
Parameter-free Statistically Consistent Interpolation: Dimension-independent Convergence Rates for Hilbert kernel regression

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Previously, statistical textbook wisdom has held that interpolation of noisy training
2 data will lead to poor generalization. However, recent work has shown that this
3 is not true and that good generalization can be obtained with function fits that
4 interpolate training data. This could explain why overparameterized deep nets with
5 zero or small training error do not necessarily overfit and could generalize well.
6 Data interpolation schemes have been exhibited that are provably Bayes optimal in
7 the large sample limit and achieve the theoretical lower bounds for excess risk (Sta-
8 tistically Consistent Interpolation) in any dimension. These interpolation schemes
9 are non-parametric Nadaraya-Watson style estimators with singular kernels, which
10 exhibit statistical consistency in any data dimension for large sample sizes. The
11 recently proposed weighted interpolating nearest neighbors scheme (wiNN) is in
12 this class, as is the previously studied Hilbert kernel interpolation scheme. In
13 the Hilbert scheme, the regression function estimator for a set of labelled data
14 pairs, $(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$, $i = 0, \dots, n$, has the form $\hat{f}(x) = \sum_i y_i w_i(x)$, where
15 $w_i(x) = \|x - x_i\|^{-d} / \sum_j \|x - x_j\|^{-d}$. This interpolating function estimator is
16 unique in being entirely free of parameters and does not require bandwidth se-
17 lection. While statistical consistency was previously proven for this scheme, the
18 precise convergence rates for the finite sample risk were not established. Here, we
19 carry out a comprehensive study of the asymptotic finite sample behavior of the
20 Hilbert kernel regression scheme and prove a number of relevant theorems. We
21 prove under broad conditions that the excess risk of the Hilbert regression estimator
22 is asymptotically equivalent pointwise to $\sigma^2(x) / \ln(n)$ where $\sigma^2(x)$ is the noise
23 variance. We also show that the excess risk of the plugin classifier is upper bounded
24 by $2|f(x) - 1/2|^{1-\alpha} (1 + \varepsilon)^\alpha \sigma^\alpha(x) (\ln(n))^{-\frac{\alpha}{2}}$, for any $0 < \alpha < 1$, where f is
25 the regression function $x \mapsto \mathbb{E}[y|x]$. Our proofs proceed by deriving asymptotic
26 equivalents of the moments of the weight functions $w_i(x)$ for large n , for instance
27 for $\beta > 1$, $\mathbb{E}[w_i^\beta(x)] \sim_{n \rightarrow \infty} ((\beta - 1)n \ln(n))^{-1}$. We further derive an asymptotic
28 equivalent for the Lagrange function and explicitly exhibit the nontrivial extrapola-
29 tion properties of this estimator. Notably, the convergence rates are independent of
30 data dimension and the excess risk is dominated by the noise variance. The bias
31 term, for which we also give precise asymptotic estimates, is always subleading
32 when the density of data at the considered point is strictly positive. If this local
33 density is zero, we show that the bias term does not vanish in the limit of a large
34 data set and we compute its limit explicitly. Finally, we present heuristic arguments
35 for a universal w^{-2} power-law behavior of the probability density of the weights
36 in the large n limit.

1 Introduction

Data interpolation and statistical regression of noisy data are both classical subjects but their domain of application have been disjoint until recently. Scattered data interpolation techniques [1] are generally used for clean data. On the other hand, when supervised learning or statistical regression techniques are applied to noisy data, in general smoothing or regularization methods are applied to prevent training data interpolation, as the latter is believed to lead to poor generalization [2]. However, accumulating empirical evidence from overparameterized deep networks has shown that data interpolation (equivalently, zero error on the training set) does not automatically imply poor generalization [3, 4]. This has in turn given rise to a rapidly growing body of theoretical work to understand how and why noisy data interpolation can still lead to good generalization [5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15].

A key observations in this regard is the phenomenon of Statistically Consistent Interpolation [16], i.e., regression function estimation that interpolates training data but also generalizes as well as possible by achieving the Bayes limit for expected generalization error (risk) when the sample size becomes large. This hints at a rich set of theoretical questions at the interface between the disciplines of scattered data interpolation and supervised learning, that have only begun to be addressed. In particular, there has been comparatively little study of the generalization error or risk of interpolating learners. Computation of generalization error bounds in machine learning often relies on the capacity of the class of fitting functions [17], however such model complexity based bounds are not tight enough to be useful for interpolating learners [4]. For nonparametric interpolation approaches such as that considered here, it is also not clear what model complexity means. Thus, there is a need for other approaches to understanding the generalization behavior of nonparametric interpolating learners, including more direct treatments of the generalization error for specific interpolation schemes so as to gain better theoretical understanding. The current paper addresses this need.

We present a detailed analysis of the finite-sample risk of an interpolating learner with intriguing theoretical properties, the Hilbert kernel estimator (Devroye *et al.* [18]). A unique property of this Nadaraya-Watson (NW) style estimator [19, 20] is that it is fully parameter-free and does not have any bandwidth or scale parameter. It is global and uses all data points for each estimate: the associated kernel is a power law and thus scale-free. Although statistical consistency of this estimator was proven [18] when it was proposed, there has been no systematic analysis of the associated convergence rates and asymptotic finite sample behavior. We provide this analysis in the present study.

Related work The only other interpolation scheme we are aware of, that is proven to be statistically consistent in arbitrary dimensions under general conditions, is the recently proposed weighted interpolating nearest neighbors method (wiNN) [7], which is also a NW estimator utilizing a singular power law kernel of a very similar form but with two important differences: a finite number of neighbors k is utilized (rather than all data points), and the power law exponent δ of the NW kernel satisfies $0 < \delta < d/2$ rather than $\delta = d$. To achieve consistency k has to scale appropriately with sample size. Despite the superficial resemblance, the wiNN and Hilbert Kernel estimators have quite different convergence rates, as we will see from the results of this paper. Also worth mentioning is the Shepard interpolation scheme [21] originally proposed for interpolation of 2D geospatial data sets, also a NW style interpolating estimator, though used in the context of scattered data interpolation. In scattered data interpolation [1], the focus is generally on the approximation error (corresponding to the “bias” term in our analysis below). The approximation error of the Shepard scheme has been analyzed [22] but as we will see below the risk for Hilbert kernel interpolation is dominated by the noise or “variance” term. In contrast with wiNN or Hilbert kernel interpolation, other interpolating learning methods such as simplex interpolation [7] or ridgeless kernel regression [11] are generally not statistically consistent in fixed finite dimension [8].

