.** There is a constant C > 0 such that for any estimator  $\hat{f}$  taking input  $X_1, ..., X_n, y_1, ..., y_m$  and  $\hat{\pi}$ ,

$$\inf_{\widehat{f}} \sup_{f \in \mathcal{F}_{s,L}} \mathbb{E} \left( f(X) - \widehat{f}(X) \right)^2 \ge C(m^{-\frac{2}{3}} + \min\{ \nu^2, m^{-\frac{2}{d+2}} \} + n^{-2s/d} ).$$

Comparing this result with our result in Remark 3 following Theorem 3, our upper and lower bounds differ by the gap between  $\sqrt{\nu}$  and  $\nu^2$ , in the case of Lipschitz functions (s=1).

Proof Sketch. We focus on the dependence on  $\nu$ , as the other parts are identical to Theorem 2. We construct a packing set of Lipschitz functions, and we subsequently construct a noisy comparison oracle  $\widehat{\pi}$  which provides no additional information beyond the labeled samples. The construction of our packing set is inspired by the construction of standard lower bounds in non-parametric regression (see Figure 2), but we modify this construction to ensure that  $\widehat{\pi}$  is uninformative. In the classical construction we divide  $[0,1]^d$  into  $u^d$  grid points, with  $u=m^{1/(d+2)}$  and add a "bump" at a carefully chosen subset of the grid points. Here we instead divide  $[0,t]^d$  into a grid with  $u^d$  points, and add an increasing function along the first dimension, where t is a parameter we choose in the sequel.

We now describe the ranking oracle which generates the permutation  $\widehat{\pi}$ : we simply rank sample points according to their first coordinate. This comparison oracle only makes an error when both x, x' lies in  $[0, t]^d$ , and both  $x_1, x'_1$  lie in the same grid segment [tk/u, t(k+1)/u] for some  $k \in [u]$ . So the Kendall-Tau error of the comparison oracle is  $(t^d)^2 \times ((1/u)^2 \times u) = ut^{2d}$ . We choose t such that this value is less than  $\nu$ . Once again we complete the proof by lower bounding the cardinality of the packing-set for our stated choice of t, upper bounding the Kullback-Leibler divergence between the induced distributions and appealing to Fano's inequality.

# 5. Regression with Noisy Pairwise Comparisons

In this section we focus on the setting where the ordinal information is obtained in the form of noisy pairwise comparisons, following Equation (2). We investigate a natural strategy of aggregating the pairwise comparisons to form a consensus ranking  $\hat{\pi}$  and then applying the R<sup>2</sup> algorithm with this estimated ranking. We build on results from theoretical computer science, where such aggregation algorithms are studied for their connections to sorting with noisy comparators. In particular, Braverman & Mossel (2009) study noisy sorting algorithms under the noise model described in (2) and establish the following result:

**Theorem 6** ((Braverman & Mossel, 2009)). Let  $\alpha > 0$ . There exists a polynomial-time algorithm using noisy pairwise comparisons between n samples, that with probability  $1 - n^{-\alpha}$ , returns a ranking  $\hat{\pi}$  such that for a constant  $c(\alpha, \lambda) > 0$  we have that:

$$\sum_{i,j\in[n]}\mathbb{I}[(\pi(i)-\pi(j))(\widehat{\pi}(i)-\widehat{\pi}(j))<0]\leq c(\alpha,\lambda)n.$$

Furthermore, if allowed a sequential (active) choice of comparisons, the algorithm queries at most  $O(n \log n)$  pairs of samples.

