

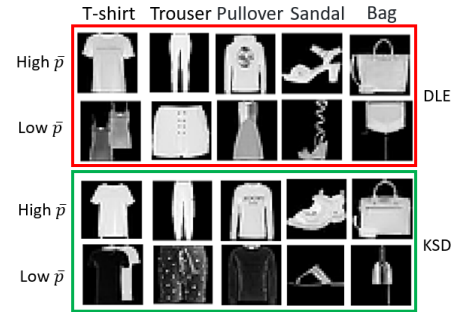
1 We thank reviewers for their insightful comments and here we would like to address some questions raised in the review.

2 **R1: “Consistency results are given but these assume the parameter space is compact (and other not so simple**
 3 **to check assumptions).”** The compactness condition is only used to define a domain on which Assumption 2 and
 4 4 hold. If Assumption 2 and 4 are defined over a neighbourhood region of the true parameter, we can remove the
 5 compactness condition by adding an extra proof which shows $(\hat{\delta}, \hat{\theta})$ eventually fall in such a neighbourhood, but doing
 6 so would introduce further technical complications. The compactness condition is among a set of conditions commonly
 7 used in classic consistency proofs (see e.g., Wald’s Consistency Proof, 5.2.1, van der Vaart, 1998). It is possible to
 8 derive weaker conditions given specific choices of f or $p(x; \theta)$. However, in the current manuscript, we only focus on
 9 more *generic* settings and conditions that would give rise to estimation consistency and useful asymptotic theories.

10 **R1: “... though this further assumes a (fairly strong) condition of uniform convergence ...”** The uniform
 11 convergence on Hessian is needed to control the residual of the asymptotic expansion (eq. 25, 26) and is a slight
 12 modification of a classical regularity condition on the uniformly bounded third order derivative (5.3, van der Vaart,
 13 1998). Again, this assumption may be weakened given specific choices of f and $p(x; \theta)$ but we focus on investigating
 14 generic settings where specific choices of f and $p(x; \theta)$ are not available.

15 **“R1: it would be good to compare DLE to for example KSD on a**
 16 **complex model. R2: more empirical examples on non-toy datasets...”**

17 We run the same typical/outlier image detection task in Section 6.2 on
 18 Fashion MNIST dataset and compare DLE and KSD (see the figure). Both
 19 methods work well and seem to assign high likelihood to similar images.
 20 However, the tails of the fitted densities seem different, judging from low
 21 likelihood images. We observe similar results on MNIST dataset and will
 22 provide analysis into the differences of tail behaviors in the revision.



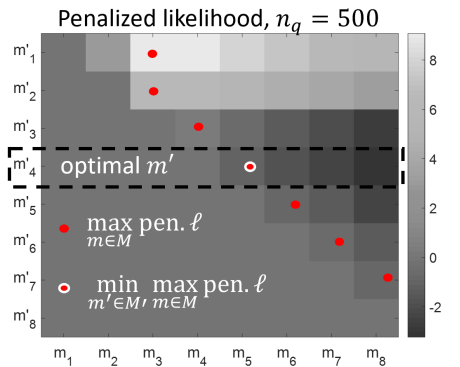
23 **“R2: ... but it didn’t compare to other methods like Contrastive Diver-**
 24 **gence or Noise Contrastive Estimation (NCE).”**

25 Contrastive divergence was commonly used for restricted Boltzman Machine (RBM) where Gibbs sampling can be
 26 efficiently done. However, we consider a much wider family of models. In our MNIST experiment, the model is more
 27 complicated than an RBM and Gibbs sampling is hard. Other MCMC methods such as Metropolis-Hasting are unlikely
 28 to succeed as it is also difficult to design a proposal distribution in a 784-dimensional space. Due to the difficulties of
 29 applying MCMC in high dimensional tasks, we restrict our discussion on sampling-free methods for their computational
 30 efficiency and reliability in those applications. We tried NCE in the MNIST experiment but cannot find a good noise
 31 distribution which would give comparable results to DLE and KSD. Those results will be presented in revision.

32 **R2:“... how a practitioner could select the Stein features to use?” R3:**
 33 **“...some guidelines or heuristics for how to select the feature..”**

34 Section 4.3 provides an information-criterion based model selection method.
 35 Suppose M is a set of different choices of Stein features. Given a parameter
 36 θ , one should select the features $\hat{m}(\theta) := \arg \max_{m \in M} \mathbb{E}_q[\ell(\theta, \hat{\delta}(m))]$ as
 37 this choice would minimize $\text{KL}[q || r_{\hat{\delta} p_{\theta}}]$ (see eq. 2).

38 If we have a set of candidate density models M' , we can jointly select
 39 density model and Stein feature at the same time: $(\hat{m}, \hat{m}') :=$
 40 $\arg \min_{m' \in M'} \max_{m \in M} \mathbb{E}_q[\ell(\hat{\theta}(m'), \hat{\delta}(m))]$, where $(\hat{\theta}(m'), \hat{\delta}(m))$
 41 are estimated parameters under the model choice (m', m) . Replacing
 42 $\mathbb{E}_q[\ell(\theta, \hat{\delta}(m))]$ with the penalized likelihood derived in Section 4.3, we can
 43 get a practical model selection method. We create a numerical experiment and plot the calculated penalized likelihood
 44 using scaled colors with respect to both M and M' (see the figure on the right). It shows that our information criterion
 45 can indeed select the optimal density model. We will explain this procedure in our revision.



46 **R2: “Given these conflicting forces, how does one choose the Stein features?”** Yes, there *can* be a trade-off between
 47 efficiency and overfitting. This happens in classic settings too: MLE is an efficient estimator, but suffers from overfitting
 48 when the dataset is small. Given a small number of samples, we may have to settle for a less efficient estimator to avoid
 49 overfitting. However, the aforementioned information criterion can be used to select Stein features in this setting.

50 **R3: “...may imply that that the estimated density $p(x, \theta)$ is not positive everywhere...”** The unnormalized density
 51 model $\bar{p}(x; \theta)$ in our problem, by definition, should be non-negative everywhere for all $\theta \in \Theta$, thus the estimated
 52 density $\bar{p}(x; \hat{\theta})$ is also non-negative. The estimated *density ratio* is guaranteed to be positive only within X_q , but the
 53 estimated density is guaranteed to be positive everywhere by definition. We will clarify this in our revision.