Summary of results of this paper Notation and assumptions pertaining to this summary are defined in the problem setup section below. We prove under broad conditions that the excess risk of the Hilbert regression estimator is asymptotically equivalent pointwise to $\sigma^2(x)/\ln(n)$ where $\sigma^2(x)$ is the noise variance. We also show that the excess risk of the plugin classifier is upper bounded by $2|f(x) - 1/2|^{1-\alpha} (1 + \varepsilon)^\alpha \sigma^\alpha(x) (\ln(n))^{-\frac{\alpha}{2}}$, for any $0 < \alpha < 1$, where f is the regression function $x \mapsto \mathbb{E}[y|x]$. Our proofs proceed by deriving asymptotic equivalents of the moments of the weight functions $w_i(x)$ for large n , for instance for $\beta > 1$, $\mathbb{E}[w_i^\beta(x)] \sim_{n \rightarrow \infty} ((\beta - 1)n \ln(n))^{-1}$. We further derive an asymptotic equivalent for the Lagrange function and explicitly exhibit the nontrivial

93 extrapolation properties of this estimator. Notably, the convergence rates are independent of data
 94 dimension and the excess risk is dominated by the noise variance. The bias term, for which we also
 95 give precise asymptotic estimates, is always subleading when the density of data at the considered
 96 point is strictly positive. If this local density is zero, we show that the bias term does not vanish in the
 97 limit of a large data set and we compute its limit explicitly. Finally, we present heuristic arguments
 98 for a universal w^{-2} power-law behavior of the probability density of the weights in the large n limit.

99 2 Problem setup

100 **Notation, Definitions, Statistical Model** We model the labelled training data set
 101 $(x_0, y_0), \dots, (x_n, y_n)$ as $n + 1$ *i.i.d.* observations of a random vector (X, Y) with values in
 102 $\mathbb{R}^d \times \mathbb{R}$ for regression, and with values in $\mathbb{R}^d \times \{0, 1\}$ for binary classification. Due to the indepen-
 103 dence property, the collection X_0, \dots, X_n has the product density $\prod_{i=0}^n \rho(x_i)$. We will denote by \mathbb{E}
 104 an expectation over the collection of $n + 1$ random vectors and by \mathbb{E}_X the expectation over the col-
 105 lection X_0, \dots, X_n . An expectation over the same collection while holding $X_i = x_i$ will be denoted
 106 $\mathbb{E}_{X|x_i}$. The regression function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is defined as the conditional mean of Y given $X = x$,
 107 $f(x) := \mathbb{E}[Y | X = x]$ and the conditional variance function is $\sigma^2(x) := \mathbb{E}[|Y - f(X)|^2 | X = x]$.
 108 f minimizes the expected value of the mean squared prediction error (risk under squared loss),
 109 $f = \arg \min \mathcal{R}_{\text{sq}}(h)$ where $\mathcal{R}_{\text{sq}}(h) := \mathbb{E}[(h(X) - Y)^2]$. Given any regression estimator $\hat{f}(x)$ the
 110 corresponding risk can be decomposed as $\mathbb{E}[\mathcal{R}_{\text{sq}}(\hat{f}(X))] = \mathcal{R}_{\text{sq}}(f) + \mathbb{E}[(\hat{f}(X) - f(X))^2]$. The
 111 excess risk is given by $\mathcal{R}_{\text{sq}}(\hat{f}) - \mathcal{R}_{\text{sq}}(f) = \mathbb{E}[(\hat{f}(X) - f(X))^2]$. For a consistent estimator this
 112 excess risk goes to zero as $n \rightarrow \infty$ and we are interested in characterizing the *rate* at which it goes to
 113 zero with increasing n (note our sample size is $n + 1$ for notational simplicity but for large n this
 114 does not change the rate).

115 In the case of binary classification, $Y \in \{0, 1\}$ and $f(x) = \mathbb{P}[Y = 1 | X = x]$. Let $F: \mathbb{R}^d \rightarrow \{0, 1\}$
 116 denote the Bayes optimal classifier, defined by $F(x) := \theta(f(x) - 1/2)$ where $\theta(\cdot)$ is the Heaviside
 117 theta function. This classifier minimizes the risk $\mathcal{R}_{0/1}(h) := \mathbb{E}[\mathbb{1}_{\{h(X) \neq Y\}}] = \mathbb{P}(h(X) \neq Y)$ under
 118 zero-one loss. Given the regression estimator \hat{f} , we consider the plugin classifier $\hat{F}(x) = \theta(\hat{f}(x) - \frac{1}{2})$.
 119 The classification risk for the plugin classifier \hat{F} is bounded as $\mathbb{E}[\mathcal{R}_{0/1}(\hat{F}(x))] - \mathcal{R}_{0/1}(F(x)) \leq$
 120 $2\mathbb{E}[|\hat{f}(x) - f(x)|] \leq 2\sqrt{\mathbb{E}[(\hat{f}(x) - f(x))^2]}$.

121 Finally, we define two sequences $a_n, b_n > 0$, $n \in \mathbb{N}$, to be asymptotically equivalent for $n \rightarrow +\infty$,
 122 denoted $a_n \sim_{n \rightarrow +\infty} b_n$, if the limit of their ratio exists and $\lim_{n \rightarrow \infty} a_n/b_n = 1$.

123 In summary, our work will focus on the estimation of asymptotic equivalents for $\mathbb{E}[(\hat{f}(x) - f(x))^2]$
 124 and other relevant quantities as this determines the rate at which the excess risk goes to zero for
 125 regression, and bounds the rate at which the excess risk goes to zero for classification.

126 **Assumptions.** We define the support Ω of the density ρ as $\Omega = \{x \in \mathbb{R}^d / \rho(x) > 0\}$, the closed
 127 support $\bar{\Omega}$ as the closure of Ω , and Ω° as the interior of Ω . Our results will not assume any compactness
 128 condition on Ω or $\bar{\Omega}$. The boundary of Ω is then defined as $\partial\Omega = \bar{\Omega} \setminus \Omega^\circ$. We assume that ρ has a
 129 finite variance σ_ρ^2 . In addition, we will most of the time assume that the density ρ is continuous at the
 130 considered point $x \in \Omega^\circ$, and in some cases, $x \in \partial\Omega \cap \Omega$.

