

1 - We would like to start by thanking the authors for their comments and suggestions.

2 **Reviewer # 1:**

3 - We also believe that these two elements are the most important points that should be augmented. We will add an
4 extended discussion about the condition C1 and its relation to the conditional Fisher information matrix and a detailed
5 discussion with respect to previous work as a conclusion. This will be possible using the 9th page in the camera-ready
6 version of the paper

7 **Reviewer # 2:**

8 - In this paper, we have chosen to focus solely on the theoretical aspects of this algorithm as we believe that a serious
9 and exhaustive numerical study of the algorithm's performance deserves a dedicated study. For pairwise and binary
10 variables, [12] shows that RISE beats conditional likelihood and is fast to run (although in [12], the optimization is
11 performed using a second-order method which is much less efficient than what we propose). We expect GRISE to be a
12 very efficient algorithm in practice as it contains RISE as a special case.

13 - The efficient way to verify the identifiability condition is currently in the Appendix D1 and not in the main text. If
14 space permits, we will put it back in the main text.

15 - Thank you for the list of typos, we will amend them in the final submission.

16 **Reviewer # 3:**

17 - It may look like our paper is a straightforward generalization of [17], but only superficially. This is the result of a
18 deliberate choice on our part to find a compromise to facilitate the comprehension of this study for readers who are
19 both familiar and unfamiliar with the field. We understand that this choice may have the unfortunate consequence of
20 downplaying the innovations presented in this paper. For instance, the local learnability condition looks like it is similar
21 to the restricted strong convexity condition from [17] but this is not the case. Strong convexity (restricted or not) is
22 related to the notion of curvature of a function and is therefore linked to the standard Euclidian ℓ_2 -norm. It turns out
23 that the GRISE loss function is not strongly convex independently of the dimension of the problem (this implies in turns
24 that the number of samples cannot be $\log p$ according to the standard theory). With the local learnability condition, we
25 are showing that by looking at the problem using different types of geometry (in our case it is mostly the ℓ_∞ geometry)
26 we can bypass this issue. This approach has been demonstrated using GRISE, but it is in fact an important result in
27 itself for it is applicable to the analysis of any estimator that minimizes an empirical risk. We also provide a recipe
28 to efficiently verify this new concept of convexity that is easily applicable to other cost functions. With respect to
29 other innovations presented in this paper, we would also mention Algorithm 1 which uses a recursive call to GRISE. A
30 straightforward generalization of [17] will not work for reconstructing graphical models precisely because such a loss
31 function is not strongly convex independently of the dimension.

32 - The parameter $\hat{\gamma}$ is a prior on the strength of the parameters of the problem. This implies that it should be bigger
33 than the actual γ of the system. So unlike the λ from [17], $\hat{\gamma}$ does not depend on the size of the system or on the
34 number of samples. Ideally one would like to take $\hat{\gamma} = \gamma$, but γ is an unknown quantity a priori. This implies that in
35 practice we would chose $\hat{\gamma}$ much larger than the expected γ . Happily we show that this has only a limited impact on
36 the sample complexity and the runtime of the algorithm as they both scale polynomially with $\hat{\gamma}$ while it is exponential
37 in γ (which is the best one can hope for). Regardless of the size of the problem, in practice one would not expect γ
38 to be bigger than 15 as it means that the intrinsic randomness of the variables is less than 1ppm. There are several
39 reasons for switching from a penalty (the λ form) to a constraint (the γ form). The main reason is that the penalty is
40 designed to enforce sparsity and is adapted to graphical models with bounded degrees. By proposing a constrained
41 version, we can reconstruct graphical models irrespective of their degree while in [17] dense graphs can require up
42 to $\mathcal{O}(p^4 \log(p))$ samples. Another reason is that the very nice property of λ in [17] being independent of some of the
43 problem parameters is lost for more general models. In fact λ becomes dependant of the parameter γ for graphical
44 models that are not pairwise and binary.

45 - The graphical models considered here can have a more general form (arbitrary order, arbitrary alphabet, arbitrary
46 parameterization, arbitrary graph structure) while in [17] it is only applicable to pairwise models with bounded degrees
47 with binary variables using monomial basis functions. The tools for the analysis may looks superficially similar to [17]
48 at first glance but are in fact based on a radically new concept that does not use convexity.

49 **Reviewer # 4:**

50 - In our study we focus essentially on the theoretical aspect of GM learning (we believe that the current paper is already
51 very dense). We plan to perform an exhaustive numerical study of this algorithm in a future work.

52 - In algorithm 1, L is not exactly an input parameter for it is determined by the family of models that is under
53 consideration and is given by the highest order of the basis functions in the family.