
Fully Variational Noise-Contrastive Estimation

Christopher Zach

Chalmers University of Technology
Gothenburg, Sweden
zach@chalmers.se

Abstract

By using the underlying theory of proper scoring rules, we design a family of noise-contrastive estimation (NCE) methods that are tractable for latent variable models. Both terms in the underlying NCE loss, the one using data samples and the one using noise samples, can be lower-bounded as in variational Bayes, therefore we call this family of losses *fully variational noise-contrastive estimation*. Variational autoencoders are a particular example in this family and therefore can be also understood as separating real data from synthetic samples using an appropriate classification loss. We further discuss other instances in this family of fully variational NCE objectives and indicate differences in their empirical behavior.

1 Introduction

Estimating the parameters of a model distribution from a training set is an important research topic with applications in deep generative models (e.g. [8, 18, 24, 15, 28, 5]), out-of-distribution (OOD) or anomaly detection [32, 23, 16, 17] and representation learning [4, 22, 2, 19]. Maximum-likelihood estimation is the method of choice when the parametric model distribution is normalized and can be evaluated efficiently (which is the case for “elementary” probability distributions and for normalizing flows [24]). The expressiveness of a model distribution can be enhanced by introducing latent variables and by using an unnormalized distribution (also known as energy-based model). Both of these modifications prevent the maximum likelihood method from being applicable: latent variables often lead to intractable integrals or sums when computing the marginal likelihood, and likewise the normalization factor (also called the partition function) of an unnormalized model is typically intractable.

Latent variables are usually addressed by utilizing the evidence lower bound (ELBO) of the likelihood as in variational Bayes (e.g. [12]), and parameters of unnormalized models can be estimated from data by methods such as score matching [11] or noise-contrastive estimation (NCE, [9, 10]). NCE can intuitively be understood as learning a binary classifier separating training data from samples drawn from a fully known noise distribution. Variational NCE [26] aims to enable the estimation of unnormalized latent variable models from data by leveraging the ELBO. It succeeds only partially, since the ELBO cannot be applied on all terms in the NCE objective, and an intractable marginal remains. In this work we derive modified instances of NCE that allow the application of the ELBO on all terms, and the resulting objective is therefore free from intractable sums (or integrals). We call the resulting method *fully variational noise-contrastive estimation*. Interestingly, variational autoencoders [14, 25] are one particular (and important) instance in this family of fully variational NCE methods.

2 Background

Proper scoring rules Let $\mathcal{P} \subseteq \mathbb{R}^d$, and let $G: \mathcal{P} \rightarrow \mathbb{R}$ be a differentiable convex mapping. The Bregman divergence between $p \in \mathcal{P}$ and $q \in \mathcal{P}$ is defined as

$$D_G(p||q) \stackrel{\text{def}}{=} G(p) - (G(q) + (p - q)^\top \nabla G(q)), \quad (1)$$

i.e. $D_G(p||q)$ is the error between $G(p)$ and the linearization (first-order Taylor expansion) of G at q . Convexity of G implies that $D_G(p||q)$ is non-negative. If G is strictly convex, then $D_G(p||q) = 0$ iff $p = q$.

Now let p and q be the parameters of a categorical distribution, i.e. $P(X = k|p) = p_k$ and $P(X = k|q) = q_k$ for a categorical random variable X with values in $\{1, \dots, d\}$. The domain \mathcal{P} is therefore the probability simplex, $\mathcal{P} = \{p \in [0, 1]^d : \sum_{k=1}^d p_k = 1\}$. In this setting $D_G(p||q)$ can be stated as

$$\begin{aligned} D_G(p||q) &= G(p) - G(q) + \mathbb{E}_{X \sim p} \left[\frac{d}{dq_X} G(q) \right] - \mathbb{E}_{X \sim q} \left[\frac{d}{dq_X} G(q) \right] \\ &= G(p) + \mathbb{E}_{X \sim p} \left[\frac{d}{dq_X} G(q) - G(q) \right] - \mathbb{E}_{X' \sim q} \left[\frac{d}{dq_{X'}} G(q) \right]. \end{aligned} \quad (2)$$

