

Worst-case Error Bounds for Online Learning of Smooth Functions

Weian (Andrew) Xie

WEIANXIE@MIT.EDU

Massachusetts Institute of Technology

Editors: Steve Hanneke and Tor Lattimore

Abstract

Online learning is a model of machine learning where the learner is trained on sequential feedback. We investigate worst-case error for the online learning of real functions that have certain smoothness constraints. Suppose that \mathcal{F}_q is the class of all absolutely continuous functions $f : [0, 1] \rightarrow \mathbb{R}$ such that $\|f'\|_q \leq 1$, and $\text{opt}_p(\mathcal{F}_q)$ is the best possible upper bound on the sum of the p^{th} powers of absolute prediction errors for any number of trials guaranteed by any learner. We show that for any $\delta, \epsilon \in (0, 1)$, $\text{opt}_{1+\delta}(\mathcal{F}_{1+\epsilon}) = O(\min(\delta, \epsilon)^{-1})$. Combined with the previous results of Kimber and Long (1995) and Geneson and Zhou (2023), we achieve a complete characterization of the values of $p, q \geq 1$ that result in $\text{opt}_p(\mathcal{F}_q)$ being finite, a problem open for nearly 30 years.

We study the learning scenarios of smooth functions that also belong to certain special families of functions, such as polynomials. We prove a conjecture by Geneson and Zhou (2023) that it is not any easier to learn a polynomial in \mathcal{F}_q than it is to learn any general function in \mathcal{F}_q . We also define a noisy model for the online learning of smooth functions, where the learner may receive incorrect feedback up to $\eta \geq 1$ times, denoting the worst-case error bound as $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q)$. We prove that $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q)$ is finite if and only if $\text{opt}_p(\mathcal{F}_q)$ is. Moreover, we prove for all $p, q \geq 2$ and $\eta \geq 1$ that $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q) = \Theta(\eta)$.

Keywords: online learning, mistake-bound model, smooth functions, noisy labels, single-variable functions, worst-case error

1. Introduction

Say that a learner wants to make daily predictions for stock price values, given relevant inputs such as industry performance, market-wide economic trends, or recent company performance. The next day, the learner obtains some form of feedback on its previous prediction—the actual price value, for example—and generates better-informed future predictions. This model of machine learning is called online learning. The question of how fast the learner can use its accumulated information to generate better predictions in the worst case scenario arises naturally, and is the original motivation for the model of online learning previously studied in (Angluin, 1988; Geneson and Zhou, 2023; Littlestone, 1988; Long, 2000; Kimber and Long, 1995; Mycielski, 1988; Littlestone and Warmuth, 1989).

Worst-case error bounds for online learning were first studied on functions over a discrete domain and with output values among a finite set, in (Angluin, 1988; Littlestone, 1988; Mycielski, 1988). The model over real-valued single-variable smooth functions $f : [0, 1] \rightarrow \mathbb{R}$ that we investigate was first introduced in Kimber and Long (1995), later studied in Long (2000) and most recently Geneson and Zhou (2023). The last paper, alongside Barron (1991) and Härdle (1991), also extended the problem to multi-variable smooth functions $f : [0, 1]^d \rightarrow \mathbb{R}$.

As our problem investigates how much accuracy the learner can guarantee in the worst-case scenario, we naturally represent the learning process as a game between the learner and an adversary,

the latter of which is trying to force as much error as possible. In the discrete model, different forms of feedback have been studied, including the standard model, where the adversary tells the learner the precise function output at the queried input; the bandit model (Feng et al., 2023), where the adversary only tells the learner YES or NO based on whether the learner’s answer is correct; and the noisy model, where the adversary is allowed to give incorrect reinforcement up to η times, for some $\eta \geq 1$. The smooth function problem has so far only been studied in the context of the standard model. Indeed, the bandit model would not be interesting to study in the context of smooth functions, as the adversary can guarantee infinite error each time—through answering NO each time and simply setting the function $f \equiv C$ for some large constant C . However, the noisy model is interesting to study. For the noisy model in the smooth function setting, the adversary provides feedback aligning with a function from a target class at most of the points, but allow for up to η errors/deviations.

1.1. Definitions

In the standard model of online learning, an algorithm A tries to learn a function from a given class of functions \mathcal{F} . The class \mathcal{F} consists of functions $f : S \rightarrow \mathbb{R}$ for some fixed input set S . In the t^{th} trial, the algorithm is repeatedly given an input x_t within S and queried on the value of $f(x_t)$, whereupon it outputs a prediction \hat{y}_t . Afterwards, the true value of $f(x_t)$ is revealed to A , and the raw error of the round $e_t = |\hat{y}_t - f(x_t)|$ is recorded. The total error function is calculated as the sum of the p^{th} powers of the raw errors, for some parameter $p \geq 1$. Specifically, as per the notation of Geneson and Zhou (2023), for a fixed learning algorithm A operating on a finite sequence of inputs $\sigma = (x_0, x_1, \dots, x_m) \in S^{m+1}$, and function $f \in \mathcal{F}$, this sum is $L_p(A, f, \sigma) = \sum_{t=1}^m e_t^p = \sum_{t=1}^m |\hat{y}_t - f(x_t)|^p$. Subsequently, we define the worst-case scenario learning error for a fixed algorithm A over a class \mathcal{F} , as

$$L_p(A, \mathcal{F}) = \sup_{f \in \mathcal{F}, \sigma \in \cup_{m \in \mathbb{Z}^+} S^m} L_p(A, f, \sigma).$$

We then define the worst-case error for the best possible learning algorithm:

$$\text{opt}_p(\mathcal{F}) = \inf_A L_p(A, \mathcal{F}).$$

In Kimber and Long (1995), Long (2000), and Geneson and Zhou (2023), the family \mathcal{F} studied was the class of absolutely continuous single-variable functions $f : [0, 1] \rightarrow \mathbb{R}$ that satisfy certain smoothness constraints, in the form of their derivatives having bounded norms. Namely, for all $q \geq 1$, define \mathcal{F}_q to be the class of absolutely continuous functions $f : [0, 1] \rightarrow \mathbb{R}$ such that $\int_0^1 |f'(x)|^q \leq 1$. Naturally, an extension of this definition is \mathcal{F}_∞ , denoting the class of absolutely continuous functions $f : [0, 1] \rightarrow \mathbb{R}$ such that $\sup_{x \in (0,1)} |f'(x)| \leq 1$. As noted in Geneson and Zhou (2023) and Long (2000), \mathcal{F}_∞ contains precisely the functions $f : [0, 1] \rightarrow \mathbb{R}$ that satisfy $|f(x) - f(y)| < |x - y|$ for all $x, y \in [0, 1]$. Geneson and Zhou (2023) note that for any $q \geq 1$, the range of any function $f \in \mathcal{F}_q$ is at most 1. Furthermore, they also note that $\mathcal{F}_\infty \subseteq \mathcal{F}_q \subseteq \mathcal{F}_r$ for all $1 \leq r \leq q$ by Jensen’s Inequality, from which it follows that $\text{opt}_p(\mathcal{F}_\infty) \leq \text{opt}_p(\mathcal{F}_q) \leq \text{opt}_p(\mathcal{F}_r)$ for all $1 \leq r \leq q$.

We also carry over some notation from Kimber and Long (1995), Long (2000), and Geneson and Zhou (2023). Let the q -action of a function $f : [0, 1] \rightarrow \mathbb{R}$, denoted by $J_q[f]$, be defined as

$$J_q[f] = \int_0^1 |f'(x)|^q dx.$$

As such, \mathcal{F}_q is the set of functions whose q -action is less than or equal to 1.

Furthermore, given a set $S = \{(u_i, v_i) : 1 \leq i \leq m\}$, where each $(u_i, v_i) \in [0, 1] \times \mathbb{R}$ such that $u_i < u_{i+1}$ for all $1 \leq i \leq m - 1$, define $f_S : [0, 1] \rightarrow \mathbb{R}$ such that $f_\emptyset \equiv 0$ and

$$f_S(x) = \begin{cases} v_1 & x \leq u_1 \\ v_i + \frac{(x-u_i)(v_{i+1}-v_i)}{u_{i+1}-u_i} & x \in (u_i, u_{i+1}] \\ v_m & x > u_m \end{cases}$$

if $|S|$ is nonzero. Therefore, graphically, f_S is a continuous piecewise function composed of various line segments.

As per [Kimber and Long \(1995\)](#), we define the learning algorithm LININT using f_S . In particular, on trial 0, LININT's prediction is $\hat{y}_0 = \text{LININT}(\emptyset, x_1) = 0$, and in any subsequent trial $t > 0$, its prediction is

$$\hat{y}_t = \text{LININT}((x_0, f(x_0)), \dots, (x_{t-1}, f(x_{t-1})), x_t) = f_{\{(x_0, f(x_0)), \dots, (x_{t-1}, f(x_{t-1}))\}}(x_t).$$

Intuitively, this means that LININT makes a prediction either based on linear interpolation on the two closest surrounding points (one on each side of the input) which the algorithm knows the true value of, or simply based on its closest neighbor, if the requested input is to the left or to the right of all known points.

1.2. Smooth Functions in the Standard Model

The problem of determining the value of $\text{opt}_p(\mathcal{F}_q)$ across all $p, q \geq 1$ has been extensively explored in ([Kimber and Long, 1995](#); [Long, 2000](#); [Geneson and Zhou, 2023](#); [Cesa-Bianchi et al., 1996a](#); [Faber and Mycielski, 1991](#)). [Kimber and Long \(1995\)](#) proved that if we set $q = 1$, then for any $p \geq 1$, $\text{opt}_p(\mathcal{F}_1) = \infty$. That is, no matter what parameter p is chosen, the learner can never guarantee finite error when learning a function from \mathcal{F}_1 . Furthermore, the same work also established that $\text{opt}_1(\mathcal{F}_\infty) = \infty$, from which it follows that for any $q \geq 1$, we have $\text{opt}_1(\mathcal{F}_q) = \infty$ as well. On the other hand, they proved that $\text{opt}_p(\mathcal{F}_q) = 1$ for all $p, q \geq 2$, from which it follows that $\text{opt}_p(\mathcal{F}_\infty) = 1$ for all $p \geq 2$ as well.

[Kimber and Long \(1995\)](#) also established that $\text{opt}_{1+\epsilon}(\mathcal{F}_q) = O(\epsilon^{-1})$ for all $\epsilon \in (0, 1)$ and all $q \geq 2$. [Geneson and Zhou \(2023\)](#) improved this bound and established a lower bound differing by a constant factor, proving that $\text{opt}_{1+\epsilon}(\mathcal{F}_\infty) = \Theta(\epsilon^{-\frac{1}{2}})$ and $\text{opt}_{1+\epsilon}(\mathcal{F}_q) = \Theta(\epsilon^{-\frac{1}{2}})$ for all $\epsilon \in (0, 1)$ and all $q \geq 2$. They also established that $\text{opt}_2(\mathcal{F}_{1+\epsilon}) = \Theta(\epsilon^{-1})$ for any $\epsilon \in (0, 1)$. From this, the upper bound $\text{opt}_p(\mathcal{F}_{1+\epsilon}) = O(\epsilon^{-1})$ for all $p \geq 2$ and $\epsilon \in (0, 1)$ follows. Furthermore, for any $q > 1$, they established that for sufficiently large p , the learner can guarantee an error of at most 1. Specifically, for any $q > 1$ and $p \geq 2 + \frac{1}{q-1}$, $\text{opt}_p(\mathcal{F}_q) = 1$. For $q > 1$ and $2 \leq p < 2 + \frac{1}{q-1}$, the value $\text{opt}_p(\mathcal{F}_q)$ is still finite—this is guaranteed by the prior bound $\text{opt}_p(\mathcal{F}_{1+\epsilon}) = O(\epsilon^{-1})$ above—although it is not currently known whether it equals exactly 1 throughout this sub-range.

As such, upper bounds for $\text{opt}_p(\mathcal{F}_q)$ have been established for all $p, q > 1$ whenever at least one of $p \geq 2$ and $q \geq 2$ holds. However, no upper bounds for any instance of $p, q \in (1, 2)$ have been proved. In fact, there was no proof of finiteness for $\text{opt}_p(\mathcal{F}_q)$ for any choice of $p, q \in (1, 2)$. In this paper, we establish an upper bound on $\text{opt}_p(\mathcal{F}_q)$ for all values of $p, q \in (1, 2)$.

Theorem 1.1. *For $\delta, \epsilon \in (0, 1)$, we have $\text{opt}_{1+\delta}(\mathcal{F}_{1+\epsilon}) = O(\min(\delta, \epsilon)^{-1})$.*

As a corollary, we have the following result, which was also a conjecture by [Geneson and Zhou \(2023\)](#).

Theorem 1.2. *For all $p > 1$ and $q > 1$, $\text{opt}_p(\mathcal{F}_q)$ is finite.*

Combined with the results from [Kimber and Long \(1995\)](#) that $\text{opt}_p(\mathcal{F}_q) = \infty$ whenever either $p = 1$ or $q = 1$, we now achieve a complete characterization of the values of (p, q) for which $\text{opt}_p(\mathcal{F}_q)$ is finite (i.e. precisely when $p, q > 1$), answering a question that has been open since the model of online learning was first defined for smooth functions by [Kimber and Long \(1995\)](#).

