

1 We would like to thank all reviewers for their careful reviews and insightful comments.

2

3 **Reviewer 2: It seems that the algorithm needs to have a pretty large tree. Will it increase in exponen-**
4 **tial order of T or other order?**

5 The space grows only linearly in time because the algorithm constructs the tree dynamically, adding a new path of size
6 $\mathcal{O}(D)$ in each round. In particular, the algorithm never allocates the entire tree, but only the paths corresponding to
7 active experts.

8 **The algorithm can only adapt to one of the mentioned local regularities, and it really just depends on which**
9 **radius tuning function is used.**

10 This is correct, and we will modify the abstract to remove any ambiguity. Still, we believe that being able to adapt to
11 different local regularities just by supplying a specific radius tuning function is interesting. Note also that, even though
12 we are not aware of any online nonparametric algorithms with simultaneous adaptivity (not even to global regularities),
13 it is certainly possible to combine algorithms with different types of adaptivity using online aggregation techniques
14 such as prediction with expert advice.

15 **The paper omitted many details, which can be found in the appendix. In the Algorithm 1, the paper omitted**
16 **the update for inactive experts after line 10.**

17 This is not an omission: the algorithm does not require any updates of inactive experts, only experts along the path are
18 updated, which makes it computationally attractive. This is one of the beauties of the context tree weighting method,
19 originally described in the information theory literature (see [29]).

20 We apologize that some important material had to be placed in the appendix, we will certainly bring back the important
21 details you mention to the main body in the revised version of the paper.

22 **The local online predictor for each node is using Follow-the-Leader rule. It kind of makes me confused on how**
23 **to do that with only one node information.**

24 It would have been more appropriate to call it “follow-the-local-leader”. Indeed, each follow-the-leader instance hosted
25 in a ball learns only over the subsequence of examples that fall in that ball.

26 **Improvement in the local dimension case is doubtful — the oblivious adversary can generate the instance se-**
27 **quence that does not have lower local dimension.**

28 Note that our algorithms are deterministic (no distinction between oblivious and nonoblivious adversaries), and our
29 regret bounds hold for any sequence. We do not claim that our bound is always better. In fact, when the local dimension
30 equals the global dimension, we merely do not lose anything in the exponent as our bound becomes of order $\tilde{O}(T^{\frac{d}{1+d}})$.

31 **The paper assumes that the instance space is the unit ball. What about the case when the instance space is not**
32 **necessarily a unit ball? Can the algorithm still work in this case? What changes should we make?**

33 We gain a factor of C^d in the regret (note that this is not improvable), where C is the radius of an instance space. At the
34 same time we have to know C since it will appear in the radius tuning function (see, e.g. [13]).

35 **Reviewer 1: It seems that we need to specify the hierarchical levels (e.g Lipschitz constants, metric dimension-**
36 **ality) which sound non-trivial.**

37 The hierarchical levels are hyperparameters determining the reference or comparator class of the nonparametric problem
38 (which in turn defines the inductive bias). Note that the algorithm and its analysis work for *any* choice of these
39 hyperparameters. By choosing specific values, the user controls the trade-off between size of the reference class on one
40 side and the growth of the regret bound and space/time requirements on the other side.

41 **It is not entirely clear whether the additional (algorithm) complexity and input parameters justify the gains,**
42 **and if so, on what problems.**

43 Note that the overhead in running time is logarithmic in the size of the tree, since at each round we only query and
44 update experts along the path of an instance. As for the space requirements, the tree grows only linearly in time as the
45 algorithm constructs the tree dynamically (see response to R2).

46 **Reviewer 3: The techniques are currently only analyzed for regression with square loss and classification with**
47 **absolute loss. It would have been nice if the algorithms can be applied beyond these two settings.**

48 Our proofs can be easily extended to any exp-concave losses, as these do not require any tuning of learning rates in
49 local FTL predictors. We also believe that it is also possible to extend our approach to any convex loss by using a
50 parameter-free local learner such as Squint (Koolen and van Erven, COLT 2015).