Minimizing $D_G(p||q)$ w.r.t. q for fixed p is equivalent to

$$\begin{aligned} \arg \min_{q \in \mathcal{P}} D_G(p||q) &= \arg \min_{q \in \mathcal{P}} -G(q) - \sum_k (p_k - q_k) \frac{\partial}{\partial q_k} G(q) \\ &= \arg \max_{q \in \mathcal{P}} \mathbb{E}_{X \sim p} \left[\frac{\partial}{\partial q_X} G(q) + G(q) - \sum_k q_k \frac{\partial}{\partial q_k} G(q) \right] \\ &= \arg \max_{q \in \mathcal{P}} \mathbb{E}_{X \sim p} [S(X, q)], \end{aligned} \quad (3)$$

where we defined the *proper scoring rule* (PSR) S as follows,

$$S(x, q) \stackrel{\text{def}}{=} \frac{\partial}{\partial q_x} G(q) + G(q) - \sum_k q_k \frac{\partial}{\partial q_k} G(q). \quad (4)$$

Note that maximization w.r.t. q only requires samples from p , but does not need the knowledge of the distribution p itself. Therefore proper scoring rules are one method to estimate distribution parameters when only samples from an unknown data distribution p are available.

If G is strictly convex, then the resulting PSR is a *strictly* PSR. If e.g. G is chosen as the negated Shannon entropy, then $S(x, q) = \log q_x$ is called the *logarithmic* scoring rule underlying maximum likelihood estimation and the cross-entropy loss in machine learning. It is an instance of a *local* PSR [20], which does not depend on any value of $q_{x'}$ for $x' \neq x$ (the score matching cost [11] being another example). We refer to [7] and [3] for an extensive overview and further examples of proper scoring rules.

PSRs for binary RVs When X is a binary random variable, and therefore $x \in \{0, 1\}$, then we only need one parameter $\mu \in [0, 1]$ to characterize the corresponding Bernoulli distribution. For a differentiable convex function $G: [0, 1] \rightarrow \mathbb{R}$ the induced Bregman divergence between $\mu \in [0, 1]$ and $\nu \in [0, 1]$ is given by

$$D_G(\mu||\nu) = G(\mu) - G(\nu) - (\mu - \nu)G'(\nu) \quad (5)$$

and

$$\begin{aligned} \arg \min_{\nu \in [0, 1]} D_G(\mu||\nu) &= \arg \max_{\nu \in [0, 1]} G(\nu) + (\mu - \nu)G'(\nu) \\ &= \arg \max_{\nu \in [0, 1]} \mathbb{E}_{x \sim \text{Ber}(\mu)} [G(\nu) + (x - \nu)G'(\nu)]. \end{aligned} \quad (6)$$

The resulting PSR S is therefore

$$S(1, \nu) = G(\nu) + (1 - \nu)G'(\nu) \quad S(0, 1 - \nu) = G(\nu) - \nu G'(\nu). \quad (7)$$

G can be recovered via

$$G(\nu) = \nu S(1, \nu) + (1 - \nu)S(0, 1 - \nu) = \mathbb{E}_{x \sim \text{Ber}(\mu)} [S(x, x\nu + (1 - x)(1 - \nu))]. \quad (8)$$

Noise-contrastive estimation Noise-contrastive estimation (NCE, [9, 10]) ultimately casts the estimation of parameters of an unknown data distribution as a binary classification problem. Let $\Omega \subseteq \mathbb{R}^n$ and X be a n -dimensional random vector. Let p_d the (unknown) data distribution, p_θ a model distribution (with parameters θ) and p_n a user-specified noise distribution. Let Z be a (fair) Bernoulli RV that determines whether a sample is drawn from the data (respectively model) distribution or from the noise distribution p_d .¹ NCE applies the logarithmic PSR to match the posteriors,