Combining this result with our result on the robustness of  $R^2$  we obtain an algorithm for nonparametric regression with access to noisy pairwise comparisons with the following guarantee on its performance:

**Corollary 7.** For constants  $C_1, C_2 > 0$ ,  $R^2$  with  $\widehat{\pi}$  estimated as described above produces an estimator  $\widehat{f}$  with MSE

$$\mathbb{E}(\widehat{f}(X) - f(X))^{2} \le C_{1} m^{-2/3} \log^{2} n \log m + C_{2} \max\{n^{-2s/d}, n^{-1/2} \log^{2} n \log m\}.$$

**Remarks:** (1) From a technical standpoint this result is an immediate corollary of Theorems 3 and 6, but the extension is important from a practical standpoint. The ranking error of O(1/n) from the noisy sorting algorithm leads to an additional  $\widetilde{O}(1/\sqrt{n})$  term in the MSE. This error is dominated by the  $n^{-2s/d}$  term if  $d \ge 4s$ , and in this setting the result in Theorem 7 is also optimal up to log factors (following the lower bound in Section 3.2).

(2) We also note that the analysis in (Braverman & Mossel, 2009) extends in a straightforward way to a setting where only a randomly chosen subset of the pairwise comparisons are obtained.

### 6. Experiments & Simulations

To verify our theoretical results and test  $R^2$  in practice, we perform three sets of experiments. First, we conduct experiments on simulated data, where the noise in the labels and ranking can be controlled separately. Second, we test  $R^2$  on UCI datasets, where the rankings are simulated using labels. We present these results in Appendix A. Finally, we

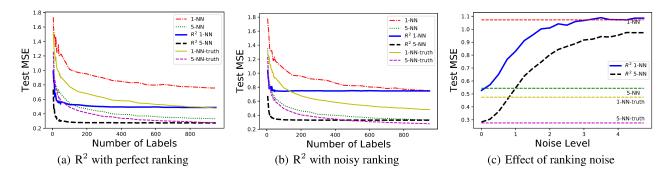


Figure 3. Experiments on simulated data. 1-NN and 5-NN represents algorithms using noisy label data only;  $R^2$  1-NN and  $R^2$  5-NN uses noisy labels as well as rankings; 1-NN-truth and 5-NN-truth uses perfect label data only.

consider a practical application of predicting people's age from portraits and we test  $R^2$  on two realistic estimation tasks

We compare  $R^2$  with k-NN algorithms in all experiments. We choose k-NN methods because they are near-optimal theoretically, and are widely used in practice. Theoretical guidelines suggest using the tuning parameter  $k_m = m^{\frac{2}{d+2}}$  when we have access to m labeled samples; however for all m,d values we considered,  $m^{\frac{2}{d+2}}$  is very small (< 5). Instead we choose a range of different constant values of k (that do not change with m) in our experiments. We repeat each experiment 20 times and report the average MSE<sup>2</sup>.

#### 6.1. Simulated Data

**Data Generation.** We generate simulated data following Härdle et al. (2012). Let d=8, and sample X uniformly random from  $[0,1]^d$ . Our target function is  $f(x)=\sum_{i=1}^d f^{(d \bmod 4)}(x_d)$ , where  $x_d$  is x's d-th dimension, and

$$f^{(1)}(x) = px - 1/2, f^{(2)}(x) = px^3 - 1/3,$$
  
$$f^{(3)}(x) = -2\sin(-px), f^{(4)}(x) = e^{-px} + e^{-1} - 1$$

with p sampled uniformly random in [0,10]. We rescale f(x) so that it has 0 mean and unit variance. The labels are generated as  $y=f(x)+\varepsilon$  where  $\varepsilon\sim\mathcal{N}(0,0.5^2)$ . We generate a training and a test set of n=1000 samples respectively. At test time, we compute the MSE  $\frac{1}{n}\sum_{i=1}^n(f(X_i^{\text{test}})-\widehat{f}(X_i^{\text{test}}))^2$  for all test data  $X_1^{\text{test}},...,X_n^{\text{test}}$ .

**Variants of R** $^2$ . We consider two variants of R $^2$ . The first variant is exactly Algorithm 1 using 1-NN as the final estimator; the second variant uses 5-NN as the final estimator. Using 5-NN does not change the asymptotic behavior of our bounds. However, we find that using 5-NN improves our estimator empirically.

