

1 **#ToR1. Q1: ["Weaknesses...-METHOD"].** **A1:** Thanks for your valued and constructive comments. The main
2 motivation of our manuscript is to propose robust additive models (with theoretical analysis and applications) for
3 realizing nonlinear estimation and structure discovery simultaneously, even data contaminated with complex noise
4 and without priori knowledge of variable structure. There are **key differences** between our MAM and [11,37,34]:
5 1) *Hypothesis space (nonlinear VS linear).* The proposed MAM employs the spline-based additive hypothesis space,
6 which is more flexible than the linear assumption for regression function [11,37,34], and leads to better performance
7 for characterizing the nonlinear relationship in data (e.g., [13,15,21,29,41]). 2) *Sparsity within group ($\nu \in [0, 1]^P$ VS*
8 $\nu \equiv \mathbf{I}_P := (1, \dots, 1)^T$). The flexible selection of ν in our Outer Problem can unveil main effect variables across all
9 tasks, which is useful to remove ambiguity for model identifiability. Indeed, [11] ignores the sparsity within group
10 and [37,34] don't consider the structure discovery. 3) *Optimization (bilevel non-convex optimization VS bilevel convex*
11 *optimization).* The non-convexity of outer and inner problems is the main challenge for MAM's optimization, which
12 can not be tackled directly by the optimization strategy for bilevel convex optimization problem[11]. A computing
13 algorithm for bilevel non-convex scheme is proposed by developing an effective optimization method for the non-convex
14 non-smooth inner problem. 4) *Application.* As we known, MAM is the first attempt to the interpretable CMEs analysis
15 by data-driven structure discovery strategy. In short, our work is new for structure prediction, and can't be established
16 by easily combining [11] and [37]. We will add these remarks after Table 1 (line 69).

17 **Q2: ["Weaknesses...-ALGORITHM"].** **A2:** The non-convexity of outer and inner problem is the main difficulty in
18 the optimization. To tackle this problem, an optimization strategy is proposed based on building blocks (ProxSAGA
19 [30], HQ optimization[24] and DFBB[36]). Theorem 1 in Supp C guarantees that our computing algorithm can obtain
20 the local optimal solution. Indeed, it's always a challenge to give the complexity analysis on bilevel optimization (we
21 leave this direction for future work). Following your suggestion, Remark 2 (line 126 in Supp) is added to illustrate the
22 extra complexity (computation time and space) of our optimization algorithm (compared with BiGL [11]).

23 **Q3: ["Correctness...-EVALUATION"].** **A3:** Thanks for your valued suggestions. Some competitors (BiGL[11],
24 Lasso[35], RMR[37], SpAM[29]) with detailed analysis are added in Table 3 and Supp D. The AAE for each task is
25 BiGL(11.09h, 59.75km/s, 63.51km/s, 95.97km/s, 4.77nT), Lasso(12.16h, 62.56km/s, 59.81km/s, 95.34km/s,
26 4.38nT), RMR(12.02h, 58.65km/s, 49.03km/s, 98.13km/s, 4.20nT) and SpAM(10.02h, 50.48km/s, 54.97km/s,
27 72.33km/s, 3.88nT), respectively. We can observe that MAM (in Table 3) can achieve smaller AAE than other
28 competitors due to its ability of robust nonlinear approximation and structure discovery.

29 **Q4: ["Clarity-It..."].** **A4:** We have polished the manuscript following your suggestions and will proofread it carefully.

30 **Q5: ["Relation to .."].** **A5:** In Table 1, MAM has been related to [11,37,26,41] in terms of hypothesis space, evaluation
31 criterion, robustness and sparsity. Additionally, we have enriched Table 1 by adding the (non)convexity of objective
32 function, and added remark to discuss the extra complexity of MAM compared with methods[11,37].

33 **Q6: ["Additional feedback"].** **A6: [-When...]:** mGAM with an oracle variable structure is the baseline of MAM.
34 For easy reading, we have added Remark 3 (pg.5) to state the relation of mGAM and MAM. **[-There...]:** τ_j is the
35 weight for each group, by controlling which we can emphasize some groups that are known a priori. **[-Line 53]:** "In
36 theory,..." -> "In theory, we provide the convergence analysis of the proposed optimization algorithm". **[-Line 143-144]:**
37 The first term is used for robust estimation. The second term guarantees the strong convexity of transformed inner
38 objective (See Eq.4 in Supp). The third term incorporates the variable structure ϑ . **[-Line 173]:** We have corrected $\varepsilon \rightarrow \mu$,
39 where $\varepsilon = n\sigma^2\mu$ is the penalty parameter in the transformed inner problem (line 26 in Supp). Note that ε is different
40 from the noise ϵ . **[-It...]:** A description (line 152) is added to explain that the second set of gradients is computed by
41 backpropagating through the iterations of a smooth algorithm. **[-What...]:** The information appearing in parenthesis is
42 standard deviation. We have revised the title of Table 2. **[-In...]:** We have revised Equation 2.

43 **#ToR3. Q1: ["Weaknesses..."].** **A1:** Thanks for the constructive comments. In Supp, Section E (pg.12) is added
44 to discuss the generalization bound based on algorithmic stability. Firstly, we proved that the mode-induced metric
45 satisfies Quadratic Growth condition [Charles 2018] under mild assumptions. Based on the definition of multi-task
46 uniform stability and Theorem 2 [Zhang 2015], we quantify the stability of MAM and establish its generalization bound.
47 In addition, a brief comparison between overlapping group lasso[16] and our MAM is provided in Supp.

48 **Q2: ["Clarity..."].** **A2:** 1) We have added more key parameters in Table 4 and a flowchart of optimization before
49 starting with Section B.1. 2) The description ("Even...") is added in line 196. 3) Some competitors are added in CMEs
50 experiment (See A3 of #ToR1). 4) We have corrected these typos and will proofread the whole manuscript.

51 **#ToR5. Q1: ["Weaknesses..."].** **A1:** A detailed description of the selection of group number L and $|V|$ is added
52 in line 171(pg.6). In fact, we have discussed the impact of group number L in Fig 8, which indicates a satisfactory
53 result can also be obtained even if L is set to be larger than oracle group number. Moreover, our method can realize
54 data-driven variable selection without the inactive variables $|V|$ being set initially (e.g., the red pixel in Figs 2 and 3).

55 **Q2: ["..discuss..."].** **A2:** In hyper-parameter selection problem (e.g., variable structure ϑ), bilevel framework has better
56 generalization than traditional hyper-parameter selection method (e.g, Cross-validation), since the hyper-parameter is
57 tuned over a continuous space by minimizing the outer problem [K.P. Bennett et al., IJCNN 2006]. According to your
58 suggestions, a further discussion on the advantage of bilevel framework is provided in Introduction (line 42).