

1 **Reviewer #1:** We thank the reviewer for his/her review and suggestions. **Practicality and claims:** In term of practical
2 use of the EM training we agree and explicitly acknowledged that it is computationally demanding. However, the
3 obtained analytical results hold for any depth/with and nonlinearities (as long as they are piecewise affine); the results
4 of the paper are thus general and can be used to gain in depth theoretical understanding of generative networks and
5 their learning dynamics (from the explicit M step). The obtained analytical forms allow (i) to better design VAEs (now
6 knowing the target posterior that variational inference approximates), (ii) to guide the design of variational distributions
7 (for example favoring full covariance and multimodal posterior) as well as (iii) interpreting the learned parameters from
8 the M-step. Those insights are gained despite the computational limits of the practical EM learning as they rely on
9 the analytical derivations only. We will add further analysis and discussions on this. On that note, we also believe
10 that tremendous future work can also be done to derive faster EM learning leveraging the obtained formula either by
11 providing principled approximations of the per region Gaussian integrals or by approximation of the partitions; we
12 believe that this paper is only the beginning of such research directions. **Direct log-likelihood maximization:** As in
13 any missing variable model (here z is unobserved, only x is observed) one can not directly minimize the negative
14 log likelihood and must first infer z . EM is one common strategy to do so based on the posterior $p(z|x)$; once z
15 is inferred, one can then do the maximization of the (now estimated) log-likelihood. **Adding cost analysis of the**
16 **algorithm:** We will add in the appendix exact computation times and further details for each of the experiment and
17 different architectures for the EM learning as well as each step involved (partition finding, region triangulation, per
18 region integration). **Comparison to VAE and figures:** In Fig. 4 we compared the negative log-likelihood of both
19 models. While a VAE can be trained using the standard variational inference strategy, we evaluate its NLL after training
20 and compare with the generative deep network trained with EM. We will explicit this in the caption. We will also make
21 the legend and labels clearer in the figures.

22 **Reviewer #2:** We thank the reviewer for their appreciation of the paper. We will correct the typos and explicit the
23 pseudo-code as well as providing exact link with the implementation. **Computational limitations:** Indeed, the current
24 analytical EM learning is computationally demanding, we believe that future work can be done on this point by (1)
25 providing analytical form of gaussian integration on a convex polytope (this would remove the need of triangulation and
26 then inclusion-exclusion formula) or by (2) providing principled approximation of those integrals. Note that our main
27 contributions are the analytical derivations of the probability distributions and EM formula, the practical EM learning
28 demonstrates the usefulness of those derivations. **Gaussian prior and piecewise affine nonlinearities:** The review
29 is correct; this only applies to Gaussian prior and output distributions and with DN employing spline operators like
30 ReLU, leaky-ReLU, abs. value, . . . which includes a large part of current generative network architectures. Also, the
31 proposed method (with exact partition and per region derivation) can be employed to different distributions as long
32 as they are conjugate priors. We will add this note in the paper. **Constant covariance:** Indeed, this case covers the
33 practical cases of training in current generative models, however more general cases could be considered and even
34 different distributions. We believe that the proposed methodology (per region derivation) provides a general framework
35 and as long as the prior and output distributions are conjugate priors, analytical forms should be obtainable. We will
36 add this discussion in the paper.

37 **Reviewer #3:** we thank the reviewer for their careful review and appreciation of the paper. **Previous work:** We thank
38 the reviewer for this relevant reference (which we denote by ICML2019 thereafter). ICML2019 relates linear VAEs
39 to PPCA and propose a mode approximation of the posterior in turn producing a novel type of VAEs (Laplacian
40 VAEs). ICML2019 also provides insights into the manifold geometry (piecewise affine) of ReLU VAEs. We will
41 add this reference and detailed review in the background section. However we believe that none of our contributions
42 is over-shadowed by ICML2019 since: (i) we extend the PPCA link of linear VAES to nonlinear VAEs resulting in
43 MPPCA; (ii) we extend their geometrical insights to piecewise affine nonlinearities (not only ReLU) which consequently
44 also allow to apply ICML2019 approximation methods to a broader class of VAEs; (iii) in ICML2019, no analytical
45 (explicit) form is given for the probability distributions of a nonlinear VAE as the motivation of the paper was to provide
46 a mode approximation based on a linearization of the network to tackle large scale tasks. We will also discuss the paper
47 approximation method in the future work section as such posterior mode estimation could be employed and potentially
48 improved with the proposed distributions. **Lemma 2 ReLU assumption:** you are correct, Lemma 2 holds for more
49 general DGNs (as long as there is no surjectivity), we will add this note and discuss such cases in the paper.

50 **Reviewer #4:** We thank the reviewer for his/her appreciation of the paper and we agree that providing exact methods
51 even with demanding computational cost is crucial to exactly measure the impact of current approximation methods
52 in VAEs. **Computational complexity discussions:** indeed, the computational bottleneck comes from the number of
53 regions that then need to be triangulated. We will add computational time of each of the involved steps in the appendix:
54 (i) computation of the partition, (ii) triangulation of each region (on average) and (iii) integration on a region. We will
55 provide those statistics for the few different topologies that were used in the paper. Concerning the rate of growth of the
56 number of regions in a real network, we will add citations to the following papers: "Complexity of Linear Regions
57 in Deep Networks", "On the Number of Linear Regions of Deep Neural Networks" and "A Spline Theory of Deep
58 Networks" with discussions.