On the Value of Infinite Gradients in Variational Autoencoder Models

Anonymous Author(s) Affiliation Address email

Abstract

A number of recent studies of continuous variational autoencoder (VAE) models 1 have noted, either directly or indirectly, the tendency of various parameter gradients 2 to drift towards infinity during training. Because such gradients could potentially з 4 contribute to numerical instabilities, and are often framed as a problematic phenomena to be avoided, it may be tempting to shift to alternative energy functions 5 that guarantee bounded gradients. But it remains an open question: What might 6 the unintended consequences of such a restriction be? To address this issue, we 7 examine how unbounded gradients relate to the regularization of a broad class of 8 autoencoder-based architectures, including VAE models. Our main finding is that, 9 if the ultimate goal is to simultaneously avoid over-regularization (high reconstruc-10 tion errors, sometimes referred to as posterior collapse) and under-regularization 11 (excessive latent dimensions are not pruned from the model), then an autoencoder-12 based energy function with infinite gradients around optimal representations is 13 provably required per a certain technical sense we carefully detail. Given that both 14 over- and under-regularization can directly lead to poor generated sample quality 15 or suboptimal feature selection, this result suggests that heuristic modifications to 16 or constraints on the VAE energy function may be ill-advised, and large gradients 17 should be accommodated to the extent possible. 18

19 1 Introduction

Suppose we have access to continuous variables $x \in \chi$ that are drawn from ground-truth measure μ_{gt} . This measure assigns probability mass $\mu_{gt}(dx)$ to the infinitesimal dx residing within $\chi \subset \mathbb{R}^d$ such that we have $\int_{\chi} \mu_{gt}(dx) = 1$. This formalism allows us to consider data that may lie on or near an *r*-dimensional manifold embedded in \mathbb{R}^d (implying r < d), capturing the notion of low-dimensional structure relative to the high-dimensional ambient space.

25 Because of the possibility of an unknown latent manifold, it is common to approximate the corre-26 sponding ground-truth measure via a density model parameterized as

$$p_{\theta}(\boldsymbol{x}) = \int p_{\theta}(\boldsymbol{x}|\boldsymbol{z}) p(\boldsymbol{z}) d\boldsymbol{z}.$$
 (1)

In this expression θ are trainable parameters and $z \in \mathbb{R}^{\kappa}$ serves as a low-dimensional latent representation, with fixed prior $p(z) = \mathcal{N}(z; \mathbf{0}, I)$ and ideally $\kappa \geq r$. If some θ^* were available such that $\int_A p_{\theta^*}(x) dx \approx \int_A \mu_{gt}(dx)$ for any measurable $A \subseteq \chi$, then the model would adequately reflect the intrinsic underlying distribution. Of course we will generally not know in advance the value of θ^* , but in principle we might consider minimizing $-\log p_{\theta}(x)$ averaged across a set of training samples $\{x^{(i)}\}_{i=1}^n$ drawn from μ_{gt} , i.e., minimize $\frac{1}{n}\sum_{i=1}^n -\log \left[p_{\theta}(x^{(i)})\right] \approx \int -\log \left[p_{\theta}(x)\right] \mu_{gt}(dx)$

Submitted to 35th Conference on Neural Information Processing Systems (NeurIPS 2021). Do not distribute.

33 over θ . Unfortunately though, the marginalization required to produce $p_{\theta}(x^{(i)})$ is generally in-

tractable for models of sufficient representational power. To circumvent this issue, the variational autoencoder (VAE) [Kingma and Welling, 2014, Rezende et al., 2014] instead optimizes the tractable

autoencoder (VAE) [Kingma and Welling, 2014, Rezende et al., 2014] instead optimizes the tractable variational bound $\mathcal{L}(\theta, \phi) \triangleq$

$$\frac{1}{n}\sum_{i=1}^{n}\left\{-\mathbb{E}_{q_{\phi}\left(\boldsymbol{z}|\boldsymbol{x}^{(i)}\right)}\left[\log p_{\theta}\left(\boldsymbol{x}^{(i)}|\boldsymbol{z}\right)\right]+\mathbb{KL}\left[q_{\phi}(\boldsymbol{z}|\boldsymbol{x}^{(i)})||p(\boldsymbol{z})\right]\right\}\geq\frac{1}{n}\sum_{i=1}^{n}-\log\left[p_{\theta}\left(\boldsymbol{x}^{(i)}\right)\right].$$
(2)

Here $q_{\phi}(\boldsymbol{z}|\boldsymbol{x})$ represents a variational approximation to $p_{\theta}(\boldsymbol{z}|\boldsymbol{x})$ with additional parameters ϕ gov-37 erning the tightness of the bound. It is commonly referred to as an encoder distribution since it 38 quantifies the mapping from x to the latent code z. For analogous reasons, $p_{\theta}(x|z)$ is labeled 39 as the *decoder* distribution. When combined, the data-dependent factor $-\mathbb{E}_{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \left[\log p_{\theta}(\boldsymbol{x}|\boldsymbol{z})\right]$ 40 can be viewed as instantiating a form of stochastic autoencoder (AE) structure, which attempts to 41 assign high probability to accurate reconstructions of each x; if $q_{\phi}(z|x)$ is Dirac delta function, 42 then a regular deterministic AE emerges with loss dictated by the decoder negative log-likelihood 43 $-\log p_{\theta}(\boldsymbol{x}|\boldsymbol{z})$. Beyond this, $\mathbb{KL}[q_{\phi}(\boldsymbol{z}|\boldsymbol{x})||p(\boldsymbol{z})]$ serves as a regularization factor that pushes the 44 encoder distribution towards the prior. The bound (2) can be minimized over $\{\theta, \phi\}$ using SGD and a 45 simple reparameterization trick [Kingma and Welling, 2014, Rezende et al., 2014]. 46