$$P_{d,n}(Z = 1|X = x) = \frac{p_d(x)}{p_d(x) + p_n(x)} \quad P_{\theta,n}(Z = 1|X = x) = \frac{p_\theta(x)}{p_\theta(x) + p_n(x)}, \quad (9)$$

which yields the NCE objective

$$J_{\text{NCE}}(\theta) = \mathbb{E}_{X \sim p_d} \left[\log \frac{p_\theta(X)}{p_\theta(X) + p_n(X)} \right] + \mathbb{E}_{X \sim p_n} \left[\log \frac{p_n(X)}{p_\theta(X) + p_n(X)} \right]. \quad (10)$$

After introducing $r_\theta(x) \stackrel{\text{def}}{=} p_\theta(x)/p_n(x)$ this reads as

$$J_{\text{NCE}}(\theta) = \mathbb{E}_{X \sim p_d} [-\log(1 + r_\theta(X)^{-1})] + \mathbb{E}_{X \sim p_n} [-\log(1 + r_\theta(X))], \quad (11)$$

establishing the connection to logistic regression. At first glance this is superficially similar to GANs [8], but it lacks e.g. the problematic min-max structure of GANs. In contrast to e.g. maximum likelihood estimation, NCE is applicable even when the model distribution is unnormalized, i.e.

$$p_\theta(x) = \frac{1}{Z(\theta)} p_\theta^0(x) \quad (12)$$

for an unnormalized model $p_\theta^0(x)$ and an intractable partition function $Z(\theta) = \sum_x p_\theta^0(x)$.² NCE allows to estimate the value of the partition function $Z(\theta)$ for the obtained model parameters θ by augmenting the parameter vector to (θ, Z) and use the relation $p_\theta(x) = p_\theta^0(x)/Z$. Extensions to the basic NCE framework are discussed in [21] and [1].

NCE is not directly applicable to latent variable models, where the joint density $p_\theta(X, Z)$ is specified, but the induced marginal $p_\theta(X)$ is only indirectly given via

$$p_\theta(x) = \sum_z p_\theta(x, z) = \sum_z p_\theta(x|z) p_Z(z), \quad (13)$$

where we use a generative model for the joint $p_\theta(X, Z)$.

Using latent variable models greatly enhances the expressiveness of model distributions, but exact computation of the marginal $p_\theta(x)$ is often intractable. By noting that the term under the first expectation in Eq. 11 is concave w.r.t. $r_\theta(x)$, Variational NCE [26] proposes to apply the evidence lower bound (ELBO) to obtain a tractable variational lower bound for the first term in Eq. 11. Unfortunately, the second term in Eq. 11 is convex in r_θ and the ELBO does not apply here. Importance sampling is leveraged instead to estimate the intractable expectation inside the second term. In the following section we show how the ELBO can be applied on both terms in a slightly generalized version of NCE.

3 Fully Variational NCE

First, we generalize the NCE objective (Eq. 10) to arbitrary strictly proper scoring rules for binary random variables,

$$J_{S\text{-NCE}}(\theta) = \mathbb{E}_{X \sim p_d} \left[S \left(1, \frac{r_\theta(x)}{1 + r_\theta(x)} \right) \right] + \mathbb{E}_{X \sim p_n} \left[S \left(0, \frac{1}{1 + r_\theta(x)} \right) \right], \quad (14)$$

where r_θ is the density ratio, $r_\theta(x) \stackrel{\text{def}}{=} p_\theta(x)/p_n(x)$. $J_{S\text{-NCE}}$ is maximized w.r.t. the parameters θ in this formulation. Recall that $r_\theta(x)/(1 + r_\theta(x))$ is the posterior of x being a sample drawn from the model p_θ , and $1/(1 + r_\theta(x))$ is the posterior for x being a noise sample. Our aim is to determine a convex function G such that both mappings

$$f_1(r) = S(1, r/(1 + r)) \quad \text{and} \quad f_0(r) = S(0, 1/(1 + r)) \quad (15)$$

¹We omit the possibility of using general Bernoulli RV for notational simplicity.