131 For the regression function f , we will obtain results assuming either of the following conditions

- 132 • C_{Cont}^f : f is continuous at the considered x ,
- 133 • C_{Holder}^f : for all $x \in \Omega^\circ$, there exist $\alpha_x > 0$, $K_x > 0$, and $\delta_x > 0$, such that
 134 $x' \in \Omega$ and $\|x - x'\| \leq \delta_x \implies |f(x) - f(x')| \leq K_x \|x - x'\|^{\alpha_x}$
 135 (local Hölder smoothness condition),

136 where condition C_{Holder}^f is obviously stronger than C_{Cont}^f . In addition, we will always assume a
 137 growth condition for the regression function f :

- 138 • C_{Growth}^f : $\int \rho(y) \frac{f^2(y)}{1+\|y\|^{2\alpha}} d^d y < \infty$.

139 As for the variance function σ , we will obtain results assuming either that σ is bounded or satisfies a
 140 growth condition similar to the one above

- 141 • C_{Bound}^σ : there exists $\sigma_0^2 \geq 0$, such that, for all $x \in \Omega$, we have $\sigma^2(x) \leq \sigma_0^2$,
- 142 • C_{Growth}^σ : $\int \rho(y) \frac{\sigma^2(y)}{1+\|y\|^{2\alpha}} d^d y < \infty$.

143 When we will assume condition C_{Growth}^σ (obviously satisfied when σ^2 is bounded), we will also
 144 assume a continuity condition C_{Cont}^σ for σ at the considered x .

145 Note that all our results can be readily extended in the case where $x \in \partial\Omega = \bar{\Omega} \setminus \Omega^\circ$ but keeping
 146 the condition $\rho(x) > 0$ (i.e., $x \in \partial\Omega \cap \Omega$), and assuming the continuity at x of ρ as seen as a
 147 function restricted to Ω , i.e., $\lim_{y \in \Omega \rightarrow x} \rho(y) = \rho(x)$. Useful examples are when the support Ω of
 148 ρ is a d -dimensional sphere or hypercube and x is on the surface of Ω (but still with $\rho(x) > 0$). To
 149 guarantee these results for $x \in \partial\Omega \cap \Omega$, we need also to assume the continuity at x of f , and assume
 150 that Ω is smooth enough near x , so that there exists a strictly positive local solid angle ω_x defined by

$$\omega_x = \lim_{r \rightarrow 0} \frac{1}{V_d \rho(x) r^d} \int_{\|x-y\| \leq r} \rho(y) d^d y = \lim_{r \rightarrow 0} \frac{1}{V_d r^d} \int_{y \in \Omega / \|x-y\| \leq r} d^d y, \quad (1)$$

151 where $V_d = S_d/d = \pi^{d/2}/\Gamma(d/2 + 1)$ is the volume of the unit ball in d dimensions, and the second
 152 inequality results from the continuity of ρ at x . If $x \in \Omega^\circ$, we have $\omega_x = 1$, while for $x \in \partial\Omega$, we
 153 have $0 \leq \omega_x \leq 1$. For instance, if x is on the surface of a sphere or on the interior of a face of a
 154 hypercube (and in general, when the boundary near x is locally an hyperplane), we have $\omega_x = \frac{1}{2}$. If x
 155 is a corner of a hypercube, we have $\omega_x = \frac{1}{2^d}$. From our methods of proof presented in the appendix,
 156 it should be clear that all our results for $x \in \Omega^\circ$ perfectly generalize to any $x \in \partial\Omega \cap \Omega$ for which
 157 $\omega_x > 0$, by simply replacing V_d whenever it appears in our different results by $\omega_x V_d$.

158 **Hilbert kernel interpolating estimator and Bias-Variance decomposition.** The Hilbert kernel
 159 regression estimator $\hat{f}(x)$ is a Nadaraya-Watson style estimator employing a singular kernel:

$$w_i(x) = \frac{\|x - x_i\|^{-d}}{\sum_{j=0}^n \|x - x_j\|^{-d}}, \quad (2)$$

$$\hat{f}(x) = \sum_{i=0}^n w_i(x) y_i. \quad (3)$$

160 The weights $w_i(x)$ are also called Lagrange functions in the interpolation literature and satisfy the
 161 interpolation property $w_i(x_j) = \delta_{ij}$, where $\delta_{ij} = 1$, if $i = j$, and 0 otherwise. At any given point
 162 x , they provide a partition of unity so that $\sum_{i=0}^n w_i(x) = 1$. The mean squared error between the
 163 Hilbert estimator and the true regression function has a bias-variance decomposition (using the *i.i.d*
 164 condition and the earlier definitions)

$$\hat{f}(x) - f(x) = \sum_{i=0}^n w_i(x) [f(x_i) - f(x)] + \sum_{i=0}^n w_i(x) [y_i - f(x_i)], \quad (4)$$

$$\mathbb{E}[(\hat{f}(x) - f(x))^2] = \mathcal{B}(x) + \mathcal{V}(x), \quad (5)$$

$$(\text{Bias}) \mathcal{B}(x) = \mathbb{E}_X \left[\left(\sum_{i=0}^n w_i(x) [f(x_i) - f(x)] \right)^2 \right], \quad (6)$$

$$(\text{Variance}) \mathcal{V}(x) = \mathbb{E} \left[\sum_{i=0}^n w_i^2(x) [y_i - f(x_i)]^2 \right] = \mathbb{E}_X \left[\sum_{i=0}^n w_i^2(x) \sigma^2(x_i) \right]. \quad (7)$$

165 The present work derives asymptotic behaviors and bounds for the regression and classification risk
 166 of the Hilbert estimator for large sample size n . These results are derived by analyzing the large n
 167 behaviors of the bias and variance terms, which in turn depend on the behavior of the moments of the
 168 weights or the Lagrange functions $w_i(x)$. For all these quantities, asymptotically equivalent forms
 169 are derived. The proofs exploit a simple integral form of the weight function and details are provided
 170 in the appendix, while the body of the paper provides the results and associated discussions.