**Baselines.** We compare  $R^2$  with the following baselines: i) 1-NN and 5-NN using noisy labels (x, y). Since  $R^2$  uses

ordinal data in addition to labels, it should have lower MSE than 1-NN and 5-NN. ii) 1-NN and 5-NN using perfect labels (x, f(x)). Since these algorithms use perfect labels, when m = n they serve as a benchmark for our algorithms.

 ${\bf R}^2$  with perfect rankings. In our first experiment,  ${\bf R}^2$  had access to the ranking over all 1000 training samples, while the k-NN baseline algorithms only had access to labeled samples. We varied the number of labeled samples for all algorithms from m=5 to m=1000. The results are depicted in Figure 3(a).  ${\bf R}^2$  1-NN and  ${\bf R}^2$  5-NN exhibited better performance than their counterparts using only labels, whether using noisy or perfect labels; in fact,  ${\bf R}^2$  1-NN and  ${\bf R}^2$  5-NN performed nearly the same as 1-NN or 5-NN using all 1000 perfect labels, while only requiring around 50 labeled samples.

 ${\bf R}^2$  with noisy rankings. We then consider noisy rankings; particularly in Figure 3(b), the input ranking of  ${\bf R}^2$  is obtained from noisy labels. This eliminates the need for isotonic regression in Algorithm 1, but we find that the ranking still provides useful information for the unlabeled samples. In this setting  ${\bf R}^2$  outperformed the 1-NN and 5-NN counterparts using noisy labels. However,  ${\bf R}^2$  was outperformed by algorithms using perfect labels when n=m. As expected,  ${\bf R}^2$  and k-NN with noisy labels achieved identical MSE when n=m.

**Effect of ranking noise.** We also consider the effect of ranking noise in Figure 3(c). We fixed the number of labeled/ranked samples to 100/1000, and varied the noise level of ranking. For noise level  $\sigma$ , the ranking is generated from

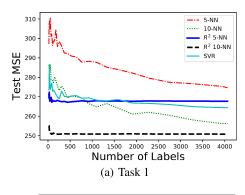
$$y' = f(x) + \varepsilon'$$

where  $\varepsilon' \sim \mathcal{N}(0, \sigma^2)$ . We varied  $\sigma$  from 0 to 5 and plotted the MSE. As  $\sigma$  goes up, the error of both variants of the  $R^2$  algorithm increases as expected.

#### **6.2. Predicting Ages from Portraits**

To further validate  $R^2$  in practice, we consider the task of estimating people's age from portraits. We use the APPA-

<sup>&</sup>lt;sup>2</sup>Our plots are best viewed in color.



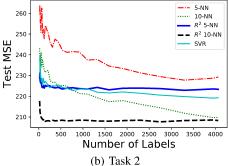


Figure 4. Experiments on age prediction.

REAL dataset (E Agustsson, 2017), which contains 7,591 images, where each image is associated a biological age and an apparent age. The biological age is the person's actual age, whereas the apparent ages are collected by crowdsourcing. Estimates from (on average 38 different) labelers are averaged to obtain the apparent age. APPA-REAL also provides the standard deviation of the apparent age estimates. The images are divided into 4113 train, 1500 validation and 1978 test samples, and we only use the train and validation samples for our experiments.

**Features and Models.** We extract the 128-dim feature for each image using the last layer of FaceNet (Schroff et al., 2015). We rescale the features so that every  $X \in [0,1]^d$ . We use 5-NN and 10-NN in this experiment. To further show the effectiveness of  $\mathbb{R}^2$ , we also compare to kernelized support vector regression (SVR). We used the standard parameter configuration in scikit-learn (Pedregosa et al., 2011), using penalty parameter C = 1, RBF kernel, and tolerance of 0.1.