The latter requires that we assume specific functional forms for the encoder and decoder distributions. In this regard, it is common to select $q_{\phi}(\boldsymbol{z}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{z}|\boldsymbol{\mu}_{z}, \operatorname{diag}[\boldsymbol{\sigma}_{z}]^{2})$, where the Gaussian moment vectors $\boldsymbol{\mu}_{z}$ and $\boldsymbol{\sigma}_{z}$ are functions of model parameters ϕ and the random variable \boldsymbol{x} , i.e., $\boldsymbol{\mu}_{z} \equiv \boldsymbol{\mu}_{z}(\boldsymbol{x};\phi)$, and $\boldsymbol{\sigma}_{z} \equiv \boldsymbol{\sigma}_{z}(\boldsymbol{x};\phi)$. Similarly, for continuous data the decoder model is conventionally parameterized as $p_{\theta}(\boldsymbol{x}|\boldsymbol{z}) = \mathcal{N}(\boldsymbol{x};\boldsymbol{\mu}_{x},\gamma \boldsymbol{I})$, with mean defined analogously as $\boldsymbol{\mu}_{x} \equiv \boldsymbol{\mu}_{x}(\boldsymbol{z};\theta)$ and scalar variance parameter $\gamma > 0$. The functions $\boldsymbol{\mu}_{z}(\boldsymbol{x};\phi)$, $\boldsymbol{\sigma}_{z}(\boldsymbol{x};\phi)$, and $\boldsymbol{\mu}_{x}(\boldsymbol{z};\theta)$ are all instantiated using deep neural network layers. Given this definitions, (2) can be expressed in the more transparent form

$$\mathcal{L}(\theta,\phi) \equiv \frac{1}{n} \sum_{i=1}^{n} \left\{ \mathbb{E}_{q_{\phi}}(\boldsymbol{z}|\boldsymbol{x}^{(i)}) \left[\frac{1}{\gamma} \| \boldsymbol{x}^{(i)} - \boldsymbol{\mu}_{x}(\boldsymbol{z};\theta) \|_{2}^{2} \right] + d \log \gamma + \left\| \boldsymbol{\sigma}_{z}\left(\boldsymbol{x}^{(i)};\phi \right) \right\|_{2}^{2} - \log \left| \operatorname{diag} \left[\boldsymbol{\sigma}_{z}\left(\boldsymbol{x}^{(i)};\phi \right) \right]^{2} \right| + \left\| \boldsymbol{\mu}_{z}\left(\boldsymbol{x}^{(i)};\phi \right) \right\|_{2}^{2} \right\}.$$
(3)

Although VAE models have been successfully applied to a variety of practical problems [Li and She, 55 2017, Schott et al., 2018, Walker et al., 2016], at times they exhibit potentially problematic behavior 56 that is not fully understood. For example, a number of recent works have mentioned that if a trainable 57 decoder variance parameter γ is included within a Gaussian VAE as in (3), then the optimal value 58 may converge to zero, resulting in infinite or unbounded gradients and potential instabilities [Dai 59 and Wipf, 2019, Mattei and Frellsen, 2018, Rezende and Viola, 2018, Takahashi et al., 2018]. While 60 unbounded gradients may indeed be troublesome from an optimization perspective, in this work we 61 will reframe such gradients as an integral part of any successful autoencoder-based energy function 62 designed to model continuous data arising from a low-dimensional manifold. 63

To accomplish this, our analysis is split into three parts. First, in Section 2 we detail how unbounded 64 65 gradients contribute to an optimal, balanced form of regularization, allowing the VAE to capture low-dimensional manifold structure via a maximally parsimonious latent representation. Such 66 representations turn out to be critical for tasks such as generating non-blurry samples that resemble 67 the training data [Tolstikhin et al., 2018], or for using autoencoder-based models in general to robustly 68 screen outliers [An and Cho, 2015, Xu et al., 2018]. Of course it is natural to consider whether these 69 same goals could not be achieved using an alternative energy function with strictly bounded gradients. 70 The second and primary component of our contribution answers this question in the negative. More 71

⁷² concretely, our main result from Section 3 proves that canonical autoencoder-based architectures ⁷³ will necessarily require unbounded gradients to guarantee the type of maximally parsimonious latent ⁷⁴ representation mentioned above. Thirdly, in Section 4 we elucidate the benefits of learning γ during ⁷⁵ training, even in situations where we know that the optimal value will be at or near zero and contribute ⁷⁶ to arbitrarily-large gradients. In particular, we argue that (at the very least) learning γ localizes

- troublesome unbounded gradients to narrow regions around minima of (3), while simultaneously
- ⁷⁸ smoothing the VAE objective across optimization trajectories prior to convergence.

Overall, our contribution can be viewed as complementary to the wide body of work analyzing what 79 is commonly-referred to as *posterior collapse* in VAE models [He et al., 2019, Razavi et al., 2019]. 80 The latter can be related to the situation where γ is too large (either implicitly [Dai et al., 2020] or 81 explicitly [Lucas et al., 2019]) and along all or most latent dimensions the posterior $q_{\phi}(z_i|\boldsymbol{x}^{(i)})$ 82 collapses to the prior $\mathcal{N}(0,1)$ leading to high reconstruction errors. In contrast, we direct our attention 83 herein to the *opposite* condition whereby γ is arbitrarily small and unbounded gradients invariably 84 ensue. In this regime, the resulting latent representations obtained from bad local minimizers can 85 potentially be under-regularized in an underappreciated sense that will be described in subsequent 86 87 sections.

88 2 Optimal Low-Dimensional Structure via Unbounded VAE Gradients

As alluded to previously, the VAE objective will experience unbounded gradients if $\gamma \rightarrow 0$ as has sometimes been observed (at least approximately) during training. But perhaps counter-intuitively, this phenomena nonetheless serves a critical purpose in the context of modeling data with lowdimensional manifold structure as described in Section 1. To quantify this assertion, we first precisely define what type of low-dimensional or sparse latent representations will be considered optimal for our present analysis; later we link this definition to practical VAE/AE applications.

