

1 **Reviewer #1** Thanks for your comments; we’ll clarify our usage of the following terms in the revised paper.

- 2 • We use “**consistency**” to mean  $\mathbb{E}[L_{\mathcal{D}}(\hat{w}) - L_{\mathcal{D}}(w^*)] \rightarrow 0$ . Traditionally, as in [27], this limit means  $n \rightarrow \infty$  for a
- 3 fixed problem, but in that setting linear models do not interpolate. Instead, for asymptotic interpolation we study
- 4 a sequence of distributions changing with  $n$ , with the noise magnitude  $\lambda_n$  possibly increasing. In a more typical
- 5 “high-dimensional” regime,  $p$  would also increase with  $n$ , e.g.  $p = \gamma n$  in [13]; we instead take  $p \rightarrow \infty$  for each  $n$ .
- 6 • By “**interpolation learning**” we mean achieving “good”  $L_{\mathcal{D}}(w)$  while  $L_S(w) = 0$  in a noisy, non-realizable setting.

7 **Reviewer #2** Thanks for your feedback; we’ll add more intuition, details, and reorganize proofs in revision.

- 8 • **Min-risk interpolator:** Thm. 4.5 decomposes as risk of one interpolator, plus gap to worst;  $\hat{w}_{MR}$  minimizes risk.
- 9 • **Restricted eigenvalue:** It arises naturally from (7)’s dual; it measures how of  $\Sigma$  is unobserved by  $X$ , and is the
- 10 generalization gap for  $y = 0_n$ ,  $B = 1$ . It also relates to [3]: the “malignant” covariance  $I_p$  has  $\kappa_X(I_p) \stackrel{a.s.}{=} 1$ , while
- 11 the benign covariance of Setting B has  $\kappa_X(\Sigma) \approx \lambda_n/n \rightarrow 0$ . We expect it might play the role of  $\xi_n$  in  $(\star)$ .
- 12 • **Finite degrees of freedom:** It is true that Setting B is simple in this way. Our approach also allows for analysis
- 13 where  $d_S$  increases with  $n$ , though we know it must be  $o(n)$  for consistency to be possible.
- 14 • **Consistency  $\rightarrow$  1-sided unif. conv.:** Take  $\mathcal{S}_{n,\delta} = \{(X, y) : L_{\mathcal{D}}(A(X, y)) \leq \sigma^2 + \epsilon_{n,\delta}\}$ . (We’ll clarify footnote 5.)

15 **Reviewer #3** Thanks for your writing suggestions (converting some discussion into lemmas, substantially re-focusing

16 the abstract, and clarifying e.g. line 230), which we agree will improve the presentation.

- 17 • **Portable insights:** The main takeaway we believe to be broadly relevant is that when analyzing using “uniform
- 18 convergence,” especially in the context of interpolation learning, it is important to use “relative” or “optimistic”
- 19 bounds which take  $L_S$  into account. Our approach of bounding the generalization gap via duality may also be widely
- 20 applicable: even in complex settings without strong duality, upper bounds should still be possible from weak duality.
- 21 We will emphasize these more throughout the paper.
- 22 • **Comparison to [3]/[23]:** While prior work almost fully characterizes consistency in this class of problems, it is
- 23 quite different from most existing work in statistical learning theory. Our theorem 4.5 attempts to be more like
- 24 popular Rademacher bounds, although to develop this connection further (and compare with existing conditions),
- 25 more calculations are required in general – even if the speculative bound  $(\star)$  holds. We’ll increase our discussion of
- 26 the relationship to the benign/weakly benign conditions, e.g. with the examples above. Our approach also explains
- 27 non-minimal-norm predictors, and it may be easier to numerically check  $\kappa_X(\Sigma)$  and  $\|X\|$  in practice.
- 28 • **1- vs 2-sided uniform convergence:** For predictors with  $L_S(w) = 0$ , these modes of convergence are indeed
- 29 identical. These restricted uniform bounds sidestep entirely the two-sided failure mode of Section 3.2, with high  $L_S$
- 30 but low  $L_{\mathcal{D}}$ . This is not the only difference between the standard and restricted settings, however: we strongly expect
- 31 that norm balls do not exhibit one-sided uniform convergence either (line 125), due to cases where  $L_S$  is large but
- 32  $L_{\mathcal{D}}$  is even larger. We will add more discussion of this relationship in the revision.
- 33 • **Restricted eigenvalue under interpolation:** We are not aware of any previous use of  $\kappa_X(\Sigma)$  in the literature.
- 34 • **Is low norm key?** As any low-norm interpolator generalizes, we believe we’ve shown that the answer to this question
- 35 is “yes.” We agree that this belongs in the abstract and should be highlighted more in the paper body as well.
- 36 • **Restricted convergence bounds:** As we mention around line 177, bounds like (7) are very standard in realizable
- 37 PAC analyses. Generally, (7) will *always* be small for consistent predictors – even if, as in Section 3, unrestricted
- 38 bounds fail – because taking  $B = \|\hat{w}_{MN}\|$  makes (7) just  $L_{\mathcal{D}}(\hat{w}_{MN})$ . The questions are whether we can usefully
- 39 bound the analogue of (7), and how large  $B$  can be; we answer these questions for Setting B in Section 4.1.
- 40 • **DCT in footnote 3:** Since  $L_{\mathcal{D}}$  is also an expectation, there are two interchanges of limit and expectation, and finding
- 41 a dominator seems nontrivial; it seems to essentially boil down to the proof of Proposition 4.6.

42 **Reviewer #4** Thank you for your questions and suggestions, which we will emphasize in revision.

- 43 • **Connection to deep learning:** High-dimensional linear models can serve as a simpler test bed to help develop
- 44 methods useful for deep nets: we need to walk before we can run. Knowing which techniques can explain many of
- 45 the surprising phenomena of deep learning, e.g. double descent, in linear models, helps us narrow down which tools
- 46 to try in the harder setting. (See also the *portable insights* comment to Reviewer 3.)
- 47 • **Why uniform convergence?** (1) It is in many ways *the* standard toolkit in statistical learning theory. (2) A direct
- 48 bound on  $L_{\mathcal{D}}(\hat{w}_{MN})$  may not tell us *why*  $\hat{w}_{MN}$  works; a uniform bound based on norm strongly indicates norm is the
- 49 “reason.” (3) In practice we may not find the exactly minimal-norm interpolator; uniform bounds are more “robust.”
- 50 • **Restricted problem setting:** Indeed, Theorem 4.1 is limited to a very particular setting, but we use it mainly to
- 51 demonstrate the success of our style of analysis. We emphasize that Theorem 4.5 holds quite generally.
- 52 • **Comparison to LASSO:** Here, we simply make the point that in a sparse setting, there exists a consistent learning
- 53 rule, but *no* interpolation method – including the minimal  $l_1$  norm interpolator – can be consistent for  $p = \mathcal{O}(n)$ .
- 54 • **Related statistics papers:** If you have any particular work in mind, we are eager to consider it.