²For brevity we use sums to refer to marginalization of RV, but these sums should always be understood as the appropriate Lebesgue integrals.

are concave. If this is the case, then

$$\begin{aligned} f_k(r_\theta(x)) &= f_k\left(\frac{p_\theta(x)}{p_n(x)}\right) = f_k\left(\frac{\sum_z p_\theta(x,z)}{p_n(x)}\right) = f_k\left(\frac{\sum_z p_\theta(x,z)q_k(z|x)}{p_n(x)q_k(z|x)}\right) \\ &\geq \sum_z q_k(z|x) f_k\left(\frac{p_\theta(x,z)}{p_n(x)q_k(z|x)}\right) = \mathbb{E}_{z \sim q_k(Z|x)} \left[f_k\left(\frac{p_\theta(x,z)}{p_n(x)q_k(z|x)}\right) \right] \end{aligned}$$

for $k \in \{0, 1\}$. $q_k(Z|X)$ is a posterior corresponding to the encoder part. Overall, $J_{S\text{-NCE}}$ in Eq. 14 can be lower bounded as follows,

$$J_{S\text{-NCE}}(\theta) = \mathbb{E}_{x \sim p_d} [f_1(r_\theta(x))] + \mathbb{E}_{x \sim p_n} [f_0(r_\theta(x))] \geq \max_{q_1, q_0} J_{S\text{-fvNCE}}(\theta, q_1, q_0) \quad (16)$$

with the r.h.s. defined as the *fully variational* NCE loss,

$$J_{S\text{-fvNCE}}(\theta, q_1, q_0) \stackrel{\text{def}}{=} \mathbb{E}_{x \sim p_d, z \sim q_1(Z|x)} \left[f_1\left(\frac{p_\theta(x,z)}{p_n(x)q_1(z|x)}\right) \right] + \mathbb{E}_{x \sim p_n, z \sim q_0(Z|x)} \left[f_0\left(\frac{p_\theta(x,z)}{p_n(x)q_0(z|x)}\right) \right]. \quad (17)$$

Note that we allow in principle two separate encoders, q_1 and q_0 , since the ELBO is applied at two places independently. For brevity we introduce the following short-hand notations for the joint distributions,

$$p_{d,k}(x, z) \stackrel{\text{def}}{=} p_d(x)q_k(z|x) \quad p_{n,k}(x, z) \stackrel{\text{def}}{=} p_n(x)q_k(z|x), \quad (18)$$

resulting in a more compact expression for $J_{S\text{-fvNCE}}$,

$$J_{S\text{-fvNCE}}(\theta, q_1, q_0) = \mathbb{E}_{(x,z) \sim p_{d,1}} \left[f_1\left(\frac{p_\theta(x,z)}{p_{n,1}(x,z)}\right) \right] + \mathbb{E}_{(x,z) \sim p_{n,0}} \left[f_0\left(\frac{p_\theta(x,z)}{p_{n,0}(x,z)}\right) \right]. \quad (19)$$

From $p_\theta(x)p_\theta(z|x) = p_\theta(x,z)$ we deduce that the lower bound is tight, i.e. $J_{S\text{-NCE}}(\theta) = \max_{q_1, q_0} J_{S\text{-fvNCE}}(\theta, q_1, q_0)$ when the encoders q_1 and q_0 are equal to the model posterior, $q_1(Z|X) = q_0(Z|X) = p_\theta(Z|X)$ a.e. $J_{S\text{-fvNCE}}$ in Eq. 17 is formulated as a population loss, but the corresponding empirical risk can be immediately obtained by sampling from p_d , p_n and the encoder distributions.

Now the question is whether such concave mappings f_1 and f_0 satisfying Eq. 15 for a PSR S exist. Since common PSRs such as the logarithmic and the quadratic PSR violate these properties, existence of such a PSR is not obvious. The next section discusses how to construct such PSRs and provides examples.