95 2.1 Optimal Sparse Representations

96 **Definition 1** An autoencoder-based architecture (VAE or otherwise) produces an **optimal sparse** 97 *representation* of a training set X if the following two conditions simultaneously hold:

(i) The reconstruction error is zero.¹ For a stochastic VAE model this requirement entails that

$$\frac{1}{n}\sum_{i=1}^{n}\mathbb{E}_{q_{\phi}\left(\boldsymbol{z}|\boldsymbol{x}^{(i)}\right)}\left[\left\|\boldsymbol{x}^{(i)}-\boldsymbol{\mu}_{x}\left[\boldsymbol{z};\boldsymbol{\theta}\right]\right\|_{2}^{2}\right]=0.$$
(4)

In contrast, for an AE with encoder function $\boldsymbol{\mu}_{z}(\boldsymbol{x};\phi)$ and decoder $\boldsymbol{\mu}_{x}(\boldsymbol{z};\theta)$, we analogously require that the now deterministic reconstruction satisfies $\frac{1}{n}\sum_{i=1}^{n} \|\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_{x}[\boldsymbol{\mu}_{z}(\boldsymbol{x}^{(i)};\phi);\theta]\|_{2}^{2} = 0.$

(ii) Conditioned on achieving perfect reconstructions per criteria (i) above, the number of latent dimensions of z containing no information about X is maximal. More specifically, for the VAE we say that the j-th latent dimension contains no information regarding X if q_φ (z_j|x⁽ⁱ⁾) = N(0,1) for all i, i.e., the posterior is pushed to the prior along this dimension. Likewise, for an AE with encoder μ_z (x; φ), the corresponding requirement can be relaxed to μ_z (x⁽ⁱ⁾; φ)_j = 0 for all i. In either case, a latent dimension so-defined provides no benefit in reducing the reconstruction error and could in principle be removed from the model.²

Conceptually, this definition is merely describing the most parsimonious latent representation of the training data that nonetheless allows us to obtain perfect reconstructions. And when combined with the low-dimensional manifold assumption from Section 1, it readily follows that an optimal sparse representation of X will generally involve $\kappa - r$ uninformative dimensions (assuming $\kappa \ge r$). As a simple illustrative example, for data generated by a low-dimensional linear subspace model, PCA can be trivially applied to obtain the corresponding optimal sparse representation, in this case defined by the smallest subspace containing all of the data variance.

In broader contexts involving nonlinear low-dimensional manifolds, the VAE can achieve something analogous when granted sufficient encoder/decoder capacity, at least assuming that the global optimum

¹If insufficient capacity or other modeling errors are a factor, we can of course relax this definition to allow for reconstruction errors within some tolerance level.

²It could also be argued that if $q_{\phi}(z_j | \boldsymbol{x}^{(i)})$ is set to another arbitrary distribution not equal to the prior $p(\boldsymbol{z})$, it would similarly contain no information about \boldsymbol{X} . However, this possibility is most because if $q_{\phi}(z_j | \boldsymbol{x}^{(i)})$ is not useful for predicting \boldsymbol{X} , then it will be set to the prior $p(\boldsymbol{z})$ to optimize the KL-divergence term, not some arbitrary distribution. Analogous reasoning holds for the deterministic AE case, where typical sparsity penalties (e.g., [Fan and Li, 2001, Rao et al., 2003]) are used to push noninformative dimensions to zero.



Figure 1: The importance of optimal sparse representations in screening outliers. In this example, the simple 2D principal subspace obtainable by PCA can perfectly reconstruct the inlier manifold shown in red. But this requires using two separate informative dimensions, allowing both inliers *and* outliers to be reconstructed with zero error within this subspace. In contrast, it is only by recovering the curved 1D inlier manifold, which relies on a single informative dimension, that inliers and outliers can be differentiated. Please see supplementary for practical example using real data.

of (3) can be found [Dai and Wipf, 2019]. This capability requires that the VAE avoid both over-118 or under-regularization of the latent representations. To be more precise, VAE over-regularization 119 (sometimes loosely referred to as latent posterior collapse [He et al., 2019, Razavi et al., 2019]) 120 occurs when too many latent dimensions are non-informative (i.e., the latent posterior along these 121 dimensions is close to the non-informative prior) such that the reconstruction error is high and criteria 122 (i) is violated. In contrast, with *under-regularized* solutions criteria (i) may be satisfied, and yet 123 in reducing the reconstruction error towards zero, an excessive number of latent dimensions are 124 informative in violation of criteria (ii). 125

In avoiding both of these suboptimal scenarios, it can be shown that the VAE explicitly relies on $\gamma \rightarrow 0$ and the attendant unbounded gradients that follow [Dai and Wipf, 2019]. From an intuitive standpoint, we might expect that achieving criteria (*i*) would require an unbounded gradient given that, if we minimize (3) over γ in isolation, the optimal value satisfies

$$\gamma^* = \frac{1}{dn} \sum_{i=1}^n \mathbb{E}_{q_{\phi}\left(\boldsymbol{z} | \boldsymbol{x}^{(i)}\right)} \left[\left\| \boldsymbol{x}^{(i)} - \boldsymbol{\mu}_x \left[\boldsymbol{z}; \boldsymbol{\theta} \right] \right\|_2^2 \right].$$
(5)

If we then plug this value back into the $d \log \gamma$ term from (3), the result is unbounded from below as the reconstruction error goes to zero. Of course to actually achieve near-zero reconstruction errors, at least some dimensions of σ_z must be pushed towards zero, which can also lead to infinite gradients within the KL-divergence term. See [Dai and Wipf, 2019] for more details.