1.3. Online Learning of Polynomials

[Geneson and Zhou \(2023\)](#) first explored the online learning of special families of smooth functions with further imposed restrictions, such as polynomials. Placing extra restrictions on smooth functions is natural as most functions modeling real-life phenomena belong to families that have additional properties beyond merely than the smoothness constraints that were extensively studied. As such, this direction can provide new results that are more specifically applicable to real-life learning scenarios.

Carrying over the notation from [Geneson and Zhou \(2023\)](#), let $\mathcal{P}_q \subseteq \mathcal{F}_q$ denote the family of polynomials f such that $f \in \mathcal{F}_q$. We prove a conjecture from [Geneson and Zhou \(2023\)](#) in the affirmative, establishing polynomials in \mathcal{F}_q are as hard to learn as \mathcal{F}_q .

Theorem 1.3. *For all $p > 0$ and $q \geq 1$, we have $\text{opt}_p(\mathcal{P}_q) = \text{opt}_p(\mathcal{F}_q)$.*

1.4. Smooth Functions with Noisy Feedback

The online learning of functions with noisy feedback has previously been studied in the discrete setting for classifiers, in ([Cesa-Bianchi et al., 1996b](#); [Auer and Long, 1999](#); [Filmus et al., 2023](#)), and we extend it here to the context of smooth functions. Indeed, suppose that a certain phenomenon to be learned cannot be perfectly represented by any function within a class, and the goal is to simply find a function from the class that can model the phenomenon as well as possible. As such, these inaccuracies correspond to the noisy feedback made by the adversary in our model, which we informally call ‘lies’.

Assume the same learning format as the standard scenario. In the t^{th} trial, the algorithm is queried on the value of f at an input x_t , whereupon it outputs a prediction \hat{y}_t and the adversary subsequently reveals the actual value of $f(x_t)$. However, for a fixed positive integer $\eta \geq 1$, the adversary is allowed to lie about the value of $f(x_t)$ for up to η trials t . We define the error function similarly to our definition in the standard case, as the sum of the p^{th} powers of the raw errors of each trial $e_t = |\hat{y}_t - f(x_t)|$.

In the standard model of the learning of smooth functions, the error made by the learner on the first trial, e_0 , is not counted in the error function, as otherwise the adversary can simply set the function f identically equal to some constant C sufficiently far away from the learner’s prediction to generate an arbitrarily large error. As such, the learner needs to know the value of the function for at least one input to guarantee finite error on its first prediction. In this noisy version, where the learner is allowed to lie up to η times, getting feedback at one input is not sufficient to guarantee finite error on the next prediction for the learner, as the adversary can lie. There is no point in studying the worst-case error of the learner when the adversary can always force infinite error. As such, it is

necessary that we give the learner more initial rounds of feedback. We establish the following result about the number of initial rounds feedback necessary for the learner to guarantee finite error on its first real prediction.

Theorem 1.4. *For any integer $\eta \geq 1$, if incorrect feedback can be given up to η times, then at least $2\eta + 1$ initial rounds must be thrown out for the learner to guarantee finite error on its first prediction that counts toward the error evaluation.*

Accordingly, let the error function calculate the sum of the p^{th} powers of the raw errors starting from the $2\eta + 2^{\text{nd}}$ trial. In particular, for a fixed learning algorithm A operating on a finite sequence of inputs $\sigma = (x_0, x_1, \dots, x_m) \in [0, 1]^{m+1}$ and fixed function $f \in \mathcal{F}$, define the error function for the noisy model to be $L_{p,\eta}^{\text{nf}}(A, f, \sigma) = \sum_{t=2\eta+1}^m e_t^p = \sum_{t=2\eta+1}^m |\hat{y}_t - f(x_t)|^p$.

Similar to the definition of the standard model, we denote the worst-case error for a fixed algorithm A over a class \mathcal{F} of functions $f : [0, 1] \rightarrow \mathbb{R}$ as

$$L_{p,\eta}^{\text{nf}}(A, \mathcal{F}) = \sup_{f \in \mathcal{F}, \sigma \in \cup_{m \in \mathbb{Z}^+} [0, 1]^m} L_{p,\eta}^{\text{nf}}(A, f, \sigma).$$

Finally, we define the worst-case error for the best possible learning algorithm over a class \mathcal{F} , where the adversary is allowed to lie up to η times:

$$\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}) = \inf_A L_{p,\eta}^{\text{nf}}(A, \mathcal{F}).$$

We first establish a characterization for when $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q)$ is finite.

Theorem 1.5. *For any integer $\eta \geq 1$, the value of $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q)$ is finite if and only if $p > 1$ and $q > 1$.*

For $p, q \geq 2$ and any $\eta \geq 1$, we show that the noisy worst-case error is precisely on the order of η .

Theorem 1.6. *For any $\eta \geq 1, p, q \geq 2$ we have $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q) = \Theta(\eta)$.*

1.5. Order of Results

In Section 2, we work with the standard single-variable model, sketching the proof of Theorem 1.1 and its corollary Theorem 1.2. In Section 3, we study the online learning of smooth polynomials, establishing Theorem 1.3. In Section 4, we focus on the noisy learning scenario, defining its setup and establishing some preliminary bounds on $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q)$. In Section 5, we discuss open problems, conjectures, and future directions of research. The Appendix will contain the detailed proofs of several Lemmas used throughout the paper.

2. An Upper Bound on $\text{opt}_p(\mathcal{F}_q)$ for $p, q \in (1, 2)$

In this section, we prove that for all $\delta, \epsilon \in (0, 1)$, we have $\text{opt}_{1+\delta}(\mathcal{F}_{1+\epsilon}) = O(\min(\delta, \epsilon)^{-1})$. To prove this upper bound, we will use the LININT learning algorithm first defined by Kimber and Long (1995), which has previously been used to prove various upper bounds on $\text{opt}_p(\mathcal{F}_q)$ in Kimber and Long (1995), Long (2000), and Geneson and Zhou (2023). Specifically, we will first prove that $L_{1+\epsilon}(\text{LININT}, \mathcal{F}_{1+\epsilon}) = O(\epsilon^{-1})$ for $\epsilon \in (0, 1)$. We will then use this to bound $\text{opt}_{1+\delta}(\mathcal{F}_{1+\epsilon})$ for all $\epsilon, \delta \in (0, 1)$.

The main idea for the first part of our proof that $L_{1+\epsilon}(\text{LININT}, \mathcal{F}_{1+\epsilon}) = O(\epsilon^{-1})$ is similar to the main idea of previous proofs of upper bounds on $L_p(\text{LININT}, \mathcal{F}_q)$. Specifically, we will compare changes in $J_{1+\epsilon}[f_S]$ as new points are added to the input-output set S to powers of $|y - f_S(x)|$, the raw absolute errors generated by the corresponding rounds, culminating in Lemma 2.3. This approach is similar to that of Lemma 10 in Kimber and Long (1995) and Lemma 2.10 in Geneson and Zhou (2023). The main difference is that we will use novel inequality tactics, such as the binomial series expansion, to prove stronger inequalities than the inequalities previously proved. In the second part of our proof, we establish Lemma 2.5, a novel inequality that holds for an arbitrary number of rounds. Our result then follows upon combining Lemmas 2.3 and 2.5.

We emphasize where the difficulty of the regime $p, q \in (1, 2)$ lies. The learning algorithm LININT itself is unchanged from prior work (Kimber and Long, 1995; Long, 2000; Geneson and Zhou, 2023); the novelty is entirely in the analysis. Prior techniques were able to resolve $\text{opt}_p(\mathcal{F}_q)$ only when at least one of $p \geq 2$ and $q \geq 2$ held, because both the L_p loss and the q -action become much better behaved in that range, yielding clean error inequalities (for instance, $q \geq 2$ makes the relevant convexity arguments and second-order comparisons go through directly). When $p, q \in (1, 2)$, these inequalities degenerate, and no finiteness proof was previously known for any such pair. Our core contribution is the structural setup that bypasses these roadblocks: we introduce a tailored potential function $H_{p,S}$ that relates the smoothness constraint directly to the incurred error (used alongside the q -action J_q), and we use generalized binomial series expansions to establish the stronger inequalities (Lemmas 2.1 and 2.2) that drive the argument into the previously intractable regime. While the final algebra is elementary, it is precisely this conceptual reduction that makes a nearly 30-year-old problem tractable, and we expect these tools to be useful more broadly.

From this point onwards, we assume without loss of generality that the learning algorithm is never queried on the same input more than once, as the algorithm can always guarantee zero error upon being queried on the same input the second time. First, we prove a modification of the inequality that is Corollary 2.9 from Geneson and Zhou (2023). In particular, we specify that $0 < a \leq b < 1$ and expand the domain of x that the inequality works for from all $x \notin (a, b)$ to all $x \notin (-a, a)$, at the cost of weakening the inequality by a constant factor.

Lemma 2.1. *For reals $0 < a \leq b < 1$ such that $a + b \leq 1$, $q \in (1, 2)$, and any $|x| \geq a$, we have*

$$a \left| \frac{x}{a} + 1 \right|^q + b \left| \frac{x}{b} - 1 \right|^q - (a + b) \geq \frac{(q-1)|x|^q}{3}.$$

Proof Sketch. [See Appendix B.1 for full proof.] For any $x \in (-\infty, -a] \cup [b, \infty)$, the inequality directly follows from Corollary 2.9 from Geneson and Zhou (2023), so we only care about $x \in [a, b)$. Fix a, b , and set the LHS subtracted by RHS as $f(x)$. Using generalized binomial series, we first check that $f'(x) \geq 0$ for $x \geq a$, and then verify that $f(a) \geq 0$. \diamond

We now prove a modification (somewhat stronger version) of Lemma 2.7 from Geneson and Zhou (2023), while narrowing the x -domain of $(-a, b)$ to $(-a, a)$, as we dealt with the case of $x \in [a, b)$ in Lemma 2.1.

Lemma 2.2. *For reals $0 < a \leq b$, $q \in (1, 2)$, and $x \in (-a, a)$, we have*

$$a \left(1 + \frac{x}{a} \right)^q + b \left(1 - \frac{x}{b} \right)^q - (a + b) \geq \frac{q(q-1)}{3a} \cdot x^2.$$

Proof Sketch. [See Appendix B.2 for full proof.] Again, we expand the left side by generalized binomial series expansion, considering some casework for whether $\frac{x}{a} < 0$ or $\frac{x}{a} \geq 0$. The result follows from comparing terms of the series and the right side. \diamond

Combining the above two inequalities, we make the following key claim, which is essentially a strengthened and more specific version of Lemma 2.10 in Genson and Zhou (2023).

Lemma 2.3. *For a fixed $q \in (1, 2)$, a nonempty set $S = \{(u_i, v_i) : 1 \leq i \leq k\}$ of points in $[0, 1] \times \mathbb{R}$ with $u_i < u_{i+1}$ for each $1 \leq i \leq k - 1$, and another point $(x, y) \in [0, 1] \times \mathbb{R}$ such that $x \neq u_i$ for any i , we must have either*

$$J_q[f_{S \cup \{(x,y)\}}] - J_q[f_S] \geq \frac{q-1}{3} \cdot |y - f_S(x)|^q$$

or

$$J_q[f_{S \cup \{(x,y)\}}] - J_q[f_S] \geq \frac{(q-1)}{3|m|^{2-q} \cdot d} \cdot (y - f_S(x))^2,$$

where we let $d = \min_i |x - u_i|$ and $m = f'_S(x)$, the slope of the linear interpolation function at x .

Proof Sketch. [See Appendix B.3 for full proof.] Assuming $u_i < x < u_{i+1}$, we substitute $a = x - u_i$, $b = u_{i+1} - x$, $c = y - f_S(x)$, and do casework on whether the quantity $\frac{c}{m} \in (-a, a)$ or not, and apply Lemma 2.2 and Lemma 2.1, respectively. In particular, after substitution, $J_q[f_{S \cup \{(x,y)\}}] - J_q[f_S]$ turns out exactly in the form of the left side of Lemmas 2.2 and 2.1. \diamond

Now, we proceed to the second part of our proof, where we prove another bound, Lemma 2.5, which we will combine at the end with Lemma 2.3 using Hölder's Inequality to prove our desired result. To do this, we first make the following definition.

Definition 2.1. *For a fixed $p \geq 1$ and a set $S = \{(u_i, v_i) : 1 \leq i \leq k\}$ of at least two points in $[0, 1] \times \mathbb{R}$ with $u_i < u_{i+1}$ for each $1 \leq i \leq k - 1$, define*

$$H_{p,S} = \sum_{i=1}^{k-1} |v_{i+1} - v_i| (1 - |u_{i+1} - u_i|^{p-1}).$$

This seemingly contrived expression has the property of always being between 0 and 1, and being able to be re-written in several ways. Specifically, refer to A.3 and A.4 in the Appendix.

We prove that the variable $H_{p,S}$ never decreases when we add new points into set S , and prove a lower bound on the amount that $H_{p,S}$ increases by upon the addition of a new point into S in terms of p , the slope of f_S at the x -coordinate of the new point, and the closest x -coordinate difference between the new point and a point in S . Note that these are precisely the quantities involved in the second inequality in Lemma 2.3.