171 **3 Results**

172 **3.1 The weights, variance and bias terms**

173 **3.1.1 Moments of the weights: large n behavior**

174 In this section, we consider the moments and the distribution of the weights $w_i(x)$ at a given point x .
 175 The first moment is simple to compute. Since the weights sum to 1 and X_i are *i.i.d.*, it follows that
 176 $\mathbb{E}_{X_i|x_i}[w_i(x)]$ are all equal and thus $\mathbb{E}_{X_i|x_i}[w_i(x)] = (n+1)^{-1}$. The other moments are much less
 177 trivial to compute and we prove the following theorem in the appendix A.2:

178 **Theorem 3.1.** *For $x \in \Omega^\circ$ (so that $\rho(x) > 0$), we assume ρ continuous at x . Then, the moments of*
 179 *the weight $w_0(x)$ satisfy the following properties:*

- 180 • For $\beta > 1$:

$$\mathbb{E} \left[w_0^\beta(x) \right] \underset{n \rightarrow +\infty}{\sim} \frac{1}{(\beta - 1)n \ln(n)}. \quad (8)$$

- 181 • For $0 < \beta < 1$: defining $\kappa_\beta(x) := \int \frac{\rho(x+y)}{\|y\|^{|\beta|d}} d^d y < \infty$, we have

$$\mathbb{E} \left[w_0^\beta(x) \right] \underset{n \rightarrow +\infty}{\sim} \frac{\kappa_\beta(x)}{(V_d \rho(x) n \ln(n))^\beta}. \quad (9)$$

- 182 • For $\beta < 0$: all moments for $\beta \leq -1$ are infinite, and the moments of order $-1 < \beta < 0$
 183 satisfy

$$\mathbb{E} \left[w_0^\beta(x) \right] \leq 1 + n \kappa_{|\beta|}(x) \kappa_\beta(x), \quad (10)$$

184 so that a sufficient condition for its existence is $\kappa_\beta(x) = \int \rho(x+y) \|y\|^{|\beta|d} d^d y < \infty$.

185 Heuristically, the behavior of these moments are consistent with the random variable $W = w_0(x)$
 186 having a probability distribution satisfying a scaling relation $P(W) = \frac{1}{W_n} p\left(\frac{W}{W_n}\right)$, with the scaling
 187 function p having the universal tail (i.e., independent of x and ρ), $p(w) \underset{w \rightarrow +\infty}{\sim} w^{-2}$, and a scale W_n
 188 expected to vanish with n , when $n \rightarrow +\infty$. With this assumption, we can determine the scale W_n by
 189 imposing the exact condition $\mathbb{E}[W] = 1/(n+1) \sim 1/n$:

$$\mathbb{E}[W] = \frac{1}{W_n} \int_0^1 p\left(\frac{W}{W_n}\right) W dW = W_n \int_0^{\frac{1}{W_n}} p(w) w dw \quad (11)$$

$$\sim W_n \int_1^{\frac{1}{W_n}} \frac{dw}{w} \sim -W_n \ln(W_n) \sim \frac{1}{n}, \quad (12)$$

190 leading to $W_n \sim \frac{1}{n \ln(n)}$. Then, the moment of order $\beta > 1$ is given by

$$\mathbb{E}[W^\beta] = \frac{1}{W_n} \int_0^1 p\left(\frac{W}{W_n}\right) W^\beta dW \sim W_n \int_0^1 W^{\beta-2} dW \underset{n \rightarrow +\infty}{\sim} \frac{1}{(\beta - 1)n \ln(n)}, \quad (13)$$

191 which indeed coincides with the first result of Theorem 3.1. Our heuristic argument also suggests
 192 that in the case $0 < \beta < 1$, we have

$$\mathbb{E}[W] = \frac{1}{W_n} \int_0^1 p\left(\frac{W}{W_n}\right) W^\beta dW \underset{n \rightarrow +\infty}{\sim} \frac{\int_0^{+\infty} p(w) w^\beta dw}{(n \ln(n))^\beta}, \quad (14)$$

193 where the last integral converges since $p(w) \underset{w \rightarrow +\infty}{\sim} w^{-2}$ and $\beta < 1$. This result is perfectly consistent
 194 with Eq. (9) in Theorem 3.1, and suggests that $\int_0^{+\infty} p(w) w^\beta dw = \frac{\kappa_\beta(x)}{(V_d \rho(x))^\beta}$. Interestingly, for
 195 $0 < \beta < 1$, and contrary to the case $\beta > 1$, we find that the large n equivalent of the moment is not
 196 universal and depends explicitly on x and the density ρ . As for moments of order $-1 < \beta < 0$, we
 197 conjecture that they are still given by Eq. (9) (and equivalently, by Eq. (14)) provided they exist, and
 198 that the sufficient condition for their existence $\kappa_\beta(x) < \infty$ is hence also necessary, since $\kappa_\beta(x)$ also
 199 appears in Eq. (9). The fact that moments for $\beta \leq -1$ do not exist strongly suggests that $p(0) > 0$.

200 In fact, Eq. (14)) also suggests that all moments for $-1 < \beta < 0$ exist if and only if $0 < p(0) < \infty$.
 201 In the Fig. 2 of the appendix, we present numerical simulations confirming our scaling ansatz, the
 202 fact that $p(w) \underset{w \rightarrow +\infty}{\sim} w^{-2}$, and the quantitative prediction for W_n .

203 It is shown in Devroye *et al.* [18] that the Hilbert kernel regression estimate does not converge
 204 almost surely (*a.s.*) by giving a specific example. Insight can be gained into this lack of almost sure
 205 convergence by considering the weight function $w_0(x)$, for a sequence of independent training sample
 206 sets of increasing size $n + 1$. Let the corresponding sequence of weights be denoted as $\omega_n \in [0, 1]$.
 207 From Theorem 3.1, it is clear that ω_n converges to zero in probability, since the following Chebyshev
 208 bound holds (analogous to the bound on the regression risk):

$$\mathbb{P}(\omega_n > \varepsilon) \leq \frac{1 + \delta}{\varepsilon^2 n \ln(n)}, \quad (15)$$

209 for arbitrary $\varepsilon > 0$ and $\delta > 0$, and for n larger than some constant $N_{x,\delta}$. Alternatively, one can
 210 exploit the fact that $\mathbb{E}[\omega_n] = \frac{1}{n+1}$, leading to $\mathbb{P}(\omega_n > \varepsilon) \leq \frac{1}{\varepsilon n}$, which is less stringent than Eq. (15)
 211 as far as the n -dependence is concerned, but is more stringent for the ε -dependence of the bounds.