134 2.2 Relevance to Typical VAE Usage Regimes

Obtaining minimalist latent representations as distilled by Definition 1 can serve a variety of practical 135 downstream applications, such as feature extraction [Bengio et al., 2013, Ng, 2011], compression 136 [Ballé et al., 2018, Donoho, 2006, Minnen et al., 2018], manifold learning [Silva et al., 2006], 137 corruption removal [Dai et al., 2018], or even the generation of realistic samples. With respect to the 138 latter, it has been shown in [Dai and Wipf, 2019] that what we have above defined as an optimal sparse 139 representation can be viewed as a necessary (albeit not sufficient) condition for generating samples 140 using a continuous-space VAE that match the training distribution. In this context, the unneeded 141 latent dimensions are simply set to the uninformative Gaussian prior to optimize the KL regularizer; 142 however, this white noise can be filtered out by the decoder so as not to impact the reconstructions 143 allowing both criteria (i) and (ii) of Definition 1 to be satisfied. In principle, a deterministic AE 144 architecture capable of producing optimal sparse representations can also be leveraged to generate 145 realistic samples; this would simply involve first discarding the uninformative dimensions and then 146 applying the same analysis from [Dai and Wipf, 2019]. In fact, variants of this strategy have been 147 previously considered in [Ghosh et al., 2019, Tolstikhin et al., 2018]. 148

And as a final motivational example, any AE-based architecture capable of producing optimal sparse representations can naturally be applied to screening outliers by squeezing the latent space to the minimal number of informative dimensions needed for reconstructing inliers. In doing so, we reduce the risk that outlier points $x^{(out)}$ can be accurately reconstructed by exploiting the superfluous latent

flexibility. Here we are assuming that $x^{(out)} \sim \mu_{out} \neq \mu_{gt}$ for some outlier distribution μ_{out} . Figure 153 1 contains an illustration of the basic rationale.³ 154

155 Additionally, in the supplementary we demonstrate that indeed, if the inlier data (in this case Fashion 156 MNIST samples) come from a low-dimensional manifold, outlier points (MNIST samples) can be reliably differentiated, provided that $\kappa \geq r$ and the VAE has sufficient capacity and the learned γ can 157 converge to near zero. And because of the VAE's propensity to find optimal sparse representations 158 where possible, even as κ is raised such that $\kappa \gg r$, unneeded dimensions are shut off to reduce the 159 risk of outliers masquerading as inliers (see supplementary). 160

Can we Reliably Obtain Optimal Sparse Representations without 161 3 **Unbounded Gradients?** 162

As discussed in Section 2, given data originating from a low-dimensional manifold, optimal sparse 163 representations are a necessary requirement (at least approximately) for various tasks such as gener-164 ating non-blurry samples aligned with the ground-truth distribution or alternatively, screening for 165 outliers. We have also discussed how the divergent gradients associated with $\gamma \rightarrow 0$, allow VAE 166 global minima to achieve such optimal sparse representations. But what about alternatives that 167 circumvent such unbounded gradients altogether? For example, could we not consider a regularized 168 AE model that, while encouraging sparse latent representations [Ng, 2011], explicitly relies on energy 169 function terms with bounded gradients? Despite this conceptual possibility, per the analysis that 170 follows, the answer turns out to be unequivocally no. Or more specifically, if we wish to guarantee 171 an optimal sparse representation, then even arbitrary AE-based objectives will necessarily require 172 penalty terms with infinite gradients around optimal solutions. 173

3.1 A Generic AE-based Objective for Optimal Sparse Representations 174

Consider the constrained objective function $\mathcal{L}_h(\theta, \phi) \triangleq$ 175

$$h\left(\frac{1}{dn}\sum_{i=1}^{n}\left\|\boldsymbol{x}^{(i)}-\boldsymbol{\mu}_{x}\left(\boldsymbol{z}^{(i)};\theta\right)\right\|_{2}^{2}\right)+\frac{1}{d}\sum_{k=1}^{\kappa}h\left(\frac{1}{n}\left\|\boldsymbol{z}_{k}\right\|_{2}^{2}\right),$$

s.t. $\boldsymbol{z}^{(i)}=\boldsymbol{\mu}_{z}\left(\boldsymbol{x}^{(i)};\phi\right)$ $\forall i, \theta \in \Theta,$ (6)

176

where $Z \triangleq \{z^{(i)}\}_{i=1}^n \in \mathbb{R}^{\kappa \times n}$ and z_k denotes the k-th row of Z. This expression can be viewed as 177 characterizing a typical regularized AE with a generic penalty function $h: \mathbb{R}^+ \to \mathbb{R}$ on the norm 178 across training samples of each latent dimension. The multipliers 1/n, 1/d, and 1/(dn) ensure a 179 form of proportional regularization within energy functions composed of multiple penalty factors of 180 varying dimension designed to favor sparsity [Wipf and Wu, 2012]. The square-root Lasso can be 181 viewed as a special case of this strategy that emerges when h is a square-root function [Belloni et al., 182 2011]. We adopt this formalism to avoid distracting complications from tunable trade-off parameters; 183 however, our central conclusions still hold even when such a parameter is introduced. And finally, 184 the constraint $\theta \in \Theta$ is included in (6) to prevent the trivial solution $Z \to 0$, which could occur if 185 each $z^{(i)}$ is pushed to zero while the decoder μ_x includes an unconstrained compensatory factor that 186 grows towards infinity such that the error $\| \boldsymbol{x}^{(i)} - \boldsymbol{\mu}_x (\boldsymbol{z}^{(i)}; \theta) \|_2$ can still be minimized to zero. Any 187 regularized AE must include such constraints to avoid trivial solutions, or else additional penalty 188 terms on θ that serve a similar purpose. 189

We can also relate (6) to various VAE instantiations as follows: 190

Lemma 2 Let $\mu_x(z; \theta) = Wz + b$ for some $W \in \mathbb{R}^{d \times \kappa}$ and $b \in \mathbb{R}^d$, and $\sigma_z(x; \phi) = s$ for any arbitrary $s \in \mathbb{R}^{\kappa}$. Then in the limit $\gamma \to 0$, the VAE loss from (3) is such that $\min_{\sigma_z(x;\phi)} \mathcal{L}(\theta, \phi) \equiv s$ 191 192 $\min_{\mathbf{s}} \mathcal{L}(\theta, \phi)$ reduces to (6) with $h(\cdot) = \log(\cdot)$, excluding irrelevant constant factors.

193

³The only exception to this line of reasoning would be adversarial outliers that follow the exact same low-dimensional structure as the inliers, meaning μ_{out} and μ_{qt} both apply all of their probability mass to the same low-dimensional manifold. In this scenario, we would need to exploit differences between μ_{out} and μ_{gt} within the manifold to reliably screen outliers, a regime in which Definition 1 is not directly applicable. That being said, differentiating μ_{out} and μ_{at} once a shared low-dimensional manifold has been modeled is far easier than doing so in the original ambient space.