Lemma 2.4. *Consider a fixed real parameter $p \geq 1$, a set $S = \{(u_i, v_i) : 1 \leq i \leq k\}$ of at least two points in $[0, 1] \times \mathbb{R}$ with $u_1 < \dots < u_k$, and another $(x, y) \in [0, 1] \times \mathbb{R}$ with $x \neq u_i$ for any i . If we let $d = \min_{1 \leq i \leq k} |x - u_i|$ and let $m = f'_S(x)$, the slope of the linear interpolation of S at x , then we have*

$$H_{p,S \cup \{(x,y)\}} - H_{p,S} \geq (p-1)|m|d^p.$$

Proof Sketch. [See B.5 for detailed proof.] The proof is primarily inequality manipulation, where we invoke A.4 to rewrite the left side. In particular, we first establish a general inequality, Lemma B.4, which we apply to the specific value $\lambda = \frac{u_{i+1}-u_i}{d}$, where $u_i < x < u_{i+1}$, to get the result; we deal with when $x > u_n$ or $x < u_i$ separately. \diamond

Now, we establish the following key inequality.

Lemma 2.5. *Choose a sequence $(u_1, v_1), (u_2, v_2), \dots$ of points in $[0, 1] \times \mathbb{R}$ with $u_i \neq u_j$ for any $i \neq j$, and define the set $S_n = \{(u_i, v_i) : 1 \leq i \leq n\}$ for every positive integer n . For each $i > 1$, let $d_i = \min_{j < i} |u_i - u_j|$ and let $m_i = f'_{S_{i-1}}(u_i)$, the slope of the linear interpolation of the first $i - 1$ points of the sequence at u_i . If the inequality $J_1[f_{S_i}] \leq 1$, or equivalently $f_{S_i} \in \mathcal{F}_1$, holds for every positive integer i , then for any $p > 1$, we have*

$$\sum_{i=2}^{\infty} |m_i| d_i^p \leq \frac{1}{p-1}.$$

Proof Indeed, we have $\sum_{i=2}^{\infty} (p-1)|m_i|d_i^p = \sum_{i=3}^{\infty} (p-1)|m_i|d_i^p \leq \sum_{i=3}^{\infty} (H_{p,S_i} - H_{p,S_{i-1}}) \leq 1$ where the first step follows from $|m_2| = 0$, the second step follows from Lemma 2.4, and the last step follows from $0 \leq H_{p,S_i} \leq 1$ for any i , as $J_1[f_{S_i}] \leq 1$ for any i .

We now combine Lemma 2.5 with Lemma 2.3 to prove our desired upper bound on the error of the LININT learning algorithm when $(p, q) = (1 + \epsilon, 1 + \epsilon)$ for $\epsilon \in (0, 1)$.

Theorem 2.6. *For $\epsilon \in (0, 1)$, we have $L_{1+\epsilon}(\text{LININT}, \mathcal{F}_{1+\epsilon}) = O(\epsilon^{-1})$ for $\epsilon \in (0, 1)$.*

Proof Let $1 + \epsilon = q \in (1, 2)$ and suppose that $f \in \mathcal{F}_q$ is the function to be learned by LININT. Let $\sigma = (x_0, x_1, \dots, x_n)$ be the sequence of inputs, all different from each other, with $n \geq 1$. Similar to Geneson and Zhou (2023), we define $S_i = \{(x_0, f(x_0)), \dots, (x_n, f(x_n))\}$ for each $0 \leq i \leq n$, and let $\hat{y}_i \in \mathbb{R}$ for every $1 \leq i \leq n$ be the prediction of LININT corresponding to each x_i . Define the raw error of each round $e_i = |\hat{y}_i - f(x_i)|$, for each $1 \leq i \leq n$. Also, let $d_i = \min_{j < i} |x_i - x_j|$ and $m_i = f'_{S_{i-1}}(x_i)$ for each i . By Lemma 2.3, for each round $1 \leq i \leq n$, we have either

$$J_q[f_{S_i}] - J_q[f_{S_{i-1}}] \geq \frac{q-1}{3} \cdot e_i^q \quad \text{or} \quad J_q[f_{S_i}] - J_q[f_{S_{i-1}}] \geq \frac{q-1}{3|m_i|^{2-q} \cdot d_i} \cdot e_i^2.$$

Let set $I_1 \subseteq \{1, 2, \dots, n\}$ consist of all positive integers $1 \leq i \leq n$ such that e_i satisfy the first inequality. Analogously, let set $I_2 \subseteq \{1, 2, \dots, n\}$ consist of all positive integers $1 \leq i \leq n$ such that e_i, m_i, d_i satisfy the second inequality. Note that $I_1 \cup I_2 = \{1, 2, \dots, n\}$. Now, by Lemma A.2 and 2.3, we have

$$1 \geq J_q[f] \geq J_q[f_{S_n}] = \sum_{i=1}^n (J_q[f_{S_i}] - J_q[f_{S_{i-1}}]) \geq \sum_{i \in I_1} (J_q[f_{S_i}] - J_q[f_{S_{i-1}}]) \geq \frac{q-1}{3} \cdot \sum_{i \in I_1} e_i^q.$$

Dividing both sides, we get $\sum_{i \in I_1} e_i^q \leq \frac{3}{q-1}$. Similarly, $\sum_{i \in I_2} \frac{e_i^2}{|m_i|^{2-q} \cdot d_i} \leq \frac{3}{q-1}$. By Hölder's Inequality, $\sum_{i \in I_2} e_i^q \leq \left(\sum_{i \in I_2} \frac{e_i^2}{|m_i|^{2-q} \cdot d_i} \right)^{\frac{q}{2}} \left(\sum_{i \in I_2} |m_i|^q d_i^{\frac{q}{2-q}} \right)^{\frac{2-q}{2}} \leq \left(\frac{3}{q-1} \right)^{\frac{q}{2}} \left(\sum_{i \in I_2} |m_i|^q d_i^{\frac{q}{2-q}} \right)^{\frac{2-q}{2}}$.

Note also that $|m_i| \cdot d_i^{\frac{1}{2-q}} \leq |m_i| \cdot d_i$, as $\frac{1}{2-q} > 1$. We can verify that $|m_i| \cdot d_i \leq 1$ with a geometric

argument, omitted for space. As such, $|m_i| \cdot d_i^{\frac{1}{2-q}} \leq 1$ for each i , so

$$\sum_{i \in I_2} |m_i|^q d_i^{\frac{q}{2-q}} = \sum_{i \in I_2} \left(|m_i| d_i^{\frac{1}{2-q}} \right)^q \leq \sum_{i \in I_2} |m_i| d_i^{\frac{1}{2-q}} \leq \sum_{i=1}^n |m_i| d_i^{\frac{1}{2-q}} \leq 1 + \frac{1}{\frac{1}{2-q} - 1} = \frac{1}{q-1},$$

where the fourth step follows from Lemma 2.5, which is valid because $f_{S_i} \in \mathcal{F}_q \subseteq \mathcal{F}_1$ for all i .

Therefore, we have $\sum_{i \in I_2} e_i^q \leq \left(\frac{3}{q-1}\right)^{\frac{q}{2}} \left(\frac{1}{q-1}\right)^{\frac{2-q}{2}} \leq \frac{3}{q-1}$. Putting everything together, we have

$$L_q(\text{LININT}, f, \sigma) = \sum_{i=1}^n e_i^q \leq \sum_{i \in I_1} e_i^q + \sum_{i \in I_2} e_i^q \leq \frac{6}{q-1},$$

for any choice of $f \in \mathcal{F}_q$, any positive integer n , and any sequence of inputs $\sigma \in [0, 1]^{n+1}$. As such, we have $L_{1+\epsilon}(\text{LININT}, \mathcal{F}_{1+\epsilon}) \leq \frac{6}{\epsilon}$, for any $\epsilon \in (0, 1)$.

Recall from Geneson and Zhou (2023) that for any $q > 1$ and $p' > p > 1$, we have $L_{p'}(\text{LININT}, \mathcal{F}_q) \leq L_p(\text{LININT}, \mathcal{F}_q)$. Combining this with Theorem 2.6, we can get the following.

Theorem 2.7. *For $0 < \epsilon \leq \delta < 1$, we have $L_{1+\delta}(\text{LININT}, \mathcal{F}_{1+\epsilon}) = O(\epsilon^{-1})$ for $\epsilon \in (0, 1)$.*

As such, we have the following upper bound.

Theorem 2.8. *For $0 < \epsilon \leq \delta < 1$, we have $\text{opt}_{1+\delta}(\mathcal{F}_{1+\epsilon}) = O(\epsilon^{-1})$.*

Yet, note that for any $q' > q \geq 1$, $\text{opt}_p(\mathcal{F}_{q'}) \leq \text{opt}_p(\mathcal{F}_q)$ because $\mathcal{F}_{q'} \subseteq \mathcal{F}_q$. Thus, for all $0 < \delta \leq \epsilon < 1$, we have $\text{opt}_{1+\delta}(\mathcal{F}_{1+\epsilon}) \leq \text{opt}_{1+\delta}(\mathcal{F}_{1+\delta}) = O(\delta^{-1})$. Combining with Theorem 2.8, we have the following upper bound.

Theorem 2.9. *For $\delta, \epsilon \in (0, 1)$, we have $\text{opt}_{1+\delta}(\mathcal{F}_{1+\epsilon}) = O(\min(\delta, \epsilon)^{-1})$.*

Corollary 2.10. *For all $p > 1$ and $q > 1$, $\text{opt}_p(\mathcal{F}_q)$ is finite.*

As $\text{opt}_p(\mathcal{F}_q) = \infty$ whenever $p = 1$ or $q = 1$, we have completely classified the regions of $p, q \geq 1$ where $\text{opt}_p(\mathcal{F}_q)$ is finite, a problem open since the model of Kimber and Long (1995) was defined. Our result confirms a conjecture by Geneson and Zhou (2023).

Theorem 2.11. *The worst-case learning error $\text{opt}_p(\mathcal{F}_q)$ is finite if and only if $p, q > 1$.*

3. Online Learning of Polynomials

We prove the conjecture, raised by Geneson and Zhou (2023), that for all $p > 0$ and $q \geq 1$, $\text{opt}_p(\mathcal{P}_q) = \text{opt}_p(\mathcal{F}_q)$. Intuitively, this means that given a fixed q -action restriction on a smooth function, having the extra restriction of the function being a polynomial does not decrease the learner's error.

Our proof uses the following well-known result on the polynomial approximation of continuous real-valued functions.

Theorem 3.1 (Weierstrass Approximation Theorem). *For any continuous real-valued function f defined on a real interval $[a, b]$ and any $\epsilon > 0$, there exists a polynomial P such that for all $x \in [a, b]$, $|f(x) - P(x)| < \epsilon$.*

Using the Weierstrass Approximation Theorem, we prove that given any set of known points S , the adversary can find a polynomial P that is at most ϵ away from the points in S and has q -action bounded above by ϵ more than the q -action of f_S for any $\epsilon > 0$.

Lemma 3.2. *Given any $q \geq 1$, $\epsilon > 0$, and set $S = \{(u_i, v_i) : 1 \leq i \leq m\}$ of m points in $[0, 1] \times \mathbb{R}$ such that $u_1 < \dots < u_m$ and its corresponding linear interpolation function f_S , there exists a polynomial $P : [0, 1] \rightarrow \mathbb{R}$ such that $|P(u_i) - v_i| = |P(u_i) - f_S(u_i)| < \epsilon$ for all $1 \leq i \leq m$, and $J_q[P] < J_q[f_S] + \epsilon$.*

Proof Sketch. [See Appendix C.1 for full proof.] As we require not just P itself to be close enough to f_S , but also its q -action, we want the derivative of P to be close to that of f_S as well. Thus, we use Weierstrass's Theorem to construct a polynomial Q that approximates f'_S (or rather, a continuous analog of it), and then integrate Q to get P . \diamond

For a finite set of points $S \subseteq [0, 1] \times \mathbb{R}$, we now prove a fact about constructing a function that passes through all points in S by taking a weighted average of a special set of $2^{|S|}$ functions that each do not necessarily pass through the points in S . We will later combine this with Lemma 3.2 to implicitly construct a polynomial passing through all points in S by taking a weighted average of $2^{|S|}$ polynomials that each approximate but do not necessarily pass through the points in S . We use the convention that $[m] := \{1, 2, \dots, m\}$.

Lemma 3.3. *Suppose there is a set $S = \{(u_i, v_i) : 1 \leq i \leq m\}$ of m points in $[0, 1] \times \mathbb{R}$ such that $u_1 < \dots < u_m$. If we have 2^m functions $f_X : [0, 1] \times \mathbb{R}$, corresponding to each subset $X \subseteq [m]$, such that for any $X \subseteq [m]$, it holds for any integer $1 \leq i \leq m$ with $i \in X$ that $f_X(u_i) > v_i$, and it holds for any integer $1 \leq i \leq m$ with $i \notin X$ that $f_X(u_i) < v_i$, then there exists a weighted average of the functions $f \equiv \sum_{X \subseteq [m]} w_X f_X$, with $0 \leq w_X \leq 1$ for all $X \subseteq [m]$ and $\sum_{X \subseteq [m]} w_X = 1$, such that $f(u_i) = v_i$ for all $1 \leq i \leq m$.*

Proof Sketch. [See Appendix C.2 for full proof.] The result follows from a simple induction argument on m . Specifically, we first construct a weighted average aligning with the first $m - 1$ points but exceeding the last point; then, we construct a weighted average aligning with the first $m - 1$ points but exceeding the last point. Lastly, we take a weighted average of these functions. \diamond

We now establish that a weighted average of functions with q -action < 1 has q -action < 1 .