212 Let us show heuristically that ω_n does not converge *a.s.* to zero. Consider the infinite sequence of
 213 events $\mathcal{E}_n \equiv \{\omega_n > \varepsilon\}$, $n \in \mathbb{N}$, and the corresponding infinite sum $\sum_n \mathbb{P}(\mathcal{E}_n) = \sum_n \mathbb{P}(\omega_n > \varepsilon)$.
 214 Exploiting our previous heuristic argument for the scaling form of the distribution of weights, we
 215 obtain

$$\mathbb{P}(\omega_n > \varepsilon) = \int_{\varepsilon}^1 \frac{1}{W_n} p\left(\frac{W}{W_n}\right) dW \sim \int_{\varepsilon n \ln n}^{n \ln n} \frac{dw}{w^2} \sim \frac{1 - \varepsilon}{\varepsilon n \ln(n)}. \quad (16)$$

216 Since $\sum_{n=2}^N \frac{1}{n \ln(n)} \sim \ln(\ln(N))$ is a divergent series, a Borel-Cantelli argument suggests that an
 217 infinite number of the events \mathcal{E}_n (i.e., $\omega_n > \varepsilon$) must occur, which implies that ω_n does not converge
 218 *a.s.* to 0. Note that the weights are equal to 1 at the data points due to the interpolation condition, so
 219 that large weights occasionally occur, causing the lack of *a.s.* convergence.

220 3.1.2 Lagrange function: scaling limit

221 The expected value of the Lagrange functions $w_i(x)$ have a simple form in the large n limit. Due
 222 to the *i.i.d.* condition the indices i are exchangeable and we set $i = 0$ for the computation of the
 223 expected Lagrange function $L_0(x) = \mathbb{E}_{X|x_0}[w_0(x)]$. Thus, one of the sample points (denoted x_0) is
 224 held fixed and the other ones are averaged over in computing the expected Lagrange function. For
 225 $x_0 \neq x$ kept fixed, we have $\lim_{n \rightarrow \infty} L_0(x) = 0$. However, we show in the appendix A.3 that $L_0(x)$
 226 takes a very simple form when taking a specific scaling limit:

227 **Theorem 3.2.** *For $x \in \Omega^\circ$, we assume ρ continuous at x . Then, in the limit (denoted by \lim_Z), $n \rightarrow$
 228 $+\infty$, $\|x - x_0\|^{-d} \rightarrow +\infty$ (i.e., $x_0 \rightarrow x$), and such that $z_x(n, x_0) = V_d \rho(x) \|x - x_0\|^d n \log(n) \rightarrow Z$,
 229 the Lagrange function $L_0(x) = \mathbb{E}_{X|x_0}[w_0(x)]$ converges to a proper limit,*

$$\lim_Z L_0(x) = \frac{1}{1 + Z}. \quad (17)$$

230 The proof of this theorem shows that the relative error between $L_0(x)$ and $\frac{1}{1+Z}$ for finite but large n
 231 and large $\|x - x_0\|^{-d}$, such that $z_x(n, x_0)$ remains close to Z , is $O(1/\ln(n))$.

232 Exploiting Theorem 3.2, we can use a simple heuristic argument to estimate the tail of the distribution
 233 of the random variable $W = w_0(x)$. Indeed, approximating $L_0(x)$ for finite but large n by its
 234 asymptotic form $\frac{1}{1+z_x(n, x_0)}$, with $z_x(n, x_0) = V_d \rho(x) n \log(n) \|x - x_0\|^d$, we obtain

$$\int_W^1 P(W') dW' \sim \int \rho(x_0) \theta\left(\frac{1}{1 + V_d \rho(x) n \log(n) \|x - x_0\|^d} - W\right) d^d x_0, \quad (18)$$

$$\sim V_d \rho(x) \int_0^{+\infty} \theta\left(\frac{1}{1 + V_d \rho(x) n \log(n) u} - W\right) du, \quad (19)$$

$$\sim \frac{1}{n \ln(n) W} \implies P(W) \sim \frac{1}{n \ln(n) W^2}, \quad (20)$$

235 where $\theta(\cdot)$ is the Heaviside function. This heuristic result is again perfectly consistent with our guess
 236 of the previous section that $P(W) = \frac{1}{W_n} p\left(\frac{W}{W_n}\right)$, with the scaling function p having the universal

237 tail, $p(w) \underset{w \rightarrow +\infty}{\sim} w^{-2}$, and a scale $W_n \sim \frac{1}{n \ln(n)}$. Indeed, in this case and in the limit $n \rightarrow +\infty$, we
 238 obtain that $P(W) \sim \frac{1}{W_n} \left(\frac{W_n}{W} \right)^2 \sim \frac{W_n}{W^2} \sim \frac{1}{n \ln(n) W^2}$, which is identical to the result of Eq. (20).

239 3.1.3 The variance term

240 A simple application of the result of Theorem 3.1 for $\beta = 2$ (see appendix A.4) allows us to bound
 241 the variance term $\mathcal{V}(x) = \mathbb{E} \left[\sum_{i=0}^n w_i^2(x) [y_i - f(x_i)]^2 \right]$ for a bounded variance function σ^2 :

242 **Theorem 3.3.** *For $x \in \Omega^\circ$, ρ continuous at x , $\sigma^2 \leq \sigma_0^2$, and for any $\varepsilon > 0$, there exists a constant*
 243 *$N_{x,\varepsilon}$ such that for $n \geq N_{x,\varepsilon}$, we have*

$$\mathcal{V}(x) \leq (1 + \varepsilon) \frac{\sigma_0^2}{\ln(n)}. \quad (21)$$

244 Relaxing the boundedness condition for σ , but assuming the continuity of σ^2 at x along with a growth
 245 condition, allows us to obtain a precise asymptotic equivalent of $\mathcal{V}(x)$, when $n \rightarrow +\infty$:

246 **Theorem 3.4.** *For $x \in \Omega^\circ$, $\sigma(x) > 0$, $\rho \sigma^2$ continuous at x , and assuming the condition C_{Growth}^σ ,*
 247 *i.e., $\int \rho(y) \frac{\sigma^2(y)}{1 + \|y\|^{2d}} d^d y < \infty$, we have*

$$\mathcal{V}(x) \underset{n \rightarrow +\infty}{\sim} \frac{\sigma^2(x)}{\ln(n)}. \quad (22)$$

248 Note that if the mean variance $\int \rho(y) \sigma^2(y) d^d y < \infty$, which is in particular the case when σ^2 is
 249 bounded over Ω , then the condition C_{Growth}^σ is in fact automatically satisfied.