Lemma 3 For any arbitrary $\mu_x(z;\theta)$ and $\theta \in \Theta$, if we enforce $\sigma_z(x;\phi) \to 0$ for all x and apply a log transformation to each $||z_k||_2^2$, then the VAE loss from (3) collapses to (6) with $h(\cdot) = \log(\cdot)$, excluding irrelevant constant factors.

¹⁹⁷ Collectively, these results point to a close affiliation between (6), with *h* set to a log function, and the ¹⁹⁸ VAE loss, especially given that $\gamma \to 0$ and $\sigma_z(\mathbf{x}; \phi) \to \mathbf{0}$ along many dimensions are characteristics ¹⁹⁹ of VAE global optima [Dai and Wipf, 2019]. Hence it is natural to consider more general selections ²⁰⁰ of *h* in the context of optimal sparse representations.

201 3.2 Main Result: Unbounded Gradients Cannot be Avoided

Given a generic AE architecture as in (6), it is natural to examine what possible functions h are such that any global minimum of $\mathcal{L}_h(\theta, \phi)$ is guaranteed to produce an optimal sparse representation. This can be addressed as follows:

Theorem 4 Assume the constraint $\theta \in \Theta$ and data $\mathbf{X} = {\mathbf{x}^{(i)}}_{i=1}^n \in \mathbb{R}^{d \times n}$ are such that to achieve $\mathbf{x}^{(i)} = \boldsymbol{\mu}_{\mathbf{x}} \left(\mathbf{z}^{(i)}; \theta \right) \forall i (i.e., perfect reconstruction) requires that <math>\|\mathbf{z}_k\|_2 > 0$ for at least r < d rows of \mathbf{Z} . Then to guarantee (without further assumptions on \mathbf{X}) that minimization of $\mathcal{L}_h(\theta, \phi)$ achieves zero reconstruction error using at most r nonzero rows of \mathbf{Z} (i.e., active dimensions), h must have an unbounded gradient around zero.

Note that a similar result can be obtained by replacing the reconstruction penalty with the additional constraint $\sum_{i=1}^{n} || \mathbf{x}^{(i)} - \boldsymbol{\mu}_{x} (\mathbf{z}^{(i)}; \theta) ||_{2}^{2} = 0$, in which case no trade-off parameter, fixed or otherwise, need be included. We also emphasize that Theorem 4 effectively implies that, to guarantee every global minima corresponds with an optimal sparse reconstruction per our definition, the constituent penalty functions must have an unbounded gradient around zero. This can be viewed as a necessary, albeit not sufficient condition, for optimal sparsity, as sufficiency requires additional care taking limits around zero, e.g., $\gamma \rightarrow 0$ in the case of the VAE.

Consequently, we cannot simply replace a VAE model with any possible AE architecture to somehow guarantee optimal sparse representations devoid of infinite surrounding gradients. Rather, optimal sparse representations and infinite gradients go hand-in-hand unless further restrictive assumptions are imposed on the training data.

221 3.3 High-Level Intuition Behind Theorem 4

While the proof is predicated on a nuanced counterexample designed with a specific technical 222 purpose in mind (see supplementary file), we can nonetheless loosely convey the basic idea through 223 a toy illustration shown in Figure 2. Here we are assuming that the data points $\{x^{(i)}\}_{i=1}^{n}$ lie on a 224 1D manifold embedded in 2D ambient space. Moreover, we stipulate that this manifold is tightly 225 squeezed within a small non-negative $\epsilon \times \epsilon$ square near zero, represented by the blue curve on the 226 lefthand side of Figure 2. Now consider a sample point $x' = [x'_1, x'_2]^{\top}$ taken from somewhere along 227 the stated 1D manifold. We represent this point using two candidate decoder functions, both assumed 228 to be within the capacity of μ_x , as displayed in the middle of Figure 2. 229

For the simple decoder case, which is just the identity function $\mu_x(z; \theta) = z$, the values of $z_1 = z'_1$ and $z_2 = z'_2$ needed for a perfect reconstruction will both be small, i.e., $\{z'_1, z'_2\} \le \epsilon$ by design. In contrast, the optimal decoder only requires that a single dimension of z, namely z_1 , be nonzero. However, the optimal value actually needed for perfect reconstruction, denoted z^*_1 , can be arbitrarily large in controlling where along the extended, labyrinthine manifold pathway x' is located (for ease of presentation we will assume z^*_1 is also positive). Hence we can easily have that

$$z_1^* \gg \epsilon \ge \max\left(z_1', z_2'\right). \tag{7}$$

Because of this, to ensure that $z^* = [z_1^*, 0]^\top$ is preferred over the z' alternative, we require a concave penalty function h on each encoder dimension such that any infinitesimal movement away from zero incurs an arbitrarily-large cost, while increases originating from points away from zero incur only a modest additional cost (see the green curve on the righthand side of Figure 2). From this it follows that any movement of z'_1 and z'_2 away from zero, no matter how small, will be such that we can guarantee that the penalties on z^* and z' will satisfy

$$h(z_1^*) + h(0) = h(z_1^*) \approx h(z_1') \approx h(z_2') < h(z_1') + h(z_2') \approx 2[h(z_1^*) + h(0)], \quad (8)$$



Figure 2: 2D illustration of the intuition behind Theorem 4. See Section 3.3 for details.

and so z_* is preferred. The righthand side of Figure 2 motivates this relationship. Note also that if we were to explicitly bound the slope of *h* around zero, then we could always select an ϵ sufficiently small such that the inequality in (8) is reversed; hence an unbounded slope is required to achieve the stated result.