Lemma 3.4. *For any $q \geq 1$ and n continuous functions $f_1, f_2, \dots, f_n : [0, 1] \rightarrow \mathbb{R}$ with $J_q[f_i] < 1$ for all $1 \leq i \leq n$, then for any weighted average of the functions $f \equiv \sum_{i=1}^n w_i f_i$, with $0 \leq w_i \leq 1$ for all $1 \leq i \leq n$ and $\sum_{i=1}^n w_i = 1$, we have $J_q[f] < 1$.*

Proof From Minkowski's Inequality, $J_q[f] \leq \sum_{i=1}^n J_q[w_i f_i] = \sum_{i=1}^n w_i J_q[f_i] < \sum_{i=1}^n w_i = 1$.

Putting everything together, we have the following key result.

Lemma 3.5. *Given a set $S = \{(u_i, v_i) : 1 \leq i \leq m\}$ of m points in $[0, 1] \times \mathbb{R}$ with $u_1 < \dots < u_m$ and any $q \geq 1$, if there exists an absolutely continuous function $f : [0, 1] \rightarrow \mathbb{R}$ with $J_q[f] < 1$ such that $f(u_i) = v_i$ for all $1 \leq i \leq m$, then there exists a polynomial $P : [0, 1] \rightarrow \mathbb{R}$ with $J_q[P] < 1$ such that $P(u_i) = v_i$ for all $1 \leq i \leq m$.*

Proof Sketch. [See Appendix C.3 for full proof.] Using Lemma 3.2, we first create 2^m polynomial approximations of f_S with special conditions (i.e. the conditions of Lemma 3.3). Then, taking a weighted average of these with Lemma 3.3 with some analysis, we achieve our result. \diamond

Now, for any $q \geq 1$, let \mathcal{F}'_q denote the class of functions $f : [0, 1] \rightarrow \mathbb{R}$ such that $J_q[f] < 1$. As such, $\mathcal{F}'_q \subseteq \mathcal{F}_q$. Similarly, let \mathcal{P}'_q denote the class of polynomials $P : [0, 1] \rightarrow \mathbb{R}$ such that $J_q[P] < 1$. As such, $\mathcal{P}'_q \subseteq \mathcal{P}_q$ and $\mathcal{P}'_q \subseteq \mathcal{F}'_q$. We now establish a result that looks very close to our desired final result.

Lemma 3.6. *For all $p > 0$ and $q \geq 1$, we have $\text{opt}_p(\mathcal{P}'_q) = \text{opt}_p(\mathcal{F}'_q)$.*

Proof Note that by Lemma 3.5, given any finite set S of points in $[0, 1] \times \mathbb{R}$, there exists a smooth polynomial $P \in \mathcal{P}'_q$ passing through all the points in S if and only if there exists a general smooth function $f \in \mathcal{F}'_q$ passing through all the points in S .

We claim that this means that the adversary can replicate any strategy that it uses in the learning scenario of \mathcal{F}'_q to the learning scenario of \mathcal{P}'_q against the learner, and vice versa. Indeed, for any sequence of points $(x_0, y_0), \dots, (x_m, y_m)$ that the adversary reveals in the learning scenario of \mathcal{F}'_q , the same sequence of points can be revealed in the learning scenario of \mathcal{P}'_q (and vice versa as well), as the existence of a general function $f \in \mathcal{F}'_q$ passing through the sequence of points implies the existence of a polynomial $P \in \mathcal{P}'_q$ passing through the same sequence of points (and vice versa as well, trivially). So, the adversary can always employ the same strategy against the learner, based on its predictions. Thus, we have $\text{opt}_p(\mathcal{P}'_q) = \text{opt}_p(\mathcal{F}'_q)$.

Now, we prove that the following equality holds.

Lemma 3.7. *For all $p > 0$ and $q \geq 1$, we have $\text{opt}_p(\mathcal{F}'_q) = \text{opt}_p(\mathcal{F}_q)$.*

Proof For a real number $c > 0$, let the family $c\mathcal{F}_q$ contain precisely the functions $g : [0, 1] \rightarrow \mathbb{R}$ such that $g \equiv cf$ for some $f \in \mathcal{F}_q$; note that this is equivalent to $J_q[g] \leq c$.

We claim that $\text{opt}_p(c\mathcal{F}_q) = c^p \text{opt}_p(\mathcal{F}_q)$ for any $c > 0$. This follows from a scaling argument, the specifics of which we omit for space. For all $0 < c < 1$, note now that $c^p \text{opt}_p(\mathcal{F}_q) = \text{opt}_p(c\mathcal{F}_q) \leq \text{opt}_p(\mathcal{F}'_q)$, as $c\mathcal{F}_q \in \mathcal{F}'_q$. Taking $c \rightarrow 1$, we can see that $\text{opt}_p(\mathcal{F}'_q)$ cannot be any less than $\text{opt}_p(\mathcal{F}_q)$. Yet, as $\mathcal{F}'_q \in \mathcal{F}_q$, we also have $\text{opt}_p(\mathcal{F}'_q) \leq \text{opt}_p(\mathcal{F}_q)$. We thus have $\text{opt}_p(\mathcal{F}'_q) = \text{opt}_p(\mathcal{F}_q)$.

Using the same idea, we can get that similarly, $\text{opt}_p(\mathcal{P}'_q) = \text{opt}_p(\mathcal{P}_q)$, for all $p > 0$ and $q \geq 1$. Putting everything together, we get our result.

Theorem 3.8. *For any $p > 0$ and $q \geq 1$, we have $\text{opt}_p(\mathcal{P}_q) = \text{opt}_p(\mathcal{F}_q)$.*

4. Noisy Online Learning of Smooth Functions

We start with the proof of Theorem 1.4, which we split into a lower bound and an upper bound.

Theorem 4.1. *For any integer $\eta \geq 1$, if incorrect feedback can be given up to η times, then at least $2\eta + 1$ initial rounds must be thrown out for the learner to guarantee finite error on its first prediction that counts toward the error evaluation.*

Proof Sketch. [See Appendix D.1 for full proof.] We establish that 2η rounds are insufficient, with the adversary querying a single point 2η times and giving two different answers each η times. Then, we prove $2\eta + 1$ rounds is sufficient, by establishing a strategy to bound the range of f entirely. \diamond

This motivates our definition of the worst-case error $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q)$ in Section 1.4 to only start counting from the $2\eta + 2^{\text{th}}$ round. We now establish the following fundamental result.

Theorem 4.2. *For any integer $\eta \geq 1$, the value of $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q)$ is finite if and only if $p, q > 1$.*

Proof Clearly, if $p = 1$ or $q = 1$, $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q) \geq \text{opt}_p(\mathcal{F}_q) = \infty$. When $p, q > 1$, we can replicate the optimal strategy in the standard scenario for the same p, q , as if the adversary cannot lie, until something wrong occurs, indicating a lie—the total error exceeds $\text{opt}_p(\mathcal{F}_q)$ (which is finite by Theorem 2.11), or the learner gives feedback not in the allocated range $[v_{\eta+1} - 1, v_{\eta+1} + 1]$. Then, the learner forgets previous feedback, and starts all over again, and repeats. As this can happen at most $\eta + 1$ times, and we can bound the error in every stage, the total error is bounded.

We now prove Theorem 1.6. Our proof is split into two parts: an upper and lower bound. The upper bound uses a similar, but more sophisticated strategy as the previous result.

Lemma 4.3. *For any $\eta \geq 1$, $p, q \geq 2$ we have $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q) \leq 12\eta + 6$.*

Proof Sketch. [See Appendix D.2 for full proof.] Again, let the learner mimic the strategy in the standard case, until it the perceived error exceeds $\text{opt}_p(\mathcal{F}_q)$, then restart and repeat, which happens at most $\eta + 1$ times. ◇

Lemma 4.4. *For any $\eta \geq 1$, $p, q \geq 2$ we have $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q) \geq 2\eta + 1$.*

Proof Sketch. [See Appendix D.3 for full proof.] The adversary only reveals $f(0) = 0$ in the initial rounds, and then repeatedly queries the learner on $f(1)$, revealing that it is 1 and -1 each η times. Finally, set the true value to whichever costs more error. ◇

Combining the bounds, we obtain a precise bound on $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q)$ for any $p, q \geq 2$ and $\eta \geq 1$.

Theorem 4.5. *For any $\eta \geq 1$, $p, q \geq 2$, we have that $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q) = \Theta(\eta)$.*

5. Discussion and Future Work

Matching rates and the minimax characterization. With the results of this paper, we have either established upper bounds on $\text{opt}_p(\mathcal{F}_q)$ or proved that it is infinite for all $p, q \geq 1$, a problem that has been open since the model of online learning was first defined for smooth functions by Kimber and Long (1995). This completes the characterization of *when* $\text{opt}_p(\mathcal{F}_q)$ is finite, but not of its exact rate. Our bound $\text{opt}_{1+\delta}(\mathcal{F}_{1+\epsilon}) = O(\min(\delta, \epsilon)^{-1})$ is not known to be tight: we provide no matching lower bound in the regime $p, q \in (1, 2)$, and the rate does not connect continuously to the known boundary behavior, such as the $\Theta(\epsilon^{-1/2})$ rate of Geneson and Zhou (2023) as $q \rightarrow 2$. A natural and important next step is therefore to establish lower and upper bounds on $\text{opt}_p(\mathcal{F}_q)$ that match up to a constant factor for every $p, q \geq 1$, yielding the full minimax rate. It is plausible that a fundamentally different algorithmic approach, rather than LININT, is needed to obtain sharp rates throughout the interior regime.

Connections to generic online-learning techniques. A natural question is whether our bounds could instead be obtained from generic online-learning machinery, such as Follow the Regularized Leader (FTRL), Mirror Descent, the Aggregating Algorithm, or sequential covering arguments. We believe they cannot, at least not directly, and that this highlights what is distinctive about the setting. The framework of Kimber and Long (1995) measures the *absolute* cumulative error $\sum_t |\hat{y}_t - f(x_t)|^p$ under realizability, a notion strictly stronger than agnostic regret, which measures only the gap to the

best fixed function in hindsight. Generic algorithms such as FTRL and Mirror Descent are designed to guarantee regret that is sublinear in the horizon T , whereas we require a bound that is finite independent of T . Obtaining such a T -independent absolute-error bound via FTRL would require strong convexity of the loss, which the L_p loss lacks for $p \in (1, 2)$; the Aggregating Algorithm of Vovk similarly requires exp-concavity, which the L_p loss also lacks; and sequential covering constructions are highly non-constructive and yield loose constants tailored to regret rather than absolute error. By contrast, LININT is an extremely simple algorithm—it runs in $O(1)$ time per step—and yet constructively attains finite absolute error in *every* regime where finite error is possible, including the interior regime where these standard methods do not apply. We regard it as a feature of the problem that such an elementary algorithm suffices.

Beyond realizability: the agnostic setting. Our noisy model retains realizability up to a bounded number η of corruptions. If realizability is dropped entirely, finite absolute error becomes impossible: an adversary can force a constant error at every step, regardless of the learner. The natural formulation in that case is the classical regret

$$\sum_{t=1}^T |\hat{y}_t - y_t|^p - \inf_{f \in \mathcal{F}_q} \sum_{t=1}^T |f(x_t) - y_t|^p,$$

measured against the best function in \mathcal{F}_q in hindsight. Bounding this quantity independently of T is difficult for the same structural reasons described above, and even characterizing the optimal minimax regret under fully agnostic feedback appears to require constructing sequential covering nets for \mathcal{F}_q . We leave the transition from finite absolute-error bounds to agnostic regret as an interesting direction for future work.

Other function classes and higher dimensions. Having studied the online learning of polynomials and established that polynomials in \mathcal{F}_q are not any easier to learn than general functions in \mathcal{F}_q , it remains to investigate other special subsets of \mathcal{F}_q . In particular, we hope to generalize the ideas for our proof that $\text{opt}_p(\mathcal{P}_q) = \text{opt}_p(\mathcal{F}_q)$ to other special subsets \mathcal{A}_q of \mathcal{F}_q . Indeed, the crux of our proof for \mathcal{P}_q was a connection with the Weierstrass Approximation Theorem, which allowed us to use polynomials to uniformly approximate any continuous real-valued function. Yet, there exists a generalization of the Weierstrass Approximation Theorem, the Stone–Weierstrass Theorem, offering a condition under which general subalgebras uniformly approximate continuous functions. As such, the Stone–Weierstrass Theorem may be a gateway to proving that more special subsets \mathcal{A}_q of \mathcal{F}_q are hard to learn. We make conjectures about two specific subsets: it is not any easier to learn sums of exponential functions in \mathcal{F}_q , or trigonometric polynomials in \mathcal{F}_q , than \mathcal{F}_q in general. A complementary direction concerns the geometry of the underlying class. Since \mathcal{F}_q coincides with the unit ball of the Sobolev space $W^{1,q}([0, 1])$, it is natural to ask whether our characterization extends to richer smoothness classes, such as higher-order Sobolev balls or Besov spaces. Finally, the problem here is restricted to single-variable, real-valued functions on $[0, 1]$; understanding how the characterization of finiteness and the corresponding rates extend to d -dimensional inputs $f : [0, 1]^d \rightarrow \mathbb{R}$, as studied for related models in [Barron \(1991\)](#); [Härdle \(1991\)](#), is an important open direction. In the noisy learning scenario, we have yet to obtain precise values of $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q)$ at any particular values of p, q, η , which is itself an interesting research direction, alongside finding upper bounds for when at least one of $p < 2$ and $q < 2$.

