

1 **Recap: What is the paper’s goal, and why?** Bound the statistical complexity of tuning SM Kernel hyperparameters,  
2 when **noise is fully adversarial** in the **active regression** setting. While extremely useful, the SM Kernel is notoriously  
3 hard to tune (L147-148). By proving that the statistical cost of tuning the SM Kernel is low, even in this hard regime,  
4 we show that we need better algorithms for hyperparameter tuning, and not to just collect more data. Also, we provide a  
5 new mathematical framework to analyze hyperparameter tuning, which hopefully opens the door to new algorithms.

6 We thank the reviewers for their constructive and broadly positive feedback. We are pleased to see the reviewers find  
7 our problem statement interesting and relevant, find our technical results to be sound and relevant to practitioners, and  
8 find our writing to be clear. We also appreciate the reviewers asking about potential future directions of our work, either  
9 for new algorithms or for a different statistical setup.

10 One important open question identified by (R1, R3, R4) is about **finding a polynomial time algorithm**. While  
11 our paper does not have a polynomial time algorithm for SM Kernel hyperparameter tuning, **it opens the door for**  
12 **more research towards that goal**. Previously, it was unknown if exponentially many samples were needed to even  
13 *information-theoretically* tune a SM Kernel. If this were the case, then there would be no hope for finding a polynomial  
14 time algorithm. We show that this barrier does not exist, which gives more hope for a fast algorithm.

15 Along the same lines, R1 asks us to “**discuss the possible direction[s] of developing the corresponding algorithm**”.  
16 This highlights another potential benefit of our paper: prior work does not frame hyperparameter tuning as a Fourier  
17 fitting problem. By rigorously and directly connecting the two settings, our paper makes it easier for a polynomial  
18 time algorithm to originate from signal processing and benefit kernel hyperparameter tuning. **The Sparse Fourier**  
19 **Transform and Compressed Sensing literatures seem promising in this regard** – notably, they both feature fast  
20 algorithms with statistical guarantees for many Fourier fitting problems. We agree that the paper should discuss these  
21 directions explicitly, and we have added such a discussion to our conclusion.

22 R2 is concerned about the **applications and broader impact of our paper**. We study kernel ridge regression with  
23 adversarial noise. There are many applications of this setting. For instance, we cite two papers that use the SM Kernel in  
24 distinct ways: [HSSM15] uses the SM Kernel in an analysis of the lifespan of lithium-ion batteries and [WDLX15] uses  
25 the SM Kernel to model human decision-making processes. Many more applications can be found in the papers that  
26 cite [WA13], which proposed the SM Kernel. So, **the SM Kernel is relevant to many applications** despite being hard  
27 to tune. Our **broader impact depends on how the SM Kernel is used**. We agree the introduction and broader impact  
28 sections would benefit from a clearer discussion of applications of our framework, and have added that to the paper.

29 R1 also brings up the **novelty and impact of our paper, and its relationship to [AKM+19]**. [AKM+19] is a recent  
30 STOC paper about Fourier function fitting in the same adversarial noise setting as our paper. They show that kernel  
31 ridge regression provides a good interpolant when the kernel is fixed (i.e. hyperparameters are *known*). The novel  
32 technical contribution of our work is the extension of [AKM+19] to the setting with an *unknown* kernel. This allows us  
33 to frame kernel hyperparameter optimization as a Fourier fitting problem, so [AKM+19] **lets us prove the first bounds**  
34 **on the statistical complexity of learning kernels with totally adversarial noise**. Mathematically, this comparison is  
35 made clear in our introduction (Problem 1 versus Problem 2).

36 R4 asks **if our results can apply to the fixed data setting**. This is an interesting question, but answering it requires a  
37 different noise model than our totally adversarial noise model. If the data was fixed, then an adversary could arbitrarily  
38 perturb the signal  $y(t)$  only at the observed times. So, a small-norm noise signal would make approximate recovery of  
39  $y(t)$  totally impossible. Instead, perhaps the adversary could only perturb an  $\varepsilon$  fraction of the observations? We think  
40 this is an interesting direction, and it is plausible that our techniques could generalize such a fixed-kernel result into a  
41 hyperparameter optimization result, but this speculation is out of scope for our paper.

42 R2 asks for a “**intuition behind the statistical dimension**” and **how it varies between kernels**. Kernels are nonpara-  
43 metric functions that can involve infinitely many features (in our paper, a feature of a kernel is like the Fourier transform  
44 of a kernel at a given frequency). Nevertheless, in Kernel Ridge Regression, thanks in part to regularization, kernels can  
45 typically be well approximated with a finite number of features (this is the idea behind e.g. Random Fourier Features).  
46 Statistical dimension captures exactly how many features are needed. Since regression in  $d$  features requires roughly  $d$   
47 samples, this means kernel ridge regression needs roughly statistical dimension many observations. For the RBF kernel,  
48 statistical dimension is linear in the lengthscale parameter and the duration of time we interpolate over. The intuition is  
49 that if either the lengthscale or the duration of time increase, then we can represent more complex functions with the  
50 same Fourier-norm, so this necessitates more samples (i.e., larger statistical dimension).

51 R2 asks for **Clarification about  $O(1)$  and  $C$  in our bounds**. These all refer to a universal constant independent of  
52 any problem parameters such that the left hand side is bounded by the constant times the right hand side. We thank R2,  
53 have updated our paper to always use the  $C$  notation, and always describe  $C$  as a universal constant. R2 asks us to  
54 include a **better intuition for the Energy term in the introduction**, beyond what is on lines 33-35. We thank R2 and  
55 agree; we have now included a stronger mathematical intuition for the Energy term in the introduction.