To a large extent, the intuition here mirrors the basic scenario from Figure 1, and is emblematic of 246 broader situations that naturally arise in practice. For example, if we run PCA on MNIST data, we 247 find that only a 100 or so principal components are needed to achieve highly accurate reconstructions. 248 But a VAE model with only around 15 informative latent dimensions can accomplish something 249 similar [Dai et al., 2018] by closely approximating an optimal sparse representation using a nonlinear 250 decoder. Of course unless we have a objective function with a strong preference for lower-dimensional 251 structures, as instantiated through large gradients around optimal sparse representations, then the 252 network may well favor or converge to a simpler, higher-dimensional alternative (e.g., resembling a 253 PCA solution). 254

²⁵⁵ 4 Mitigating Unbounded Gradients via γ -Dependent Smoothing

While we have argued that unbounded gradients may serve a useful purpose in obtaining optimal 256 latent representations, they may nonetheless pose challenges from an optimization standpoint. In 257 addressing this concern, it is worth acknowledging that energy functions involving infinite gradients 258 and/or unbounded regions are already indispensable across a wide range of structured regression and 259 sparse estimation problems [Gorodnitsky and Rao, 1997]. This history implies that when training a 260 261 VAE or other related AE structure, we may borrow appropriate tools designed to mitigate the risk of converging to bad local solutions or regions of instability. In this vein, one effective strategy involves 262 partially minimizing what amounts to a smoothed version of the original objective function. The 263 degree of smoothness is then gradually reduced as the optimization trajectory moves towards an 264 optimum. Within the domain of underdetermined linear inverse problems, this procedure is frequently 265 used to find maximally sparse representations with minimal reconstruction error [Chartrand and Yin, 266 2008, Hu et al., 2012, Xu et al., 2013]. 267

The VAE automatically accomplishes something similar when we choose to iteratively estimate γ dur-268 ing training rather than merely setting its value to near zero as may be theoretically optimal (assuming 269 we know that there exists sufficient network capacity to achieve negligible reconstruction errors). 270 Initially, when the reconstruction cost is still high because encoder/decoder parameters have not 271 converged, the learned γ will be larger and the overall VAE energy will be relatively smooth, devoid of 272 many deep local minimizers. It is only later as the data fit $\sum_{i=1}^{n} \mathbb{E}_{q_{\phi}(\boldsymbol{z}|\boldsymbol{x}^{(i)})} \left[\left\| \boldsymbol{x}^{(i)} - \boldsymbol{\mu}_{x}(\boldsymbol{z};\theta) \right\|_{2}^{2} \right]$ 273 becomes small that γ will follow suite, and by this point it is more likely that we have already 274 approached a basin of attraction capable of producing optimal sparse reconstructions. Additionally, 275 unlike fixing $\gamma \approx 0$ for all training iterations, in which case gradients will be unbounded right from 276 the beginning, by learning γ we will likely only encounter large gradients in a narrow neighborhood 277 around minimizing solutions. This implies that in practice, we only need accommodate such gra-278 dients when the reconstruction error becomes small, at which point stability countermeasures can 279 be deployed if/when necessary, e.g., reduced step size, checks for oscillating gradient sign patterns 280 [Riedmiller and Braun, 1993], etc. 281

To help visualize these points, in Figure 3 we have plotted 1D slices through the objective function of a simple VAE model involving a single layer for both encoder and decoder, applied to data from a random low-dimensional subspace. We vary $\gamma \in \{10^{-3}, 10^{-2}, 10^{-1}, 1\}$, which exposes



Figure 3: Plots (a) and (b) show two sets of representative 1D slices through the VAE objective function (3) as the value of γ is varied. Dashed vertical lines indicate the *x*-axis location of the minimal value of each respective slice and γ setting. And for both plots (a) and (b) the 1D slices are set such that an optimal sparse representation would occur at zero on the *x*-axis when $\gamma \rightarrow 0$. It can be observed that disconnected local minima only occur when γ is small.

the increasing gradients and multi-modal nature of the objective function as γ becomes smaller. 285 Dashed vertical lines indicate the minimal value of the respective curve for each γ . Additionally, we 286 have explicitly designed this visualization such that there will exist an optimal sparse representation 287 at zero on the x-axis. Consequently, we can readily observe that as γ becomes sufficiently small, 288 the minimizing value of the VAE energy increasingly aligns with an optimal sparse representation 289 as desired. However, as γ is reduced the energy is less smooth and disconnected local minima 290 appear in both 1D slices. And local minima of the VAE loss surface can at times be risk points for 291 under-regularized representations. 292

To further explore the implications of this γ -dependent smoothing effect, we empirically compare a practical scenario whereby learning γ may be better than fixing it to an arbitrarily small value. To this effect, we first train a VAE model on CelebA data [Liu et al., 2015] and learn an appropriate small value of γ denoted γ^* (note that γ^* need not be exactly zero since with real data and limited capacity the network will generally display some nonzero reconstruction errors). Please see the supplementary for network and training details. We then retrain the same network from scratch but with $\gamma = \gamma^*$ fixed throughout all training iterations.

The resulting models are evaluated via the reconstruction error and the maximum mean discrepancy (MMD) between the aggregated posterior $q_{\phi}(z) \triangleq \frac{1}{n} \sum_{i} q_{\phi}(z|x^{(i)})$ [Makhzani et al., 2016] and the prior $p(z) = \mathcal{N}(z; 0, I)$. If too few latent dimensions are removed by swamping the appropriate channels with noise following the prior (i.e., under-regularization), then we would expect $q_{\phi}(z)$ to be confined to a low-dimensional manifold in \mathbb{R}^{κ} and the MMD to be much larger. Note that for ideal generative modeling performance via an autoencoder architecture, it is required that

$$\frac{1}{n}\sum_{i}q_{\phi}(\boldsymbol{z}|\boldsymbol{x}^{(i)})\approx\int_{\boldsymbol{\chi}}q_{\phi}(\boldsymbol{z}|\boldsymbol{x})\mu_{gt}(d\boldsymbol{x})=p(\boldsymbol{z}),$$
(9)

meaning the MMD from $\mathcal{N}(\boldsymbol{z}; \boldsymbol{0}, \boldsymbol{I})$ is ideally zero [Makhzani et al., 2016]. With manifold data this is only possible if an optimal sparse representation is produced by the VAE or autoencoder-based analogue [Tolstikhin et al., 2018].

