1 Thanks to the reviewers for their constructive comments.

2 **Reviewer #1.** Thanks for the positive feedback. We agree that a multidimensional experiment would be beneficial for
3 the main paper. Note that the paper already has a multidimensional experiment in Appendix E, that will be moved at the
4 end of the experimental section. Moreover we will add some details on the computational cost and fast algorithms to
5 solve the dual problem in Thm. 2 to the end of Sect. 3 (they are now at the end of Appendix B.5).

6 **Reviewer #2.** We thank the reviewer for the careful reading and thoughtful comments.

7 1) As correctly pointed out by Rev2, a model of the form $f_w(x) = (\phi(x)^\top w)^2$ is sufficient for Theorem 3 (and also for
8 Theorem 4). Such formulation would lead to non-convex optimization problems as already noticed by the reviewer.
9 Of course if convexity is not an issue, it is possible to restrict the model to rank-1, keeping all the results of Section
10 4. However, surprisingly, note that the dual formulation of the problem (Thm. 2) shows that the matrix $A$ also admits
11 a representation in terms of only $n$ degrees of freedom, instead of $n^2$ (given $n$ coefficients $\alpha$, the operator $A$ can be
12 recovered via Eq. 11 and then Eq. 7). Thus using the whole matrix $A$ allows on the one hand to still have a representation
13 in terms of $n$ degrees of freedom. On the other hand it leads to a convex problem and allows to control the rank of $A$
14 via elastic-net regularization. To conclude, we agree that the question raised by Rev2 can be useful to better understand
15 the value of the proposed approach. So we will add all the reasoning above as a remark right after Thm. 2.

16 2) We would like to point out that we study explicitly the variance of the proposed method in Theorem 5 where we
17 bound the Rademacher complexity of the proposed estimator. Indeed, by using a standard argument based on the
18 Rademacher complexity (see [29] Chapter 26, or [3] paragraph 4.5 and in particular Eq. 13) we can derive the following
19 learning rate. Let the population risk be defined as $R(f) = \mathbb{E}_{x,y}\ell(y, f(x))$ for some $G$-Lipschitz loss function $\ell$ and
20 $\widehat{R}_D$ be the empirical version $\widehat{R}_D(f) = \frac{1}{n}\sum_{i=1}^n \ell(y_i, f(x_i))$ for a given dataset $D$ of $n$ examples. Given a norm $\|\cdot\|_\circ$
21 (e.g., Frobenius or nuclear), a feature map $\phi$ and a radius $L$, define the class of estimators $\mathcal{F}_{\phi,L}^\circ = \{f_A \mid \|A\|_\circ \leq L\}$.
22 Denote by $\widehat{f}_{D,L} = \arg\min_{f \in \mathcal{F}_{\phi,L}^\circ} \widehat{R}_D(f)$ the empirical risk minimization solution over the set $\mathcal{F}_{\phi,L}^\circ$, then

$$\mathbb{E}_D R(\widehat{f}_{D,L}) \leq \inf_{f \in \mathcal{F}_{\phi,L}^\circ} R(f) + 2\mathbb{E}_D\Big[\sup_{f \in \mathcal{F}_{\phi,L}^\circ} |R(f) - \widehat{R}_D(f)|\Big] \leq \inf_{f \in \mathcal{F}_{\phi,L}^\circ} R(f) + 2G\,\mathcal{R}_n(\mathcal{F}_{\phi,L}^\circ),$$

23 where $\mathcal{R}_n(\mathcal{F}_{\phi,L}^\circ)$ is the Rademacher complexity of the set $\mathcal{F}_{\phi,L}^\circ$ and is bounded by $O(Lc^2/\sqrt{n})$ by Theorem 5 ($c$
24 is the bounding constant of the kernel, i.e., $c = \sup_{x \in \mathcal{X}} \|\phi(x)\|$). Now, assuming that there exists an operator
25 $A_\star$ with $\|A_\star\|_\circ$ finite (in particular it could be rank-1, i.e. $A_\star = w_\star w_\star^\top$ for some $w_\star$), such that the learning
26 problem is well posed, i.e. $\inf_{f \in C(X)} R(f) = R(f_{A_\star})$, and choosing $L = \|A_\star\|_\circ$, we obtain the learning rate
27 $\mathbb{E}_D R(\widehat{f}_{D,L}) - R(f_{A_\star}) = O(c^2 G \|A_\star\|_\circ/\sqrt{n})$, that is comparable to the one of kernel linear models [29]. We will add
28 this paragraph as a discussion after Thm. 5, to better clarify the variance and learning rates for the proposed approach.

29 3) The NCM approach, i.e., approximating a function with non-negative combination of non-negative kernel functions
30 is employed usually in kernel density estimation methods. This model is well known to be quite rigid when the kernel
31 function is non-negative, indeed it can't approximate a density with i.i.d. samples faster than $n^{-1/(d+1)}$ even if the
32 density is arbitrarily smooth (see, e.g., [33]). This happen also when the density is of the form $f_\star(x) = e^{-V(x)}$
33 with $V(x)$ an infinitely smooth potential. Instead, in this case, according to the Point 2) above, the proposed method
34 has a faster learning rate $O(\|w_\star\|^2/\sqrt{n})$, where $\phi$ is the Sobolev kernel, since by Thm. 4 there exists a $w_\star$ s. t.
35 $f_\star(x) = (w_\star^\top \phi(x))^2$. We will add this example in Section 4 to clarify the difference between NCM and our method.

36 4) As suggested by the reviewer we will add a quantitative version of the experiments in the appendix, with 50 i.i.d.
37 repetitions and the resulting error bars. In Appendix E are reported many experimental details. We will add in the main
38 text a more detailed description on how we performed cross validation on both $\sigma \in [10^{-3}, 10^3]$ and $\lambda \in \{0\} \cup [10^{-8}, 1]$
39 (logarithmic scale, 20 intervals). We would like to note that in the multivariate experiment ($d = 10$), in the appendix,
40 best parameters found by cross-validation for NCM are $\sigma = 0.5$ and $\lambda = 10^{-4}$. Since the value of $\sigma$ is not on the left
41 extreme of the range we would exclude that the result obtained by NCM in this experiment was due to oversmoothing.
42 We instead interpret such result in the light of the different learning rates $O(n^{-1/11})$ for NCM and $O(n^{-1/2})$ for our
43 method. In any case we will repeat the experiment (with 50 repetitions) and on a range $\sigma \in [10^{-10}, 10^{10}]$ (log scale
44 with 100 steps). To conclude we would like to recall that we will publish the code on GitHub (python + scipy).

45 **Reviewer #3.** We thank the reviewer for the positive feedback. As recalled in the Point 1) and 2) above, the model can
46 be expressed in terms of $n$ degrees of freedom via duality in Thm. 2. In terms of statistical complexity we achieve a
47 learning rate that is similar to the one for linear models (see Point 2). We will add these discussions after Thm. 2 and 5.
48 More details on the computational complexity of the dual formulation (Thm. 2) are at the end of Appendix B.5.

49 **Reviewer #4.** Thanks for the thoughtful comments. An extensive explanation about the relation with [4] and its intrinsic
50 limitations is already reported in Appendix F. We will move part of the content in Section 2, where we will also add a
51 short review about main result on SoS programming.