Other noise and feedback models. It would also be interesting to carry over more modes of adversary feedback previously studied in the context of classifiers to the context of smooth functions, like how we defined the noisy learning of smooth functions. Our noisy model bounds the *number* of corrupted rounds by η ; alternative and arguably more realistic models include additive (e.g. stochastic or bounded) noise on each revealed value—where one would measure an excess risk rather than absolute error—or corruption bounded by a *cumulative* budget across all rounds rather than a fixed count. A further natural relaxation, in the spirit of approximate learning, is to ask only that the learner recover the target to within some accuracy Δ given Δ -accurate feedback, rather than exactly; characterizing the dependence of the achievable error on Δ is a promising direction. Other previously studied feedback modes that could be transported to the smooth-function setting include delayed or ambiguous reinforcement and feedback on the absolute error of a prediction, to name a few.

Acknowledgments

I am extremely grateful to Dr. Jesse Geneson, for introducing me to online learning, and for his continued patience. I'm also thankful for MIT PRIMES-USA for providing this amazing research opportunity.

References

- D. Angluin, Queries and concept learning. *Machine Learning* **2** (1988) 319–342.
- P. Auer and P.M. Long. Structural results about on-line learning models with and without queries. *Machine Learning* **36** (1999) 147–181.
- A. Barron, Approximation and estimation bounds for artificial neural networks. *Workshop on Computational Learning Theory* (1991).
- N. Cesa-Bianchi, P.M. Long, and M.K. Warmuth, Worst-case quadratic loss bounds for prediction using linear functions and gradient descent. *IEEE Transactions on Neural Networks* **7** (1996) 604–619.
- N. Cesa-Bianchi, Y. Freund, D. Helmbold, and M. Warmuth. Online prediction and conversion strategies. *Machine Learning* **25** (1996) 71–110.
- V. Faber and J. Mycielski, Applications of learning theorems. *Fundamenta Informaticae* **15** (1991) 145–167.
- Y. Filmus, S. Hanneke, I. Mehalal, and S. Moran. Optimal prediction using expert advice and randomized Littlestone dimension. *Proceedings of the Thirty-Sixth Conference on Learning Theory* (2023) 773–836.
- R. Feng, J. Geneson, A. Lee, and E. Slettnes, Sharp bounds on the price of bandit feedback for several models of mistake-bounded online learning. *Theoretical Computer Science* **965** (2023) 113980.
- J. Geneson and E. Zhou. Online learning of smooth functions. *Theoretical Computer Science* **979C** (2023) 114203.
- W. Härdle, Smoothing techniques. Springer Verlag (1991).
- D. Kimber and P. M. Long, On-line learning of smooth functions of a single variable. *Theoretical Computer Science* **148** (1995) 141–156.
- N. Littlestone, Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm. *Machine Learning* **2** (1988) 285–318.
- N. Littlestone and M.K. Warmuth, The weighted majority algorithm. *Proceedings of the 30th Annual Symposium on the Foundations of Computer Science* (1989) 256–261.
- P. M. Long, Improved bounds about on-line learning of smooth functions of a single variable. *Theoretical Computer Science* **241** (2000) 25–35.
- J. Mycielski, A learning algorithm for linear operators. *Proceedings of the American Mathematical Society* **103** (1988) 547–550.

Appendix A. Facts used in paper

We state two facts regarding f_S that we invoked earlier. The first fact is a fact about the 1-action of any linear interpolation function f_S , which follows from the q -action formula.

Lemma A.1. *Given a set $S = \{(u_i, v_i) : 1 \leq i \leq m\}$, with $(u_i, v_i) \in [0, 1] \times \mathbb{R}$ for each $1 \leq i \leq m$ and $u_1 < \dots < u_m$, one useful fact about f_S is that*

$$J_1[f_S] = \sum_{i=1}^{m-1} \left(\int_{u_i}^{u_{i+1}} \left| \frac{v_{i+1} - v_i}{u_{i+1} - u_i} \right| dx \right) = \sum_{i=1}^{m-1} |v_{i+1} - v_i|.$$

Another useful fact about f_S is that it has the minimum q -action out of all functions f passing through all points in S :

Lemma A.2 (Kimber and Long (1995), Geneson and Zhou (2023)). *Let $S = \{(u_1, v_1), \dots, (u_m, v_m)\}$ be a set of m points with $(u_i, v_i) \in [0, 1] \times \mathbb{R}$ for each i , such that $u_1 < \dots < u_m$. Then, for any $q \geq 1$ and any absolutely continuous function $f : [0, 1] \rightarrow \mathbb{R}$ with $f(u_i) = v_i$ for $1 \leq i \leq m$, we have $J_q[f] \geq J_q[f_S]$.*

We also state two facts about $H_{p,S}$, which we use in the paper.

Lemma A.3. *For any $p \geq 1$ and set S of points in $[0, 1] \times \mathbb{R}$, whenever $f_S \in \mathcal{F}_1$, we have $0 \leq H_{p,S} \leq 1$.*

Proof Indeed, whenever $J_1[f_S] \leq 1$, we have

$$H_{p,S} = J_1[f_S] - \left(\sum_{i=1}^{k-1} |m_i| |u_{i+1} - u_i|^p \right) \leq J_1[f_S] \leq 1$$

and

$$H_{p,S} = J_1[f_S] - \left(\sum_{i=1}^{k-1} |m_i| |u_{i+1} - u_i|^p \right) \geq J_1[f_S] - \left(\sum_{i=1}^{k-1} |m_i| \right) = 0.$$

Moreover, $H_{p,S}$, can be rewritten in different ways, which we also use in our proofs.

Lemma A.4. *For any $p \geq 1$ and a set $S = \{(u_i, v_i) : 1 \leq i \leq k\}$ of points in $[0, 1] \times \mathbb{R}$ with $u_i < u_{i+1}$ for each $1 \leq i \leq k-1$, if we let $m_i = \frac{v_{i+1} - v_i}{u_{i+1} - u_i}$ be the slope of the line segment between points at u_i and u_{i+1} in f_S , we have that*

$$H_{p,S} = \sum_{i=1}^{k-1} (|v_{i+1} - v_i| - |m_i| |u_{i+1} - u_i|^p) = J_1[f_S] - \left(\sum_{i=1}^{k-1} |m_i| |u_{i+1} - u_i|^p \right).$$

Appendix B. Proofs for Section 2

We provide detailed proofs of the lemmas used in Section 3. First, we establish the two inequalities 2.1 and 2.2.

Lemma B.1. *For reals $0 < a \leq b < 1$ such that $a + b \leq 1$, $q \in (1, 2)$, and any $|x| \geq a$, we have*

$$a \left| \frac{x}{a} + 1 \right|^q + b \left| \frac{x}{b} - 1 \right|^q - (a + b) \geq \frac{(q-1)|x|^q}{3}.$$

Proof Note that for all $x \in (-\infty, -a] \cup [b, \infty)$, the inequality directly follows from Corollary 2.9 from [Geneson and Zhou \(2023\)](#). Therefore, we only have to deal with the case where $x \in [a, b)$. For fixed a and b , let

$$\begin{aligned} f(x) &= a \left| \frac{x}{a} + 1 \right|^q + b \left| \frac{x}{b} - 1 \right|^q - (a+b) - \frac{(q-1)|x|^q}{3} \\ &= a \left(\frac{x}{a} + 1 \right)^q + b \left(1 - \frac{x}{b} \right)^q - (a+b) - \frac{(q-1)x^q}{3} \end{aligned}$$

when $x \in [a, b]$. We prove that $f'(x) > 0$ for all $x \geq a$, so that it suffices to check that $f(a) \geq 0$. Indeed,

$$\begin{aligned} f'(x) &= q \left(\frac{x}{a} + 1 \right)^{q-1} - q \left(1 - \frac{x}{b} \right)^{q-1} - \frac{q(q-1)x^{q-1}}{3} \\ &\geq q \cdot 2^{q-1} - q - \frac{q(q-1)}{3} = q \left(2^{q-1} - 1 - \frac{q-1}{3} \right) > 0 \end{aligned}$$

for $q \in (1, 2)$, when $a \leq x \leq b$. Consider the function $g(x) = b \left(1 - \frac{x}{b} \right)^q$. Note that $g'(x) = -q \left(1 - \frac{x}{b} \right)^{q-1} \geq -q$ when $x \in [0, b]$. As $g(0) = b$, we have $g(x) \geq b - qx$ for all $x \in [0, b]$. As such, $b \left(1 - \frac{x}{b} \right)^q = g(x) \geq b - qa$. Now, we have

$$\begin{aligned} f(a) &= a \cdot 2^q + b \left(1 - \frac{a}{b} \right)^q - (a+b) - \frac{(q-1)a^q}{3} \geq a \cdot 2^q + (b - qa) - (a+b) - \frac{(q-1)a^q}{3} \\ &\geq a \cdot 2^q + (b - qa) - (a+b) - \frac{(q-1)a}{3} = a(2^q - q - 1) - \frac{(q-1)a}{3} \\ &\geq a \left(\frac{q-1}{3} \right) - \frac{(q-1)a}{3} = 0. \end{aligned}$$

Hence, we have that $f(x) \geq 0$ for all $x \in [a, b)$, as desired.

Lemma B.2. For reals $0 < a \leq b$, $q \in (1, 2)$, and $x \in (-a, a)$, we have

$$a \left(1 + \frac{x}{a} \right)^q + b \left(1 - \frac{x}{b} \right)^q - (a+b) \geq \frac{q(q-1)}{3a} \cdot x^2.$$

Proof We use the generalized binomial series expansion of $\left(1 + \frac{x}{a} \right)^q$ and $\left(1 - \frac{x}{b} \right)^q$ to rewrite the expression on the left hand side, which we will call Δ for convenience. Doing this is valid as by assumption $\left| \frac{x}{a} \right| < 1$ and $\left| \frac{x}{b} \right| < 1$, so the series expansions converge correctly. Then,

$$\begin{aligned} \Delta &= a \left(\sum_{k=0}^{\infty} \binom{q}{k} \left(\frac{x}{a} \right)^k \right) + b \left(\sum_{k=0}^{\infty} \binom{q}{k} \left(-\frac{x}{b} \right)^k \right) - (a+b) \\ &= a \left(\sum_{k=1}^{\infty} \binom{q}{k} \left(\frac{x}{a} \right)^k \right) + b \left(\sum_{k=1}^{\infty} \binom{q}{k} \left(-\frac{x}{b} \right)^k \right) \\ &= qx + a \left(\sum_{k=2}^{\infty} \binom{q}{k} \left(\frac{x}{a} \right)^k \right) - qx + b \left(\sum_{k=2}^{\infty} \binom{q}{k} \left(-\frac{x}{b} \right)^k \right) \\ &= a \left(\sum_{k=2}^{\infty} \binom{q}{k} \left(\frac{x}{a} \right)^k \right) + b \left(\sum_{k=2}^{\infty} \binom{q}{k} \left(-\frac{x}{b} \right)^k \right). \end{aligned}$$

We claim that $\left(\sum_{k=2}^{\infty} \binom{q}{k} \left(\frac{x}{a}\right)^k\right) \geq \frac{q(q-1)}{3} \left(\frac{x}{a}\right)^2$. The inequality is trivial if $\frac{x}{a} = 0$. If $\frac{x}{a} < 0$, note that every term of the sum is nonnegative, because $\binom{q}{k}$ is negative precisely when $k \geq 2$ is odd, from the condition that $q \in (1, 2)$, and $\left(\frac{x}{a}\right)^k$ is also negative precisely when k is odd. Therefore, we have $\left(\sum_{k=2}^{\infty} \binom{q}{k} \left(\frac{x}{a}\right)^k\right) \geq \frac{q(q-1)}{2} \left(\frac{x}{a}\right)^2$ clearly. If $\frac{x}{a} > 0$, note that $\binom{q}{k} \left(\frac{x}{a}\right)^k > 0$ when $k \geq 2$ is even, and $\binom{q}{k} \left(\frac{x}{a}\right)^k < 0$ when $k \geq 2$ is odd. Furthermore, for any $k \geq 2$, we have $\left|\binom{q}{k} \left(\frac{x}{a}\right)^k\right| > \left|\binom{q}{k+1} \left(\frac{x}{a}\right)^{k+1}\right|$ because $\left|\binom{q}{k+1}\right| = \left|\binom{q}{k} \cdot \frac{q-k}{k+1}\right| = \left|\binom{q}{k}\right| \cdot \left|\frac{k-q}{k+1}\right| < \left|\binom{q}{k}\right|$ when $q \in (1, 2)$ and $\left|\left(\frac{x}{a}\right)^{k+1}\right| < \left|\left(\frac{x}{a}\right)^k\right|$. As such, because $\left(\sum_{k=2}^{\infty} \binom{q}{k} \left(\frac{x}{a}\right)^k\right)$ is an alternating series with the magnitude of summands decreasing, the sum is bounded below by

$$\binom{q}{2} \left(\frac{x}{a}\right)^2 + \binom{q}{3} \left(\frac{x}{a}\right)^3 > \left(\binom{q}{2} + \binom{q}{3}\right) \left(\frac{x}{a}\right)^2 \geq \frac{2}{3} \binom{q}{2} \left(\frac{x}{a}\right)^2 \geq \frac{q(q-1)}{3} \left(\frac{x}{a}\right)^2$$

when $q \in (1, 2)$, as claimed.