250 3.1.4 The bias term

251 In appendix A.5, we prove the following three theorems for the bias term.

252 **Theorem 3.5.** *For $x \in \Omega^\circ$ (so that $\rho(x) > 0$), we assume that ρ is continuous at x , and the conditions*

- 253 • C_{Growth}^f : $\int \rho(y) \frac{f^2(y)}{1 + \|y\|^{2d}} d^d y < \infty$,
- 254 • C_{Holder}^f : *there exist $\alpha_x > 0$, $K_x > 0$, and $\delta_x > 0$, such that*
 255 *$x' \in \Omega$ and $\|x - x'\| \leq \delta_x \implies |f(x) - f(x')| \leq K_x \|x - x'\|^{\alpha_x}$*
 256 *(local Hölder condition for f).*

257 *Moreover, we define $\kappa(x) = \int \rho(x+y) \frac{f(x+y) - f(x)}{\|y\|^d} d^d y$, where we have $|\kappa(x)| < \infty$.*

258 *Then, for $\kappa(x) \neq 0$, the bias term $\mathcal{B}(x) = \mathbb{E}_X \left[\left(\sum_{i=0}^n w_i(x) [f(x_i) - f(x)] \right)^2 \right]$ satisfies*

$$\mathcal{B}(x) \underset{n \rightarrow +\infty}{\sim} \left(\mathbb{E} [\hat{f}(x)] - f(x) \right)^2, \quad \text{with} \quad \mathbb{E} [\hat{f}(x)] - f(x) \underset{n \rightarrow +\infty}{\sim} \frac{\kappa(x)}{V_d \rho(x) \ln(n)}. \quad (23)$$

259 *In the non generic case $\kappa(x) = 0$, we have the weaker result*

$$\mathcal{B}(x) = \begin{cases} O \left(n^{-\frac{2\alpha_x}{d}} (\ln(n))^{-1 - \frac{2\alpha_x}{d}} \right), & \text{for } d > 2\alpha_x \\ O \left(n^{-1} (\ln(n))^{-1} \right), & \text{for } d = 2\alpha_x \\ O \left(n^{-1} (\ln(n))^{-2} \right), & \text{for } d < 2\alpha_x \end{cases} \quad (24)$$

260 Note that $\kappa(x) = 0$ is non generic but can still happen, even if f is not constant. For instance, if Ω is a
 261 sphere centered at x or $\Omega = \mathbb{R}^d$, if $\rho(x+y) = \hat{\rho}(\|y\|)$ is isotropic around x , and if $f_x : y \mapsto f(x+y)$
 262 is an odd function of y , then we indeed have $\kappa(x) = 0$ at this symmetric point x .

263 Interestingly, for $\kappa(x) \neq 0$, Eq. (23) shows that the bias $\mathcal{B}(x)$ is asymptotically dominated by
 264 the square of $\mathbb{E}[\hat{f}(x)] - f(x)$, showing that the fluctuations of $\mathbb{E}[\hat{f}(x)] - \sum_{i=0}^n w_i(x)f(x_i)$ are
 265 negligible compared to $\mathbb{E}[\hat{f}(x)] - f(x)$, in the limit $n \rightarrow +\infty$ and for $\kappa(x) \neq 0$.

266 One can relax the local Hölder condition, but at the price of a weaker estimate for $\mathcal{B}(x)$ which will
 267 however be enough to obtain strong results for the regression and classification risks (see below):

268 **Theorem 3.6.** *For $x \in \Omega^\circ$, we assume ρ and f continuous at x , and the growth condition C_{Growth}^f :*
 269 $\int \rho(y) \frac{f^2(y)}{1+\|y\|^{2d}} d^d y < \infty$. *Then, the bias term satisfies*

$$\mathcal{B}(x) = o\left(\frac{1}{\ln(n)}\right), \quad (25)$$

270 *or equivalently, for any $\varepsilon > 0$, there exists $N_{x,\varepsilon}$, such that for $n \geq N_{x,\varepsilon}$*

$$\mathcal{B}(x) \leq \frac{\varepsilon}{\ln(n)}. \quad (26)$$

271 Let us now consider a point $x \in \partial\Omega$ for which we have $\rho(x) = 0$ (note that $x \in \partial\Omega$ does not
 272 necessarily imply $\rho(x) = 0$). In appendix A.5, we show the following theorem for the expectation
 273 value of the estimator $\hat{f}(x)$ in the limit $n \rightarrow +\infty$:

274 **Theorem 3.7.** *For $x \in \partial\Omega$ such that $\rho(x) = 0$, we assume that f and ρ satisfy the conditions*

- 275 • C_{Growth}^f : $\int \rho(y) \frac{|f(y)|}{1+\|y\|^d} d^d y < \infty$,
- 276 • C_{Holder}^ρ : *there exist $\alpha_x > 0$, $K_x > 0$, and $\delta_x > 0$, such that*
 277 $x' \in \Omega$ and $\|x - x'\| \leq \delta_x \implies |\rho(x')| \leq K_x \|x - x'\|^{\alpha_x}$
 278 *(local Hölder condition for ρ).*

279 *Moreover, we define $\kappa(x) = \int \rho(x+y) \frac{f(x+y)-f(x)}{\|y\|^d} d^d y$ ($|\kappa(x)| < \infty$ under condition C_{Growth}^f),*
 280 *and $\lambda(x) = \int \frac{\rho(x+y)}{\|y\|^d} d^d y$ ($0 < \lambda(x) < \infty$ under condition C_{Holder}^σ). Then,*

$$\lim_{n \rightarrow +\infty} \mathbb{E}[\hat{f}(x)] - f(x) = \frac{\kappa(x)}{\lambda(x)}. \quad (27)$$

281 Hence, in the generic case $\kappa(x) \neq 0$ (see Theorem 3.5 and the discussion below it) and under
 282 condition C_{Holder}^ρ , we find that the bias does not vanish when $\rho(x) = 0$, and that the estimator $\hat{f}(x)$
 283 does not converge to $f(x)$. When $\rho(x) = 0$, the scarcity of data near the point x indeed prevents the
 284 estimator to converge to the actual value of $f(x)$. In appendix A.5, we show an example of a density
 285 ρ continuous at x and such that $\rho(x) = 0$, but not satisfying the condition C_{Holder}^ρ , and for which
 286 $\lim_{n \rightarrow +\infty} \mathbb{E}[\hat{f}(x)] = f(x)$, even if $\kappa(x) \neq 0$.