**Tasks.** We considered two tasks, motivated by real-world applications.

1. In the first task, the goal is to predict biological age. The labels were biological age, whereas the ranking came from apparent ages. This is motivated by the collection process of most modern datasets where typically, an aggregated (de-noised) label is obtained through majority vote. For example, we may have the truthful biological age for a fraction of samples, but wish to collect more through crowdsourcing. In crowdsourcing, people give comparisons based on appar-

ent age instead of biological age. So we assume additional access to a ranking that comes from apparent ages.

2. In the second task, the goal is to predict the apparent age. Both labels and ranking were generated using the standard deviation provided in APPA-REAL. Labels were generated according to a Gaussian distribution with mean equal to the apparent age, and standard deviation provided in the dataset. The ranking was generated by first generating a sample of all labels using the same distribution, and ranking according to the sample. This resembles the case where we ask one single labeler for each label and comparison, to collect data for more samples. Such policy is also used in e.g., (Bi et al., 2014; Khetan et al., 2017). Note that in real applications, the ranking will have less noise than in our experiment; (Shah et al., 2015) considered exactly the same task, and showed that comparisons are more reliable than labels.

Results are depicted in Figure 4. The 10-NN version of  $R^2$  gave the best overall performance in both tasks.  $R^2$  5-NN and  $R^2$  10-NN both outperformed other algorithms when the number of labeled samples was less than 500. The performance of SVR was between 5-NN and 10-NN in our experiments. Interestingly, we observe that there is a gap between  $R^2$  and its nearest neighbor counterparts even when n=m, i.e. the ordinal information continues to be useful even when all samples are labeled, indicating the high reliability of the ordinal information for this task.

### 7. Discussion and Conclusion

We design minimax-optimal algorithms for nonparametric regression using additional ordinal information. In settings where large amounts of ordinal information are available, we find that limited direct supervision suffices to obtain accurate estimates. We provide complementary minimax lower bounds, and illustrate our proposed algorithm on real and simulated datasets. Since ordinal information is typically easier to obtain than direct labels, one might expect in these favorable settings the R<sup>2</sup> algorithm to have lower effective cost than an algorithm based purely on direct supervision.

In future work motivated by practical applications in crowd-sourcing, we hope to address the setting where both direct and ordinal supervision are actively acquired. Another possible direction is to consider *partial orders*, where we have several subsets of unlabeled data ranked, but the relation between these subsets is unknown. It would also be interesting to consider other models for ordinal information and to more broadly understand settings where indirect feedback is beneficial. Also, several recent papers (Bellec & Tsybakov, 2015; Bellec, 2018; Han et al., 2017) demonstrate the adaptivity (to complexity of the unknown parameter) of the MLE in shape-constrained problems. Understanding precise assumptions on the underlying smooth function which would induce a low-complexity isotonic regression problem is an interesting avenue for future work.

# Acknowledgements

This work is supported by AFRL grant FA8750-17-2-0212, NSF CCF-1763734, NSF DMS-1713003 and DARPA award FA8750-17-2-0130.