Similarly, we can conclude that $\left(\sum_{k=2}^{\infty} \binom{q}{k} \left(-\frac{x}{b}\right)^k\right) \geq \frac{q(q-1)}{3} \left(-\frac{x}{b}\right)^2$, which is nonnegative. Using these inequalities, we get the lower bound on our desired expression:

$$\Delta \geq a \left(\frac{q(q-1)}{3} \left(\frac{x}{a}\right)^2\right) + b \cdot 0 \geq \frac{q(q-1)}{3a} \cdot x^2.$$

Now, we combine the inequalities to establish Theorem 2.3.

Lemma B.3. *For a fixed $q \in (1, 2)$, a nonempty set $S = \{(u_i, v_i) : 1 \leq i \leq k\}$ of points in $[0, 1] \times \mathbb{R}$ with $u_i < u_{i+1}$ for each $1 \leq i \leq k-1$, and another point $(x, y) \in [0, 1] \times \mathbb{R}$ such that $x \neq u_i$ for any i , we must have either*

$$J_q[f_{S \cup \{(x,y)\}}] - J_q[f_S] \geq \frac{q-1}{3} \cdot |y - f_S(x)|^q$$

or

$$J_q[f_{S \cup \{(x,y)\}}] - J_q[f_S] \geq \frac{(q-1)}{3|m|^{2-q} \cdot d} \cdot (y - f_S(x))^2,$$

where we let $d = \min_i |x - u_i|$ and $m = f'_S(x)$, the slope of the linear interpolation function at x .

Proof First, if $x < u_1$, as established in the proof of Lemma 2.10 in Geneson and Zhou (2023),

$$J_q[f_{S \cup \{(x,y)\}}] - J_q[f_S] = (u_1 - x) \left| \frac{v_1 - y}{u_1 - x} \right|^q.$$

Because $|u_1 - x| \leq 1$, the above expression is at least $|v_1 - y|^q = |y - f_S(x)|^q \geq (q-1)|y - f_S(x)|^q$, hence satisfying the first inequality. The case of $x > u_k$ is similar.

Meanwhile, if $u_i < x < u_{i+1}$ for some integer $1 \leq i \leq k-1$, substituting $a = x - u_i$, $b = u_{i+1} - x$, $c = y - f_S(x)$, and $m = f'_S(x) = \frac{v_{i+1} - v_i}{u_{i+1} - u_i} = \frac{f_S(x) - v_i}{a} = \frac{v_{i+1} - f_S(x)}{b}$, we have

$$J_q[f_{S \cup \{(x,y)\}}] - J_q[f_S] = |m|^q \left(a \left| 1 + \frac{c}{ma} \right|^q + b \left| 1 - \frac{c}{mb} \right|^q - (a+b) \right),$$

as established in the proof of Lemma 2.10 in [Geneson and Zhou \(2023\)](#). Without loss of generality, assume that $a = |x - u_i| \leq |x - u_{i+1}| = b$, as the other case follows from symmetry. Note that $a = d = \min_i |x - u_i|$. If $\frac{c}{m} \notin (-a, a)$, applying Lemma 2.1, we have

$$J_q [f_{S \cup \{(x,y)\}}] - J_q[f_S] \geq |m|^q \left(\frac{q-1}{3} \left| \frac{c}{m} \right|^q \right) = \frac{q-1}{3} \cdot |c|^q = \frac{q-1}{3} \cdot |y - f_S(x)|^q,$$

satisfying the first inequality. On the other hand, if $\frac{c}{m} \in (-a, a)$, applying Lemma 2.2, we have

$$\begin{aligned} J_q [f_{S \cup \{(x,y)\}}] - J_q[f_S] &\geq |m|^q \left(\frac{q(q-1)}{3a} \cdot \left| \frac{c}{m} \right|^2 \right) \\ &\geq |m|^q \left(\frac{(q-1)}{3d} \cdot \left| \frac{c}{m} \right|^2 \right) = \frac{q-1}{3|m|^{2-q} \cdot d} \cdot (y - f_S(x))^2. \end{aligned}$$

As a result, in the case where $\frac{c}{m} \in (-a, a)$, the second inequality is satisfied. Hence, in any case, either the first or the second inequality needs to hold.

Before we proceed with the proof of Theorem 2.4, we first establish the following Lemma.

Lemma B.4. *For all reals $p > 1$ and $x \geq 2$, we have*

$$x^p - (x-1)^{p-1} \cdot x \geq p-1.$$

Proof We consider two cases, dependent on whether $p \geq 2$ or $1 < p < 2$. If $p \geq 2$, then for all $x \geq 2$,

$$\begin{aligned} x^p - (x-1)^{p-1} \cdot x &= x(x^{p-1} - (x-1)^{p-1}) \geq x(2^{p-1} - 1^{p-1}) \\ &\geq 2 \cdot (2^{p-1} - 1^{p-1}) = 2^p - 2 \geq p-1, \end{aligned}$$

where the second step follows from the function $f(x) = x^{p-1} - (x-1)^{p-1}$ being non-decreasing for $x \geq 2$, which can be verified by differentiation. If $1 < p < 2$, then for all $x \geq 2$, we have

$$\begin{aligned} x^p - (x-1)^{p-1} \cdot x &= x^p \left(1 - \left(1 - \frac{1}{x}\right)^{p-1} \right) = x^p \left(1 - \left(\sum_{k=0}^{\infty} \binom{p-1}{k} \left(-\frac{1}{x}\right)^k \right) \right) \\ &= x^p \left(\sum_{k=1}^{\infty} (-1)^k \binom{p-1}{k} \left(-\frac{1}{x}\right)^k \right) = x^p \left(\sum_{k=1}^{\infty} (-1)^{k+1} \binom{p-1}{k} \left(\frac{1}{x}\right)^k \right) \end{aligned}$$

using the generalized binomial notation. The binomial series expansion in the second step is valid as $\left|\frac{1}{x}\right| < 1$. Furthermore, as $0 < p-1 < 1$, $\binom{p-1}{k}$ is negative for all even positive integers k and positive for all odd positive integers k , so the expression $(-1)^{k+1} \binom{p-1}{k} \left(\frac{1}{x}\right)^k$ is positive for all positive integers k . As such, our expression is bounded below by

$$x^p \left(\binom{p-1}{1} \cdot \frac{1}{x} \right) = x^{p-1}(p-1) \geq p-1$$

when $x \geq 2$.

Now, we proceed with the proof of Lemma 2.4

Lemma B.5. Consider a fixed real parameter $p \geq 1$, a set $S = \{(u_i, v_i) : 1 \leq i \leq k\}$ of at least two points in $[0, 1] \times \mathbb{R}$ with $u_i < u_{i+1}$ for each $1 \leq i \leq k-1$, and another $(x, y) \in [0, 1] \times \mathbb{R}$ with $x \neq u_i$ for any $1 \leq i \leq k$. If we let $d = \min_{1 \leq i \leq k} |x - u_i|$ and let $m = f'_S(x)$, which is the slope of the linear interpolation of S at x , then we have

$$H_{p, S \cup (x, y)} - H_{p, S} \geq (p-1)|m|d^p.$$

Proof First, if $x < u_1$, it suffices to prove that $H_{p, S \cup (x, y)} - H_{p, S} \geq 0$ as $m = 0$. Indeed, we have

$$H_{p, S \cup (x, y)} - H_{p, S} = |y - v_1|(1 - |x - u_1|^{p-1}) \geq 0.$$

The case where $x > u_k$ resolves similarly. Now, let $u_i < x < u_{i+1}$ for some $1 \leq i \leq k-1$. Without loss of generality assume that $v_i \leq v_{i+1}$. For convenience, define

$$\Delta_1 = \sum_{j=1}^{i-1} |v_{j+1} - v_j|(1 - |u_{j+1} - u_j|^{p-1}), \quad \Delta_2 = \sum_{j=i+1}^{k-1} |v_{j+1} - v_j|(1 - |u_{j+1} - u_j|^{p-1}).$$

First, we claim that we can reduce to only proving the case where $y \in [v_i, v_{i+1}]$. Indeed, note that if $y > v_{i+1}$, we can reduce it to the case where $y = v_{i+1}$, as

$$\begin{aligned} H_{p, S \cup (x, y)} &= \Delta_1 + \Delta_2 + \sum_{j=i}^{i+1} |y - v_j|(1 - |x - u_j|^{p-1}) \\ &\geq \Delta_1 + \Delta_2 + \sum_{j=i}^{i+1} |v_{i+1} - v_j|(1 - |x - u_j|^{p-1}) = H_{p, S \cup (x, v_{i+1})}, \end{aligned}$$

where the second step follows from the fact that $|y - v_i| \geq |v_{i+1} - v_i|$ and $|y - v_{i+1}| \geq 0 = |v_{i+1} - v_{i+1}|$. In the same way, the case where $y < v_i$ can be reduced to $y = v_i$. Assume from now on that $y \in [v_i, v_{i+1}]$. Also, assume that $d = |x - u_i| \leq |x - u_{i+1}|$. The case where x is nearer to u_{i+1} than u_i can be handled similarly. As $|y - v_i| + |y - v_{i+1}| = |v_{i+1} - v_i|$ and $1 - |x - u_{i+1}|^{p-1} \leq 1 - |x - u_i|^{p-1}$,

$$\begin{aligned} H_{p, S \cup (x, y)} &= \Delta_1 + \Delta_2 + \sum_{j=i}^{i+1} |y - v_j|(1 - |x - u_j|^{p-1}) \\ &\geq \Delta_1 + \Delta_2 + (|y - v_i| + |y - v_{i+1}|)(1 - |x - u_{i+1}|^{p-1}) \\ &= \Delta_1 + \Delta_2 + |v_{i+1} - v_i|(1 - |x - u_{i+1}|^{p-1}) \\ &= \Delta_1 + \Delta_2 + \sum_{j=i}^{i+1} |v_i - v_j|(1 - |x - u_j|^{p-1}) = H_{p, S \cup (x, v_i)}, \end{aligned}$$

so it suffices to check that $H_{p, S \cup (x, v_i)} - H_{p, S} \geq (p-1)|m|d^p$. Let $\frac{|u_i - u_{i+1}|}{|x - u_i|} = \frac{|u_i - u_{i+1}|}{d} = \lambda$. As such, we can substitute $|x - u_i|$, $|x - u_{i+1}|$, and $|u_i - u_{i+1}|$ with d , $\lambda d - d$, and λd , respectively. We have $\lambda \geq 2$, from x being closer to u_i than u_{i+1} . Note that as

$$\begin{aligned} H_{p, S \cup (x, v_i)} &= \Delta_1 + \Delta_2 + |v_{i+1} - v_i|(1 - |x - u_{i+1}|^{p-1}) \\ &= \Delta_1 + \Delta_2 + \left| \frac{v_{i+1} - v_i}{u_{i+1} - u_i} \right| (|u_{i+1} - u_i| - |x - u_{i+1}|^{p-1}|u_{i+1} - u_i|) \\ &= \Delta_1 + \Delta_2 + |m|(\lambda d - (\lambda d - d)^{p-1}(\lambda d)) \end{aligned}$$

and

$$\begin{aligned} H_{p,S} &= \Delta_1 + \Delta_2 + |v_{i+1} - v_i| (1 - |u_i - u_{i+1}|^{p-1}) \\ &= \Delta_1 + \Delta_2 + |m| (|u_i - u_{i+1}| - |u_i - u_{i+1}|^p) \\ &= \Delta_1 + \Delta_2 + |m| (\lambda d - (\lambda d)^p), \end{aligned}$$

we have that indeed

$$\begin{aligned} H_{p,S \cup (x, v_i)} - H_{p,S} &= |m|(\lambda d)^p - |m|(\lambda d - d)^{p-1}(\lambda d) \\ &= |m|d^p (\lambda^p - (\lambda - 1)^{p-1} \cdot \lambda) \geq |m|d^p(p - 1), \end{aligned}$$

where the last step follows from Lemma B.4.

Appendix C. Proofs for Section 3

In this section, we will formally prove the Lemmas used to establish our result on polynomials. These proofs are typically analysis heavy, hence why we omit the details in the main body.

Lemma C.1. *Given any $q \geq 1$, $\epsilon > 0$, and a set $S = \{(u_i, v_i) : 1 \leq i \leq m\}$ of points in $[0, 1] \times \mathbb{R}$ such that $u_1 < \dots < u_m$ and its corresponding linear interpolation function f_S , there exists a polynomial $P : [0, 1] \rightarrow \mathbb{R}$ such that $|P(u_i) - v_i| = |P(u_i) - f_S(u_i)| < \epsilon$ for all $1 \leq i \leq m$, and*

$$J_q[P] < J_q[f_S] + \epsilon.$$

Proof Consider the derivative of f_S , which when considered as a graph, consists of disjoint horizontal line segments, broken off at each u_i , where it is not defined. We modify f'_S to make it continuous, so Theorem 3.1 can be used. Define d_0, d_1, \dots, d_m to be the slopes of the segments of f_S in order:

$$d_i = \begin{cases} 0 & i = 0 \\ \frac{v_{i+1} - v_i}{u_{i+1} - u_i} & 1 \leq i \leq m - 1 \\ 0 & i = m. \end{cases}$$

Now, for a sufficiently small positive $\epsilon_2 < \min_{1 \leq i < j \leq m} |u_i - u_j|$, define the function $F'_S : [0, 1] \rightarrow \mathbb{R}$ such that

$$F'_S(x) = \begin{cases} d_0 & x \leq u_1 - \epsilon_2 \\ d_{i-1} + \frac{(x - (u_i - \epsilon_2))(d_i - d_{i-1})}{\epsilon_2} & x \in (u_i - \epsilon_2, u_i] \\ d_i & x \in (u_i, u_{i+1} - \epsilon_2] \\ d_m & x > u_m. \end{cases}$$

It is easy to see that F'_S is continuous for all $x \in (0, 1)$. Indeed, intuitively, the graph of F'_S is essentially the graph of f'_S but with the disjoint horizontal segments “connected”.