287 3.2 Asymptotic equivalent for the regression risk

288 In appendix A.6, we prove the following theorem establishing the asymptotic rate at which the excess
 289 risk goes to zero with large sample size n for Hilbert kernel regression, under mild conditions that do
 290 not require f or σ to be bounded, but only to satisfy some growth conditions:

291 **Theorem 3.8.** *For $x \in \Omega^\circ$, we assume $\sigma(x) > 0$, ρ , σ , and f continuous at x , and the growth*
 292 *conditions C_{Growth}^σ : $\int \rho(y) \frac{\sigma^2(y)}{1+\|y\|^{2d}} d^d y < \infty$ and C_{Growth}^f : $\int \rho(y) \frac{f^2(y)}{1+\|y\|^{2d}} d^d y < \infty$.*

293 *Then the following statements are true:*

- 294 • *The excess regression risk at the point x satisfies*

$$\mathbb{E}[(\hat{f}(x) - f(x))^2] \underset{n \rightarrow +\infty}{\sim} \frac{\sigma^2(x)}{\ln(n)}. \quad (28)$$

295
296
297

- *The Hilbert kernel estimate converges pointwise to the regression function in probability. More specifically, for any $\delta > 0$, there exists a constant $N_{x,\delta}$, such that for any $\varepsilon > 0$, we have the following Chebyshev bound, valid for $n \geq N_{x,\delta}$*

$$\mathbb{P}[|\hat{f}(x) - f(x)| \geq \varepsilon] \leq \frac{1 + \delta}{\varepsilon^2} \frac{\sigma^2(x)}{\ln(n)}. \quad (29)$$

298 This theorem is a consequence of the corresponding asymptotically equivalent forms of the variance
299 and bias terms presented above. Note that as long as $\rho(x) > 0$, the variance term dominates over the
300 bias term and the regression risk has the same form as the variance term.

301 3.3 Rates for the plugin classifier

302 In appendix A.7, we prove the following theorem establishing the asymptotic rate at which the
303 classification risk goes to zero with large sample size n for Hilbert kernel regression:

304 **Theorem 3.9.** *For $x \in \Omega^\circ$, we assume $\sigma(x) > 0$, ρ , σ , and f continuous at x . Then, the classification
305 risk $\mathbb{E}[\mathcal{R}_{0/1}(\hat{F}(x))] - \mathcal{R}_{0/1}(F(x))$ vanishes for $n \rightarrow +\infty$.*

306 *More precisely, for any $\varepsilon > 0$, there exists $N_{x,\varepsilon}$, such that for any $n \geq N_{x,\varepsilon}$,*

$$0 \leq \mathbb{E}[\mathcal{R}_{0/1}(\hat{F}(x))] - \mathcal{R}_{0/1}(F(x)) \leq 2(1 + \varepsilon) \frac{\sigma(x)}{\sqrt{\ln(n)}}, \quad (30)$$

307 *In addition, for any $0 < \alpha < 1$, the general inequality*

$$\mathbb{E}[\mathcal{R}_{0/1}(\hat{F}(x))] - \mathcal{R}_{0/1}(F(x)) \leq 2|f(x) - 1/2|^{1-\alpha} \mathbb{E} \left[|\hat{f}(x) - f(x)|^2 \right]^{\frac{\alpha}{2}}, \quad (31)$$

308 *holds unconditionally and, for $n \geq N_{x,\varepsilon}$, leads to*

$$0 \leq \mathbb{E}[\mathcal{R}_{0/1}(\hat{F}(x))] - \mathcal{R}_{0/1}(F(x)) \leq 2|f(x) - 1/2|^{1-\alpha} (1 + \varepsilon)^\alpha \frac{\sigma^\alpha(x)}{(\ln(n))^{\frac{\alpha}{2}}}. \quad (32)$$

309 For $0 < \alpha < 1$, Eq. (32) is weaker than Eq. (30) in terms of its dependence on n , but explicitly
310 shows that the classification risk vanishes for $f(x) = 1/2$. This theorem does not require any growth
311 condition for f or σ , since both functions takes values in $[0, 1]$ in the classification context.

312 3.4 Extrapolation behavior outside the support of ρ

313 We now take the point x outside the closed support $\bar{\Omega}$ of the distribution ρ (which excludes the case
314 $\Omega = \mathbb{R}^d$). We are interested in the behavior of $\mathbb{E}[\hat{f}(x)]$ as $n \rightarrow +\infty$. In appendix A.8 we prove:

315 **Theorem 3.10.** *For $x \notin \bar{\Omega}$, we assume the growth condition $\int \rho(y) \frac{|f(y)|}{1 + \|y\|^d} d^d y < \infty$. Then,*

$$\hat{f}_\infty(x) := \lim_{n \rightarrow +\infty} \mathbb{E}[\hat{f}(x)] = \frac{\int \rho(y) f(y) \|x - y\|^{-d} d^d y}{\int \rho(y) \|x - y\|^{-d} d^d y}, \quad (33)$$

316 *and \hat{f}_∞ is continuous at all $x \notin \bar{\Omega}$.*

317 *In addition, if $\int \rho(y) |f(y)| d^d y < \infty$, and defining $d(x, \Omega) > 0$ as the distance between x and Ω , we
318 have*

$$\lim_{d(x, \Omega) \rightarrow +\infty} \hat{f}_\infty(x) = \int \rho(y) f(y) d^d y. \quad (34)$$

319 *Finally, we consider $x_0 \in \partial\Omega$ such that $\rho(x_0) > 0$ (i.e., $x_0 \in \partial\Omega \cap \Omega$), and assume
320 that f and ρ seen as functions restricted to Ω are continuous at x_0 , i.e. $\lim_{y \in \Omega \rightarrow x_0} \rho(y) =$
321 $\rho(x_0)$ and $\lim_{y \in \Omega \rightarrow x_0} f(y) = f(x_0)$. We also assume that the local solid angle $\omega_0 =$
322 $\lim_{r \rightarrow 0} \frac{1}{V_d \rho(x_0) r^d} \int_{\|x_0 - y\| \leq r} \rho(y) d^d y$ exists and satisfies $\omega_0 > 0$. Then,*

$$\lim_{x \notin \bar{\Omega} \rightarrow x_0} \hat{f}_\infty(x) = f(x_0). \quad (35)$$

323 Eq. (34) shows that far away from Ω (which is possible to realize, for instance, when Ω is bounded),
324 $\hat{f}_\infty(x)$ goes smoothly to the ρ -mean of f . Moreover, Eq. (35) establishes a continuity property for
325 the extrapolation \hat{f}_∞ at $x_0 \in \partial\Omega \cap \Omega$ under the stated conditions (remember that for $x \in \Omega^\circ$, we
326 have $\lim_{n \rightarrow +\infty} \mathbb{E}[\hat{f}(x)] = f(x)$; see Theorem 3.5, and in particular Eq. (23)).