Results are displayed in Figure 4(a), where as expected the reconstruction errors are nearly identical, but the learnable γ case leads to much lower MMD values, indicative of a better local solution with reduced under-regularization. We also plot the evolution of the gradient magnitudes $\left\|\frac{d\mathcal{L}(\theta,\phi)}{dz}\right\|_2$ in Figure 4(b) (other gradients are similar). When γ is learned, the gradient increases slowly; however, with fixed $\gamma = \gamma^*$, there exists a large gradient right from the start since γ^* is small but the reconstruction error is high. This contributes to a worse final solution per the results in Figure 4(a).



Figure 4: (a) Reconstruction error and MMD between $q_{\phi}(z)$ and $\mathcal{N}(0, I)$ on CelebA (128×128 resolution). We first train a VAE with learnable γ and obtain the optimal value γ^* . Then we fix $\gamma = \gamma^*$ and re-train the same network from scratch. While the final reconstruction errors are almost the same, the MMDs between $q_{\phi}(z)$ and the prior $\mathcal{N}(0, I)$ are significantly different. (b) The Evolution of the gradient $\left\| \frac{d\mathcal{L}(\theta,\phi)}{dz} \right\|_2$. Although both curves end up with similar final values, the large initial gradient with fixed γ is disruptive to the final solution.

Additionally, examples of using a learnable γ to improve generated sample quality based on these principles can be found in [Dai and Wipf, 2019].

317 5 Conclusion

It is not uncommon to learn the VAE decoder variance parameter in situations where the training data has a noise component that we are unable or do not wish to model. By doing so we can avoid tuning a trade-off parameter while allowing the model to adapt to the data. However, with sufficient capacity networks and relatively clean data, the risk of unbounded gradients when training γ has frequently been raised as a potentially problematic phenomena. We nonetheless provide formal justification for this choice (even in cases where γ does tend to zero) on two primary fronts:

• We prove that unbounded gradients are in fact necessary for guaranteeing that global minima 324 of canonical AE architectures will coincide with optimal spare representations, meaning high 325 fidelity reconstruction of the training data using the minimal number of informative latent 326 dimensions. Hence there is no obvious alternative if this form of parsimony is our goal. 327 Furthermore, given the value of such representations to numerous downstream tasks as described 328 in Section 2.2, our analysis suggests that heuristic modifications to or constraints on the VAE 329 energy function may be ill-advised, and large gradients should be accommodated to the extent 330 possible (e.g., reduced step size, checks for oscillating gradient sign patterns, etc.). 331

• We present compelling evidence that by learning γ , large gradients away from global minimizers, 332 as well as at least some bad local minimizers, can be mitigated or smoothed within the VAE loss 333 surface. This helps to explain observed successes learning γ in situations where the optimal 334 value turns out to be small or near zero [Dai and Wipf, 2019]. Note that as mentioned in 335 Section 1, it is already known that fixing γ too high can lead to over-regularization and the 336 widely-studied phenomena of posterior collapse [He et al., 2019, Lucas et al., 2019, Razavi et al., 337 2019]. In a similar vein, we have demonstrated the complementary yet underappreciated fact 338 that prematurely fixing γ too *low*, even to what may ultimately be the optimal value near zero, 339 can steer convergence towards under-regularized local minima and the inadvertent wasteful 340 deployment of latent degrees-of-freedom. 341

And finally, although not our focus, our results herein naturally relate to more flexible VAE models with non-Gaussian latent posteriors [Kingma et al., 2016, Rezende and Mohamed, 2015] or adaptable/trainable priors [Bauer and Mnih, 2019, Tomczak and Welling, 2018]. While these types of enhancements can be useful tools for favoring $q_{\phi}(z) \approx p(z)$, they do not circumvent the infinite gradients that will occur around optimal sparse representations. Additionally, for a brief discussion regarding the implications to β -VAE models [Higgins et al., 2017]; please see the supplementary.

348 **References**

- Jinwon An and Sungzoon Cho. Variational autoencoder based anomaly detection using reconstruction probability. *Special Lecture on IE*, 2:1–18, 2015.
- Johannes Ballé, David Minnen, Saurabh Singh, Sung Jin Hwang, and Nick Johnston. Variational image compression with a scale hyperprior. *arXiv preprint arXiv:1802.01436*, 2018.
- ³⁵³ Matthias Bauer and Andriy Mnih. Resampled priors for variational autoencoders. 2019.
- Alexandre Belloni, Victor Chernozhukov, and Lie Wang. Square-root lasso: Pivotal recovery of sparse signals via conic programming. *Biometrika*, 98(4):791–806, 2011.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new
 perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828,
 2013.
- Rick Chartrand and Wotao Yin. Iteratively reweighted algorithms for compressive sensing. *Interna- tional Conference on Accoustics, Speech, and Signal Processing*, 2008.
- Bin Dai and David Wipf. Diagnosing and enhancing VAE models. *International Conference on Learning Representations*, 2019.
- Bin Dai, Yu Wang, John Aston, Gang Hua, and David Wipf. Connections with robust PCA and
 the role of emergent sparsity in variational autoencoder models. *Journal of Machine Learning Research*, 2018.
- Bin Dai, Ziyu Wang, and David Wipf. The usual suspects? Reassessing blame for VAE posterior collapse. In *International Confernece on Machine Learning*, 2020.
- D.L. Donoho. Compressed sensing. *IEEE Trans. Information Theory*, 52(4), 2006.
- Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *J. American Statistical Association*, 96(456):1348–1360, 2001.
- Partha Ghosh, Mehdi SM Sajjadi, Antonio Vergari, Michael Black, and Bernhard Schölkopf. From variational to deterministic autoencoders. *arXiv preprint arXiv:1903.12436*, 2019.
- Irina Gorodnitsky and Bhaskar Rao. Sparse signal reconstruction from limited data using FOCUSS:
 A re-weighted minimum norm algorithm. *IEEE Transactions on signal processing*, 45(3):600–616,
 1997.
- Junxian He, Daniel Spokoyny, Graham Neubig, and Taylor Berg-Kirkpatrick. Lagging inference
 networks and posterior collapse in variational autoencoders. In *International Conference on Learning Representations*, 2019.
- Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, , and Alexander Lerchner. β -vae: Learning basic visual concepts with a constrained variational framework. In *International Conference on Learning Representations*, 2017.
- Yue Hu, Sajan Goud Lingala, and Mathews Jacob. A fast majorize–minimize algorithm for the
 recovery of sparse and low-rank matrices. *IEEE Transactions on Image Processing*, 21(2):742–753,
 2012.
- Diederik Kingma and Max Welling. Auto-encoding variational Bayes. In *International Conference* on *Learning Representations*, 2014.
- Durk Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved variational inference with inverse autoregressive flow. In *Advances in Neural Information Processing Systems 29*, pages 4743–4751. 2016.
- Xiaopeng Li and James She. Collaborative variational autoencoder for recommender systems. In
 International Conference on Knowledge Discovery and Data Mining, pages 305–314, 2017.

- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In
 IEEE International Conference on Computer Vision, pages 3730–3738, 2015.
- James Lucas, George Tucker, Roger B Grosse, and Mohammad Norouzi. Understanding posterior col-
- lapse in generative latent variable models. International Conference on Learning Representations,
 Workshop Paper, 2019.
- Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2016.
- Pierre-Alexandre Mattei and Jes Frellsen. Leveraging the exact likelihood of deep latent variables
 models. *arXiv preprint arXiv:1802.04826*, 2018.
- David Minnen, Johannes Ballé, and George D Toderici. Joint autoregressive and hierarchical priors
 for learned image compression. In *Advances in Neural Information Processing Systems 31*, pages
 10771–10780. 2018.
- 405 Andrew Ng. Sparse autoencoder. CS294A Lecture notes, 72(2011):1–19, 2011.
- Bhaskar Rao, Kjersti Engan, Shane Cotter, Jason Palmer, and Kenneth Kreutz-Delgado. Subset
 selection in noise based on diversity measure minimization. *IEEE Trans. Signal Processing*, 51(3):
 760–770, March 2003.
- Ali Razavi, Aäron van den Oord, Ben Poole, and Oriol Vinyals. Preventing posterior collapse with
 δ-VAEs. In *International Conference on Learning Representations*, 2019.
- Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. *arXiv preprint arXiv:1505.05770*, 2015.
- 413 Danilo Jimenez Rezende and Fabio Viola. Taming VAEs. arXiv preprint arXiv:1810.00597, 2018.
- ⁴¹⁴ Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and
 ⁴¹⁵ approximate inference in deep generative models. In *International Conference on Machine* ⁴¹⁶ *Learning*, 2014.
- Martin Riedmiller and Heinrich Braun. A direct adaptive method for faster backpropagation learning:
 The rprop algorithm. In *IEEE international conference on neural networks*, pages 586–591, 1993.
- Lukas Schott, Jonas Rauber, Matthias Bethge, and Wieland Brendel. Towards the first adversarially
 robust neural network model on MNIST. In *International Conference on Learning Representations*,
 2018.
- Jorge Silva, Jorge Marques, and João Lemos. Selecting landmark points for sparse manifold learning.
 Advances in Neural Information Processing Systems 18, pages 1241–1248, 2006.
- Hiroshi Takahashi, Tomoharu Iwata, Yuki Yamanaka, Masanori Yamada, and Satoshi Yagi. Student t variational autoencoder for robust density estimation. In *International Joint Conference on Artificial Intelligence*, pages 2696–2702, 2018.
- Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schoelkopf. Wasserstein auto-encoders.
 International Conference on Learning Representations, 2018.
- Jakub Tomczak and Max Welling. VAE with a VampPrior. In *International Conference on Artificial Intelligence and Statistics*, pages 1214–1223, 2018.
- Jacob Walker, Carl Doersch, Abhinav Gupta, and Martial Hebert. An uncertain future: Forecasting
 from static images using variational autoencoders. In *European Conference on Computer Vision*,
 pages 835–851, 2016.
- David Wipf and Yi Wu. Dual-space analysis of the sparse linear model. In *Advances in Neural Information Processing Systems*, pages 1745–1753, 2012.
- Haowen Xu, Wenxiao Chen, Nengwen Zhao, Zeyan Li, Jiahao Bu, Zhihan Li, Ying Liu, Youjian
 Zhao, Dan Pei, Yang Feng, et al. Unsupervised anomaly detection via variational auto-encoder for
 seasonal KPIs in web applications. In *International World Wide Web Conference*, pages 187–196,
 2018.

Li Xu, Shicheng Zheng, and Jiaya Jia. Unnatural ℓ_0 sparse representation for natural image deblurring. 440 In IEEE Conference on Computer Vision and Pattern Recognition, pages 1107–1114, 2013. 441

Checklist 442

443

468

469

471

472

473

474

475

476

477

479

480

481

482

1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's 444 contributions and scope? [Yes] 445 (b) Did you describe the limitations of your work? [Yes] We describe that our VAE/AE results 446 assume that the global minima of the objective can actually be found. Consequently, if 447 the model converges to a local solution, some of our conclusions regarding optimal sparse 448 representations will not necessarily apply. These issues are discussed in multiple places 449 in Sections 2-5. Overall, the circumscribed scope of our work is clearly articulated. 450 (c) Did you discuss any potential negative societal impacts of your work? [N/A] Our work is 451 devoted to the better understanding of an existing, well-established generative modeling 452 framework; we do not see any clear pathway to negative societal impact. 453 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? 454 [Yes] 455 2. If you are including theoretical results... 456 (a) Did you state the full set of assumptions of all theoretical results? [Yes] 457 (b) Did you include complete proofs of all theoretical results? [Yes] 458 3. If you ran experiments... 459 460 results (either in the supplemental material or as a URL)? [No] As our work is primarily 461 theoretical/analytical in nature, the experiments are secondary and involve standard 462 models/data. Even so, if accepted we can provide relevant code. 463 464 465 466 467
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] Training details are deferred to the supplementary.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [N/A]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [N/A]
- 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets... 470
 - (a) If your work uses existing assets, did you cite the creators? [N/A]
 - (b) Did you mention the license of the assets? [N/A]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
- 5. If you used crowdsourcing or conducted research with human subjects... 478
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent 483 on participant compensation? [N/A] 484