We claim that for sufficiently small ϵ_2 , $\int_0^x |F'_S(t)|^q dt$ gets arbitrarily close to $\int_0^x |f'_S(t)|^q dt$ for any $x \in [0, 1]$ and any $q \geq 1$. Indeed, $F'_S(t)$ only differs from $f'_S(t)$ when $t \in (u_i - \epsilon_2, u_i]$ for some $1 \leq i \leq m$. For a fixed i , when t is in this range, $f'_S(t)$ is precisely d_{i-1} , whereas $F'_S(t)$ is between

d_{i-1} and d_i . As such, $\left| |f'_S(t)|^q - |F'_S(t)|^q \right|$ is bounded above by $c = \max_{1 \leq i < j \leq m} (|d_i|^q - |d_j|^q)$, which is finite, for any q . As such, for any $x \in [0, 1]$,

$$\begin{aligned} \left| \int_0^x |F'_S(t)|^q dt - \int_0^x |f'_S(t)|^q dt \right| &\leq \int_0^x \left| |f'_S(t)|^q - |F'_S(t)|^q \right| dt \\ &\leq \sum_{i=1}^m \int_{u_i - \epsilon_2}^{u_i} \left| |f'_S(t)|^q - |F'_S(t)|^q \right| dt \\ &\leq \sum_{i=1}^m \int_{u_i - \epsilon_2}^{u_i} c dt = mc\epsilon_2, \end{aligned}$$

which gets arbitrarily close to 0 if we set ϵ_2 sufficiently small. We will use this inequality later on. Note that as a corollary, setting $x = 1$ above yields that $\int_0^1 |F'_S(t)|^q dt$ gets arbitrarily close to $J_q[f_S]$ when we set ϵ_2 sufficiently small—we will also use this fact later on.

Now, for sufficiently small $\epsilon_3 > 0$ (which we will specify later), consider a polynomial Q such that for all $x \in [0, 1]$, $|F'_S(x) - Q(x)| < \epsilon_3$, which is possible by Theorem 3.1 as F'_S is continuous for all $x \in [0, 1]$. Subsequently, let $P(x) \equiv \int_0^x Q(t) dt + v_1$, which is clearly a polynomial.

We claim that $J_q[P] < J_q[f_S] + \epsilon$ for any $\epsilon > 0$, given that our choices of ϵ_2 and ϵ_3 are sufficiently small. In particular, define $c_1 = \max_{x \in [0, 1]} |F'_S(x)| = \max_{0 \leq i \leq m} |d_i|$. Now, pick $\epsilon_2 < \frac{\epsilon}{2mc}$ and let ϵ_3 be sufficiently small such that $(c_1 + \epsilon_3)^q - c_1^q < \frac{\epsilon}{2}$. Note that this implies that $(|F'_S(x)| + \epsilon_3)^q - (|F'_S(x)|)^q < \frac{\epsilon}{2}$ for all $x \in [0, 1]$, as $0 \leq |F'_S(x)| \leq c_1$ and the function $(y + \epsilon_3)^q - y^q$ is increasing for all $y \geq 0$ and any positive ϵ .

As such, we have the upper bound

$$\begin{aligned} J_q[P] &= \int_0^1 |Q(x)|^q dx \leq \int_0^1 (|F'_S(x)| + \epsilon_3)^q dx \leq \int_0^1 \left(|F'_S(x)|^q + \frac{\epsilon}{2} \right) dx \\ &= \int_0^1 |F'_S(x)|^q dx + \frac{\epsilon}{2} \leq J_q[f_S] + \left(\int_0^1 |F'_S(x)|^q dx - J_q[f_S] \right) + \frac{\epsilon}{2} \\ &= J_q[f_S] + \left(\int_0^1 |F'_S(x)|^q dx - \int_0^1 |f'_S(x)|^q dx \right) + \frac{\epsilon}{2} \\ &\leq J_q[f_S] + mc\epsilon_2 + \frac{\epsilon}{2} < J_q[f_S] + \frac{\epsilon}{2} + \frac{\epsilon}{2} = J_q[f_S] + \epsilon. \end{aligned}$$

Furthermore, we also claim that $|P(u_i) - v_i| < \epsilon$ for any $\epsilon > 0$, given that our choices of ϵ_2 and ϵ_3 are sufficiently small. Indeed, as long as $\epsilon_3 < \frac{\epsilon}{2}$ and $\epsilon_2 < \frac{\epsilon}{2mc}$, we have for each $1 \leq i \leq m$ that

$$\begin{aligned} P(u_i) &= \int_0^{u_i} Q(t) dt + v_1 \leq \int_0^{u_i} (F'_S(t) + \epsilon_3) dt + v_1 \leq \int_0^{u_i} (F'_S(t)) dt + \epsilon_3 + v_1 \\ &\leq \int_0^{u_i} (f'_S(t)) dt + mc\epsilon_2 + \epsilon_3 + v_1 = (v_i - v_1) + v_1 + mc\epsilon_2 + \epsilon_3 \\ &= v_i + mc\epsilon_2 + \epsilon_3 < v_i + \frac{\epsilon}{2} + \frac{\epsilon}{2} = v_i + \epsilon \end{aligned}$$

and

$$\begin{aligned}
 P(u_i) &= \int_0^{u_i} Q(t)dt + v_1 \geq \int_0^{u_i} (F'_S(t) - \epsilon_3) dt + v_1 \geq \int_0^{u_i} (F'_S(t)) dt - \epsilon_3 + v_1 \\
 &\geq \int_0^{u_i} (f'_S(t)) dt - m\epsilon_2 - \epsilon_3 + v_1 = (v_i - v_1) + v_1 - m\epsilon_2 - \epsilon_3 \\
 &= v_i - m\epsilon_2 - \epsilon_3 < v_i - \frac{\epsilon}{2} - \frac{\epsilon}{2} = v_i - \epsilon.
 \end{aligned}$$

Therefore, for any $\epsilon > 0$, we can find a polynomial P such that $|P(u_i) - v_i| < \epsilon$ for all $1 \leq i \leq m$ and $J_q[P] < J_q[f_S] + \epsilon$.

Next, we establish Theorem 3.3, a result about constructing a function passing through all points in a set by taking a weighted average of a special set of functions that each approximate the set.

Lemma C.2. *Suppose there is a set $S = \{(u_i, v_i) : 1 \leq i \leq m\}$ of m points in $[0, 1] \times \mathbb{R}$ such that $u_1 < \dots < u_m$. If we have 2^m functions $f_X : [0, 1] \times \mathbb{R}$, corresponding to each subset $X \subseteq [m]$, such that for any $X \subseteq [m]$, it holds for any integer $1 \leq i \leq m$ with $i \in X$ that $f_X(u_i) > v_i$, and it holds for any integer $1 \leq i \leq m$ with $i \notin X$ that $f_X(u_i) < v_i$, then there exists a weighted average of the functions $f \equiv \sum_{X \subseteq [m]} w_X f_X$, with $0 \leq w_X \leq 1$ for all $X \subseteq [m]$ and $\sum_{X \subseteq [m]} w_X = 1$, such that $f(u_i) = v_i$ for all $1 \leq i \leq m$.*

Proof We proceed by induction. For the base case of $m = 1$, we are given two functions $f_\emptyset, f_{\{1\}} : [0, 1] \rightarrow \mathbb{R}$, and one point (u_1, v_1) that the desired weighted average should pass through. As $f_{\{1\}}(u_1) > v_1$ and $f_\emptyset(u_1) < v_1$, we can clearly assign two weights $0 \leq w_{\{1\}}, w_\emptyset \leq 1$ with sum 1 such that $w_{\{1\}}f_{\{1\}}(u_1) + w_\emptyset f_\emptyset = v_1$.

Proceeding with the inductive step, we assume that the result is true for $m - 1$, and prove that it is true for m as well. Let the set of points be $S = \{(u_i, v_i) : 1 \leq i \leq m\}$, and functions be $f_X : [0, 1] \times \mathbb{R}$ where $X \subseteq [m]$. Partition the 2^m functions into two groups, A and B , such that f_X goes into A if $m \in X$, and goes into B if $m \notin X$. Clearly, each of the two groups have 2^{m-1} functions. Note that for every subset $Y \subseteq [m-1]$, there exists a function in each of A and B such that the output at u_i is greater than v_i precisely when $i \in Y$, and the output at u_i is less than v_i precisely when $i \notin Y$ (specifically, the function $f_{Y \cup \{m\}}$ in A , and the function f_Y in B). By the inductive hypothesis, there exists a weighted average $F_1 \equiv \sum_{X \subseteq [m], m \in X} w_X f_X$ of the functions in A such that $F_1(u_i) = v_i$ for all $1 \leq i \leq m - 1$. Similarly, there exists a weighted average $F_2 \equiv \sum_{X \subseteq [m], m \notin X} w_X f_X$ of the functions in B such that $F_2(u_i) = v_i$ for all $1 \leq i \leq m - 1$.

Note that as all of the functions in group A satisfy $f(u_m) > v_m$, the weighted average of these functions, F_1 , must satisfy $F_1(u_m) > v_m$ as well. Similarly, as all of the functions in the functions in group B satisfy $f(u_m) < v_m$, the weighted average, F_2 , must satisfy $F_2(u_m) < v_m$. As such, there exists a weighted average of the functions F_1 and F_2 , $f \equiv W_1 F_1 + W_2 F_2$, with $0 \leq W_1, W_2 \leq 1$ and $W_1 + W_2 = 1$, such that $f(u_m) = v_m$. Furthermore, note that as $F_1(u_i) = F_2(u_i) = v_i$ for all $1 \leq i \leq m - 1$, $f(u_i) = v_i$ for all $1 \leq i \leq m - 1$ as well. Thus, $f(u_i) = v_i$ for all $1 \leq i \leq m$. Furthermore, as f is the weighted average of two weighted averages of functions f_X where $X \subseteq [m]$, f is a weighted average of functions f_X where $X \subseteq [m]$, and is hence our desired function.

Now, we establish 3.5.

Lemma C.3. *Given a set $S = \{(u_i, v_i) : 1 \leq i \leq m\}$ of m points in $[0, 1] \times \mathbb{R}$ with $u_1 < \dots < u_m$ and any $q \geq 1$, if there exists an absolutely continuous function $f : [0, 1] \rightarrow \mathbb{R}$ with $J_q[f] < 1$ such that $f(u_i) = v_i$ for all $1 \leq i \leq m$, then there exists a polynomial $P : [0, 1] \rightarrow \mathbb{R}$ with $J_q[P] < 1$ such that $P(u_i) = v_i$ for all $1 \leq i \leq m$.*

Proof Assume that $m \geq 2$, as the case where there is only one point in S is trivial. Note that by Lemma A.2, the condition that $J_q[f] < 1$ implies that $J_q[f_S] < 1$. Now, pick a small $\epsilon_1 > 0$ (which we will specify later), and define the set S_X , for each subset $X \subseteq [m]$, as follows:

$$S_X = \{(u_i, v_i + \epsilon_1) : 1 \leq i \leq m, i \in X\} \cup \{(u_i, v_i - \epsilon_1) : 1 \leq i \leq m, i \notin X\}.$$

As such, each S_X contains exactly m points, each of whose y -values are within ϵ_1 from the y -values of the corresponding points in S . For each X , we claim that setting $\epsilon_1 > 0$ sufficiently small, we can guarantee $J_q[f_{S_X}] < 1$. Indeed, if we let d_0, d_1, \dots, d_m denote the slopes of the segments of f_S in order, and define d'_0, d'_1, \dots, d'_m to be the slopes of the segments of f_{S_X} in order, notice that $d'_0 = d'_m = 0$, and for each $1 \leq i \leq m-1$, we have either $d'_i = d_i$, $d'_i = d_i + \frac{2\epsilon_1}{u_{i+1}-u_i}$, or $d'_i = d_i - \frac{2\epsilon_1}{u_{i+1}-u_i}$. As such, in any case, we have

$$|d'_i| \leq |d_i| + \frac{2\epsilon_1}{u_{i+1}-u_i} \leq |d_i| + \frac{2\epsilon_1}{\min_{1 \leq i \leq m-1} |u_{i+1}-u_i|} = |d_i| + C\epsilon_1,$$

if we substitute $C = \frac{2}{\min_{1 \leq i \leq m-1} |u_{i+1}-u_i|}$, which is fixed for a set S . Notice that for each i , we can set ϵ_1 sufficiently small such that the inequality $(|d_i| + C\epsilon_1)^q - |d_i|^q < \frac{1-J_q[f_S]}{m}$ holds, as $\frac{1-J_q[f_S]}{m}$ is positive. Hence, let ϵ_1 be sufficiently small such that the inequality $(|d_i| + C\epsilon_1)^q - |d_i|^q < \frac{1-J_q[f_S]}{m}$ holds for each $1 \leq i \leq m-1$.