327 **References**

- 328 [1] Holger Wendland. *Scattered data approximation*, volume 17. Cambridge university press, 2004.
- 329 [2] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An introduction to*
330 *statistical learning*, volume 112. Springer, 2013.
- 331 [3] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding
332 deep learning requires rethinking generalization. In *International Conference on Learning*
333 *Representations*, 2017.
- 334 [4] Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to
335 understand kernel learning. In *Proceedings of the 35th International Conference on Machine*
336 *Learning*, pages 541–549, 2018.
- 337 [5] Adele Cutler and Guohua Zhao. Pert-perfect random tree ensembles. *Computing Science and*
338 *Statistics*, 33:490–497, 2001.
- 339 [6] Abraham J Wyner, Matthew Olson, Justin Bleich, and David Mease. Explaining the success of
340 adaboost and random forests as interpolating classifiers. *Journal of Machine Learning Research*,
341 18(48):1–33, 2017.
- 342 [7] Mikhail Belkin, Daniel Hsu, and Partha Mitra. Overfitting or perfect fitting? risk bounds for
343 classification and regression rules that interpolate. *arXiv preprint arXiv:1806.05161*, 2018.
- 344 [8] Alexander Rakhlin and Xiyu Zhai. Consistency of interpolation with laplace kernels is a
345 high-dimensional phenomenon. *arXiv preprint arXiv:1812.11167*, 2018.
- 346 [9] Greg Ongie, Rebecca Willett, Daniel Soudry, and Nathan Srebro. A function space view of
347 bounded norm infinite width relu nets: The multivariate case. *arXiv preprint arXiv:1910.01635*,
348 2019.
- 349 [10] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-
350 learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy*
351 *of Sciences*, 116(32):15849–15854, 2019.
- 352 [11] Tengyuan Liang and Alexander Rakhlin. Just interpolate: Kernel" ridgeless" regression can
353 generalize. *arXiv preprint arXiv:1808.00387*, 2018.
- 354 [12] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in
355 linear regression. *arXiv preprint arXiv:1906.11300*, 2019.
- 356 [13] Andrea Montanari, Feng Ruan, Youngtak Sohn, and Jun Yan. The generalization error of
357 max-margin linear classifiers: High-dimensional asymptotics in the overparametrized regime.
358 *arXiv preprint arXiv:1911.01544*, 2019.
- 359 [14] Mina Karzand and Robert D Nowak. Active learning in the overparameterized and interpolating
360 regime. *arXiv preprint arXiv:1905.12782*, 2019.
- 361 [15] Yue Xing, Qifan Song, and Guang Cheng. Statistical optimality of interpolated nearest neighbor
362 algorithms. *arXiv preprint arXiv:1810.02814*, 2018.
- 363 [16] Partha P. Mitra. Fitting elephants in modern machine learning by statistically consistent
364 interpolation. *Nature Machine Intelligence*, 3(5):378–386, May 2021.
- 365 [17] Martin Anthony and Peter L Bartlett. *Neural Network Learning: Theoretical Foundations*.
366 Cambridge University Press, 1999.
- 367 [18] Luc Devroye, Laszlo Györfi, and Adam Krzyżak. The hilbert kernel regression estimate. *Journal*
368 *of Multivariate Analysis*, 65(2):209–227, 1998.
- 369 [19] Elizbar A Nadaraya. On estimating regression. *Theory of Probability & Its Applications*,
370 9(1):141–142, 1964.
- 371 [20] GS Watson. Smooth regression analysis. *Sankhya A: 26: 359-372, (50)*, 1964.

- 372 [21] Donald Shepard. A two-dimensional interpolation function for irregularly-spaced data. In
 373 *Proceedings of the 1968 23rd ACM national conference*, pages 517–524. ACM, 1968.
- 374 [22] Reinhard Farwig. Rate of convergence of shepard’s global interpolation formula. *Mathematics*
 375 *of Computation*, 46(174):577–590, 1986.

376 Checklist

- 377 1. For all authors...
- 378 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 379 contributions and scope? [Yes] Brief statement of the primary theorems are included in the
 380 abstract together with verbal descriptions of the principal results.
- 381 (b) Did you describe the limitations of your work? [Yes] Theorems include limiting conditions.
- 382 (c) Did you discuss any potential negative societal impacts of your work? [No] The paper
 383 consists of theorems and proofs about the convergence rates of an interpolation scheme and it
 384 is hard to conceive any negative social impact of this mathematical exercise. If anything it
 385 might have a positive impact by shedding theoretical light on the interpolation of noisy data
 386 in modern ML.
- 387 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them?
 388 [Yes] We have read the guidelines.
- 389 2. If you are including theoretical results...
- 390 (a) Did you state the full set of assumptions of all theoretical results? [Yes] We state the
 391 conditions for the theorems proven.
- 392 (b) Did you include complete proofs of all theoretical results? [Yes] Proofs of all theorems
 393 are provided in the appendix. The relevant appendix section is noted in the corresponding
 394 sections containing the theorems.
- 395 3. If you ran experiments...
- 396 (a) Did you include the code, data, and instructions needed to reproduce the main experimental
 397 results (either in the supplemental material or as a URL)? [N/A] The paper is directed to
 398 proving a number of theorems.
- 399 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were
 400 chosen)? [N/A] We do not rely on numerical experiments for the results presented in the
 401 paper.
- 402 (c) Did you report error bars (e.g., with respect to the random seed after running experiments
 403 multiple times)? [N/A]
- 404 (d) Did you include the total amount of compute and the type of resources used (e.g., type of
 405 GPUs, internal cluster, or cloud provider)? [N/A]
- 406 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 407 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 408 (b) Did you mention the license of the assets? [N/A]
- 409 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 410 (d) Did you discuss whether and how consent was obtained from people whose data you’re
 411 using/curating? [N/A]
- 412 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 413 information or offensive content? [N/A]
- 414 5. If you used crowdsourcing or conducted research with human subjects...
- 415 (a) Did you include the full text of instructions given to participants and screenshots, if applica-
 416 ble? [N/A] There are no human subjects.
- 417 (b) Did you describe any potential participant risks, with links to Institutional Review Board
 418 (IRB) approvals, if applicable? [N/A]
- 419 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on
 420 participant compensation? [N/A]