#### References

- Barlow, R. E. Statistical inference under order restrictions: The theory and application of isotonic regression. Technical report, 1972.
- Bellec, P. C. Sharp oracle inequalities for least squares estimators in shape restricted regression. *The Annals of Statistics*, 46(2):745–780, 2018.
- Bellec, P. C. and Tsybakov, A. B. Sharp oracle bounds for monotone and convex regression through aggregation. *Journal of Machine Learning Research*, 16:1879–1892, 2015.
- Bi, W., Wang, L., Kwok, J. T., and Tu, Z. Learning to predict from crowdsourced data. In *Uncertainty in Artificial Intelligence*, pp. 82–91, 2014.
- Bradley, R. A. and Terry, M. E. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- Braverman, M. and Mossel, E. Sorting from noisy information. *arXiv* preprint arXiv:0910.1191, 2009.
- Chaudhuri, K. and Dasgupta, S. Rates of convergence for the cluster tree. In *Advances in Neural Information Processing Systems*, pp. 343–351, 2010.
- Craig, C. C. On the tchebychef inequality of bernstein. *The Annals of Mathematical Statistics*, 4(2):94–102, 1933.
- Diaconis, P. and Graham, R. L. Spearman's footrule as a measure of disarray. *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 262–268, 1977.
- E Agustsson, R Timofte, S. E. X. B. I. G. R. R. Apparent and real age estimation in still images with deep residual regressors on appa-real database. In 12th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG), 2017. IEEE, 2017.
- Faber, F. A., Lindmaa, A., von Lilienfeld, O. A., and Armiento, R. Machine learning energies of 2 million Elpasolite(ABC2D6)crystals. *Physical Review Letters*, 117(13), 2016.
- Györfi, L., Kohler, M., Krzyzak, A., and Walk, H. *A distribution-free theory of nonparametric regression*. Springer Science & Business Media, 2006.

- Han, Q., Wang, T., Chatterjee, S., and Samworth, R. J. Isotonic regression in general dimensions. arXiv preprint arXiv:1708.09468, 2017.
- Härdle, W. K., Müller, M., Sperlich, S., and Werwatz, A. *Nonparametric and semiparametric models*. Springer Science & Business Media, 2012.
- Kane, D. M., Lovett, S., Moran, S., and Zhang, J. Active classification with comparison queries. *arXiv* preprint *arXiv*:1704.03564, 2017.
- Khetan, A., Lipton, Z. C., and Anandkumar, A. Learning from noisy singly-labeled data. *arXiv preprint arXiv:1712.04577*, 2017.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V.,
  Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P.,
  Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E.
  Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Poulis, S. and Dasgupta, S. Learning with feature feedback: From theory to practice. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, 2017.
- Schroff, F., Kalenichenko, D., and Philbin, J. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 815–823, 2015.
- Shah, N., Balakrishnan, S., Bradley, J., Parekh, A., Ramchandran, K., and Wainwright, M. Estimation from pairwise comparisons: Sharp minimax bounds with topology dependence. In *Artificial Intelligence and Statistics*, pp. 856–865, 2015.
- Shah, N., Balakrishnan, S., Guntuboyina, A., and Wainwright, M. Stochastically transitive models for pairwise comparisons: Statistical and computational issues. In *International Conference on Machine Learning*, pp. 11–20, 2016.
- Thurstone, L. L. A law of comparative judgment. *Psychological review*, 34(4):273, 1927.
- Tsukida, K. and Gupta, M. R. How to analyze paired comparison data. Technical report, DTIC Document, 2011.
- Tsybakov, A. B. Optimal aggregation of classifiers in statistical learning. *The Annals of Statistics*, 32(1):135–166, 2004.
- Tsybakov, A. B. Introduction to nonparametric estimation. revised and extended from the 2004 french original. translated by vladimir zaiats, 2009.

- Xu, Y., Zhang, H., Miller, K., Singh, A., and Dubrawski, A. Noise-tolerant interactive learning from pairwise comparisons with near-minimal label complexity. arXiv preprint arXiv:1704.05820, 2017.
- Xue, D., Balachandran, P. V., Hogden, J., Theiler, J., Xue, D., and Lookman, T. Accelerated search for materials with targeted properties by adaptive design. *Nature Communications*, 7, 2016.
- Zhang, C.-H. Risk bounds in isotonic regression. *The Annals of Statistics*, 30(2):528–555, 2002.
- Zou, J. Y., Chaudhuri, K., and Kalai, A. T. Crowdsourcing feature discovery via adaptively chosen comparisons. In *Proceedings of the Third AAAI Conference on Human Computation and Crowdsourcing, HCOMP 2015, November 8-11, 2015, San Diego, California.*, pp. 198, 2015.