Now, notice that

$$\begin{aligned} J_q[f_{S_X}] - J_q[f_S] &= \sum_{i=1}^{m-1} |u_{i+1}-u_i| |d'_i|^q - \sum_{i=1}^{m-1} |u_{i+1}-u_i| |d_i|^q = \sum_{i=1}^{m-1} |u_{i+1}-u_i| (|d'_i|^q - |d_i|^q) \\ &\leq \sum_{1 \leq i \leq m-1, |d'_i| > |d_i|} |u_{i+1}-u_i| (|d'_i|^q - |d_i|^q) \\ &\leq \sum_{1 \leq i \leq m-1, |d'_i| > |d_i|} (|d'_i|^q - |d_i|^q) \\ &\leq \sum_{1 \leq i \leq m-1, |d'_i| > |d_i|} ((|d_i| + C\epsilon_1)^q - |d_i|^q) < m \cdot \frac{1-J_q[f_S]}{m} = 1 - J_q[f_S], \end{aligned}$$

implying that $J_q[f_{S_X}] < 1$. Therefore, we can pick sufficiently small ϵ_1 such that for each $X \subseteq [m]$, we have $J_q[f_{S_X}] < 1$.

Now, for each $X \subseteq [m]$, pick a sufficiently small $\epsilon_2 > 0$ such that $\epsilon_2 < \epsilon_1$ and $\epsilon_2 < 1 - J_q[f_{S_X}]$, and let $P_X : [0, 1] \rightarrow \mathbb{R}$ be a polynomial such that

$$|P_X(u_i) - f_{S_X}(u_i)| < \epsilon_2 \text{ for all } 1 \leq i \leq m, \text{ and } J_q[P_X] < J_q[f_{S_X}] + \epsilon_2,$$

which must be possible by Lemma 3.2. By the first condition, we can see that whenever $1 \leq i \leq m$ satisfies $i \in X$, we have

$$P_X(u_i) > f_{S_X}(u_i) - \epsilon_2 > f_{S_X}(u_i) - \epsilon_1 = v_i,$$

and whenever $i \notin m$, we have

$$P_X(u_i) < f_{S_X}(u_i) + \epsilon_2 < f_{S_X}(u_i) + \epsilon_1 = v_i.$$

Furthermore, the second condition yields that $J_q[P_X] < 1$ as well.

As such, we have 2^m polynomials P_X , corresponding to each subset $X \subseteq [m]$, such that for each X , it holds for any integer $1 \leq i \leq m$ with $i \in X$ that $P_X(u_i) > v_i$, and it holds for any integer $1 \leq i \leq m$ with $i \notin X$ that $P_X(u_i) < v_i$. By Lemma 3.3, there exists a weighted average of the polynomials, another polynomial, $P \equiv \sum_{X \subseteq [m]} w_X P_X$, with weights $0 \leq w_X \leq 1$ for each X and $\sum_{X \subseteq [m]} w_X = 1$, such that $P(u_i) = v_i$ for all $1 \leq i \leq m$. Furthermore, by Lemma 3.4, as $J_q[P_X] < 1$ for all X , we know that $J_q[P] < 1$ as well.

Appendix D. Proofs for Section 4

In this section, we prove Lemma 4.1, 4.3, 4.4.

Theorem D.1. *For any integer $\eta \geq 1$, if incorrect feedback can be given up to η times, then at least $2\eta + 1$ initial rounds must be thrown out for the learner to guarantee finite error on its first prediction that counts toward the error evaluation.*

Proof We first show that after receiving 2η initial rounds of feedback, the learner cannot guarantee finite error on its next trial. Consider any sequence of inputs $\sigma = (x_0, x_1, \dots, x_{2\eta-1}) \in [0, 1]^{2\eta}$. Let the adversary reveal that $f(x_i) = 0$ for all $0 \leq i \leq \eta - 1$, and that $f(x_i) = C$ for all $\eta \leq i \leq 2\eta - 1$ for some arbitrarily large C . Now, suppose the learner is queried on the value of $f(x_{2\eta})$. If $\hat{y}_{x_{2\eta}} \geq \frac{C}{2}$, then let the function be $f(x) \equiv 0$, which is valid because exactly η of the feedback points did not align. Similarly, if $\hat{y}_{x_{2\eta}} \leq \frac{C}{2}$, let the function be $f(x) \equiv C$, which would also be valid. In any case, the raw error $|\hat{y}_{x_{2\eta}} - f(x_{2\eta})| \geq \frac{C}{2}$, which can be arbitrarily large.

Now, we establish that throwing out $2\eta + 1$ initial rounds suffices. We prove a stronger statement, that after $2\eta + 1$ initial rounds, the learner is able to bound the value of the function f at any input within a closed interval of length 2. Suppose that in the first $2\eta + 1$ rounds, the learner receives the set of points $S = \{(x_i, y_i) : 1 \leq i \leq 2\eta + 1\}$, where x_i and y_i denote the queried input and adversary feedback, respectively, in the i^{th} round.

Let $(v_1, v_2, \dots, v_{2\eta+1})$ be a permutation of $(y_1, y_2, \dots, y_{2\eta+1})$, such that $v_1 \leq \dots \leq v_{2\eta+1}$. Now, note that as the adversary may lie at most η times, at least $\eta + 1$ points among the $2\eta + 1$ points in S must align with the actual function. As such, at least $\eta + 1$ elements of $\{v_1, \dots, v_{2\eta+1}\}$ must be in the range of f . Note that any $\eta + 1$ -element subset of $\{v_1, \dots, v_{2\eta+1}\}$ must contain a member v_j with $j \geq \eta + 1$, so the range of f must contain some $v_j \geq v_{\eta+1}$. As the difference between the greatest and least values of f is at most 1, from the fact that $f \in \mathcal{F}_q$, we must have $f(x) \geq v_j - 1 \geq v_{\eta+1} - 1$ for all $x \in [0, 1]$. On the other hand, as any $\eta + 1$ -element subset of $\{v_1, \dots, v_{2\eta+1}\}$ must also contain a member v_k with $k \leq \eta + 1$, the range of f must also contain some $v_k \leq v_{\eta+1}$, from which it follows that $f(x) \leq v_{\eta+1} + 1$ for all $x \in [0, 1]$. As such, we have that $f(x)$ must be within $[v_{\eta+1} - 1, v_{\eta+1} + 1]$ for any $x \in [0, 1]$, as promised.

Lemma D.2. For any $\eta \geq 1$, $p, q \geq 2$ we have $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q) \leq 12\eta + 6$.

Proof Split the domain $[0, 1]$ into four intervals: $I_1 = [0, \frac{1}{4})$, $I_2 = [\frac{1}{4}, \frac{1}{2})$, $I_3 = [\frac{1}{2}, \frac{3}{4})$, $I_4 = [\frac{3}{4}, 1]$. We first claim that whenever the learner has adversary feedback on at least $2\eta + 1$ points in any interval I_j , it can determine a constant c such that for all $x \in I_j$, $f(x) \in [c - \frac{1}{2}, c + \frac{1}{2}]$.

To prove this, suppose that the learner receives the feedback set $S = \{(x_i, y_i) : 1 \leq i \leq 2\eta + 1\}$, where each $x_i \in I_j$ and $y_1 < \dots < y_{2\eta+1}$. Indeed, letting $c = y_{\eta+1}$, as there must be a true point in S with output at least c , and a true point in S with output at most c , there must exist a $X_1 \in I_j$ with $f(X_1) = c$ by the Intermediate Value Theorem.

Now, if there exists any point $X_2 \in I_j$ with $f(X_2) > c + \frac{1}{2}$, then

$$J_q[f] > J_q \left[f_{\{(X_1, c), (X_2, c + \frac{1}{2})\}} \right] = |X_2 - X_1| \left| \frac{\frac{1}{2}}{X_2 - X_1} \right|^q \geq \left(\frac{1}{4} \right)^{1-q} \left(\frac{1}{2} \right)^q \geq 1$$

as $|X_2 - X_1| \leq \frac{1}{4}$ and $q \geq 2$; contradiction. Similarly there cannot be a point $X_2 \in I_j$ with $f(X_2) < c - \frac{1}{2}$, so for all $X_2 \in I_j$, $f(X_2) \in [c - \frac{1}{2}, c + \frac{1}{2}]$, as claimed.

Now, let A be an optimal learning algorithm in the standard scenario for the same p, q , so that it always guarantees a total error value of at most $\text{opt}_p(\mathcal{F}_q) = 1$. Consider the learning algorithm \mathcal{A} trying to learn a function $f \in \mathcal{F}_q$ in the noisy scenario, following this strategy:

Upon the end of the initial $2\eta + 1$ rounds, first record the value v such that f can be bounded within $[v - 1, v + 1]$, previously shown to be possible. Now, let \mathcal{A} enter Stage 1. For any trial t requesting the value of $f(x_t)$ for $x_t \in I_j$ for some j , if the learner knows at least $2\eta + 1$ points already in I_j , predict exactly as algorithm A would; otherwise, predict v . Furthermore, call the trials of the former kind *mimicked trials*.

Now, let Stage 1 continue until a *mimicked trial*, with $x_t \in I_j$ for some j , where one of the following events happen: 1) The adversary reveals a point (x_t, y) such that $y \notin [c - \frac{1}{2}, c + \frac{1}{2}]$, where c is the constant such that f is bounded within $[c - \frac{1}{2}, c + \frac{1}{2}]$ over I_j ; 2) the algorithm A predicts \hat{y} with $\hat{y} \notin [c - \frac{1}{2}, c + \frac{1}{2}]$; or 3) the total error the learner perceives throughout all *mimicked trials* of the Stage exceeds $\text{opt}_p(\mathcal{F}_q) = 1$.

Once one of the three events happen, let \mathcal{A} exit Stage 1, forget all previous adversary feedback, predict v for the immediate next round, and enter Stage 2, repeating the same process until one of the three events happen again, whereupon it enters Stage 3, etc. Clearly, whenever one of the three events occur, a lie must have occurred since the end of the last stage. Thus, the learning process includes at most $\eta + 1$ stages.

Now, onto error bounding. As there are 4 intervals, each with at most $2\eta + 1$ *non-mimicked trials* which generate at most 1 error, the total error from *non-mimicked trials* is at most $4(2\eta + 1) = 8\eta + 4$.

Now, onto bounding *mimicked trials*. Note that every trial generates at most 1 error. For each stage besides stage $\eta + 1$, the total perceived error from the *mimicked trials*, not counting the stage's last trial, is at most $1 + \text{opt}_p(\mathcal{F}_q) = 2$. The last trial of every stage generates an actual error of at most 1.

Over all stages besides stage $\eta + 1$, there are at most η *mimicked trials* where the perceived error is not the actual error. The perceived error differs from the actual error by at most 1 in these trials, since the correct value is in $[c - \frac{1}{2}, c + \frac{1}{2}]$. Thus, the actual error from the first η stages is at most $(2 + 1)\eta + \eta = 4\eta$.

Finally, if the learning process reaches stage $\eta + 1$, the adversary cannot lie any more, so the error from that stage is at most $1 + \text{opt}_p(\mathcal{F}_q) = 2$. Therefore, in total, the sum of p^{th} powers of actual

errors for *mimicked trials* over all stages is at most $4\eta + 2$. Adding the error from *non-mimicked trials*, we get at most $(4\eta + 2) + (8\eta + 4) = 12\eta + 6$.

Lemma D.3. *For any $\eta \geq 1$, $p, q \geq 2$, we have $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q) \geq 2\eta + 1$.*

Proof Fix any algorithm A for learning \mathcal{F}_q . Consider the following adversary strategy. For the first $2\eta + 1$ rounds, repeatedly query the learner at 0, and reveal the output to be 0. Let $f(0) = 0$, so that no lies have been used up yet. For the next $2\eta + 1$ rounds, the adversary repeatedly queries the learner at 1. For the first η of these rounds, reveal the value of $f(1)$ to be -1 . Then, for the next η rounds, reveal the value of $f(1)$ to be 1. Then, on the next round, reveal the true value of $f(1)$ to be either 1 or -1 , whichever makes the total error function larger. Note that no matter what we choose, exactly η lies have been used up. We claim that we can always guarantee a total error of at least $2\eta + 1$. Indeed, note that

$$\sum_{i=1}^{2\eta+1} |\hat{y}_i + 1|^p + \sum_{i=1}^{2\eta+1} |\hat{y}_i - 1|^p = \sum_{i=1}^{2\eta+1} (|\hat{y}_i + 1|^p + |\hat{y}_i - 1|^p) \geq 2(2\eta + 1),$$

where the last step comes from the fact that for any $p \geq 2$ and any $\hat{y}_i \in \mathbb{R}$, $|\hat{y}_i + 1|^p + |\hat{y}_i - 1|^p \geq 2$. To see why this is true, notice that if $\hat{y}_i \notin [-1, 1]$ the result is obvious; otherwise, Jensen's inequality gives the result.

This means at least one of $\sum_{i=1}^{2\eta+1} |\hat{y}_i + 1|^p$ and $\sum_{i=1}^{2\eta+1} |\hat{y}_i - 1|^p$ is at least $2\eta + 1$. As the adversary can force an error of at least $2\eta + 1$ regardless of the learner's predictions, $\text{opt}_{p,\eta}^{\text{nf}}(\mathcal{F}_q) \geq 2\eta + 1$.