

1 We would like to thank all of the reviewers for their valuable time and their constructive comments. In what follows, we
2 would like to address all concerns raised.

3 **Reviewer 1:** We will incorporate the proposed minor corrections in the final version of the paper. We thank the reviewer
4 for proposing the interesting comparisons, and we will include more comprehensive discussion in the final paper. Below
5 please find some details regarding the proposed methods.

- 6 • The two-stage approach, i.e., *i*) running gradient descent to convergence, and then *ii*) projection onto sparsity set,
7 is known to be sub-optimal even for simple problems such as least-squares with sparsity constraints. In fact, the
8 results when using such approach on ℓ_2 -norm minimization are the same as the Greedy baseline: It *i*) converges to
9 the global optimum $q(\cdot)$, and then *ii*) use greedy projection to try to minimize $F[p(\cdot)] = \|p(\cdot) - q(\cdot)\|_2^2$ subject to
10 sparsity constraint, which has been shown to be inferior than IHT (subsection 4.1).
- 11 • In the general case, taking the k -heaviest coordinates of q (without assuming any structure) would result into a
12 non-valid putative solution; recall, by definition of the discrete setting, we have n coordinates, each of which takes
13 m points, leading to a m^n sample space. Simply taking the k -heaviest coordinates of that long vector would result
14 into an intermediate representation of non-zero positions that does not correspond to a probability distribution. A
15 variation of this approach is fine for the "vector-sparsity" special case.

16 On whether support set changes during iterations, we observe that in experiments (subsection 4.1) IHT changes support,
17 on average, 36.1 times for ℓ_2 -norm minimization and 6.9 times for KL minimization. We also conduct experiments
18 on fixing the support after some iterations: IHT on ℓ_2 -norm minimization (subsection 4.1) after 400 iterations, fixing
19 supports after 1, 5, 10 and 15 support sets change, give average results of 0.0026, 0.0020, 0.0018, 0.0016, respectively.

20 Regarding motivation, model compression is an exciting immediate application of our proposed approach, especially
21 since our general approach may be flexibly applied to specialized problem-specific losses. We are currently investigating
22 extensions of the current work to model/policy compression for reinforcement learning, where the loss can be constructed
23 to preserve post-compression expected reward; but the details of that approach deserve a different publication.

24 **Reviewer 2:** We thank the reviewer for the supportive and constructive review. Regarding the comment in lines
25 211-214, we chose not to include this detail as extending convergence analyses from global to local strong convexity is
26 considered fairly standard; see e.g., [Agarwal2010].

27 Regarding the comment in lines 198-202, we apologize for any confusion. This argument is a standard one in NP-
28 hardness proofs: proving the existence of corner cases is sufficient to show the hardness of the problem. Note that even
29 if *empirically* we observe that common cases are easy to solve, this does not guarantee this is the norm. Showing that
30 with high-probability the corner cases rarely appear is an interesting question by itself.

31 Regarding variance in experiments, we have observed high variance is not enough for the algorithm to get "lucky".
32 In fact, we observe that the **best** performance of Greedy is often worse than the **worst** performance of IHT. Low
33 variance is especially desirable when the algorithm is only applied a few times to save computation, as in large discrete
34 optimization problems. We believe that this makes IHT preferable in practice.

35 We also thank the reviewer for the suggestions on free energy / ELBO using IHT; we will consider these as future work.

36 **Reviewer 3:** We thank the reviewer for the constructive comments. We are happy to fully describe the structure nesting
37 i.e. why \mathcal{D}_k is included in \mathcal{D}'_k in general. QM-AM stands for Quadratic Mean-Arithmetic Mean inequality.

38 Unfortunately, we are unaware of methods that are strictly better than greedy ℓ_2 -norm projection. However, trading time
39 for performance may be an interesting topic for future work. Regarding the novelty, although we unveil a relationship
40 between the functional derivative and the gradient (l. 408), they are different in many aspects. The ℓ_2 -norm projection for
41 vector sparsity can be done optimally, but we have shown that ℓ_2 -norm projection for the general case is computationally
42 hard in general (Theorems 1 & 2). Despite the provable hardness, we provide some intuition for why greedy projection
43 is effective in our main Algorithm (Theorem 3). These are challenges of extending vector optimization to general
44 distribution optimization. We also note that the convergence analysis is for the functional setting, not the vector setting
45 (Theorem 4), which is not only more general, but also paves the way for future work on continuous distributions.

46 Regarding **why** greedy seems effective in practice, we have provided the intuition and its supporting theorem in section
47 3.4. Regarding the greedy projection (Algorithm 2), it is also possible to use KL divergence as distance metric, but our
48 convergence analysis (Theorem 4) suggests that a good projection in ℓ_2 -norm (Definition 5) is preferable.

49 We will improve the presentation as suggested in the final version of the paper.

50 [Agarwal2010] Agarwal, Alekh and Negahban, Sahand and Wainwright, Martin J, "Fast global convergence rates of
51 gradient methods for high-dimensional statistical recovery", Advances in Neural Information Processing Systems,
52 2010.