4 A Family of Suitable Proper Scoring Rules

In this section we construct a pair (f_1, f_0) of concave mappings, such that the induced functions $S(1, \cdot)$ and $S(0, \cdot)$ in Eq. 15 form a PSR. The following result provides sufficient conditions on such a pair (f_1, f_0) :

Lemma 1. *Let a pair of functions (f_0, f_1) , $f_k : (0, \infty) \rightarrow \mathbb{R}$, satisfy the following:*

1. Both f_1 and f_0 are concave,
2. f_1 and f_0 satisfy the compatibility condition

$$f_0'(r) = -r f_1'(r) \quad (20)$$

for all $r > 0$,

3. the mapping $G(\mu) = \mu f_1(\mu/(1-\mu)) + (1-\mu) f_0(\mu/(1-\mu))$ is convex in $(0, 1)$.

Then S is a PSR. Such pairs (f_1, f_0) are said to have to double ELBO property.

Proof. We abbreviate $S_1(\mu) := S(1, \mu)$ and $S_0(1-\mu) := S(0, 1-\mu)$ and recall the relations between S and G :

$$\begin{aligned} G(\mu) &= \mu S_1(\mu) + (1-\mu) S_0(1-\mu) \\ S_0(1-\mu) &= G(\mu) - \mu G'(\mu) \\ S_1(\mu) &= G(\mu) + (1-\mu) G'(\mu) = S_0(1-\mu) + G'(\mu) \end{aligned} \quad (21)$$

and therefore $G'(\mu) = S_1(\mu) - S_0(1-\mu)$. We calculate

$$G'(\mu) = S_1(\mu) - S_0(1-\mu) + \mu S_1'(\mu) - (1-\mu)S_0'(1-\mu) \quad (22)$$

Combining these relations implies that

$$\mu S_1'(\mu) - (1-\mu)S_0'(1-\mu) = 0 \iff S_0'(1-\mu) = \frac{\mu}{1-\mu} \cdot S_1'(\mu) \quad (23)$$

Now the relation between μ and r is $\mu = r/(1+r)$ and therefore $r = \mu/(1-\mu)$, which we use to express (f_1, f_0) in terms of (S_1, S_0) ,

$$f_1(r) = S_1(\mu) = S_1(r/(1+r)) \quad f_0(r) = S_0(1-\mu) = S_0(1/(1+r)). \quad (24)$$

Using $d\mu/dr = (1+r)^{-2}$ and

$$f_1'(r) = \frac{1}{(1+r)^2} S_1'(\mu) \quad f_0'(r) = -\frac{1}{(1+r)^2} S_0'(1-\mu),$$

the condition can be restated as

$$-(1+r)^2 f_0'(r) = r \cdot (1+r)^2 f_1'(r) \iff f_0'(r) = -r f_1'(r), \quad (25)$$

which is the second requirement on (f_1, f_0) . Now if (f_1, f_0) satisfy Eq. 20, then (S_1, S_0) satisfy the relations of a binary PSR in Eq. 21 for an induced function G . If G is now convex, then (S_1, S_0) is a PSR. \square

One consequence of the condition in Eq. 20 is, that f_1 is increasing and f_0 is decreasing or vice versa. This further implies that S cannot be symmetric, i.e.

$$S(1, \mu) \neq S(0, 1-\mu), \quad (26)$$

and positive and negative samples are penalized differently in the overall loss. This is in contrast to many well-known PSR, which are symmetric (such as the logarithmic PSR used in NCE). The condition also implies that

$$f_0''(r) = -f_1'(r) - r f_1''(r) \stackrel{!}{\leq} 0.$$

Since f_1 is concave and $r \geq 0$, $-r f_1''(r) \geq 0$. This has to be compensated by f_1' increasing sufficiently fast with r . Since $f_1'(r) \geq -r f_1''(r) \geq 0$, f_1 is increasing and f_0 is decreasing in $\mathbb{R}_{\geq 0}$. This observation yields some intuition on $J_{S\text{-fVNCE}}$ in Eq. 17: the first term aims to align p_θ with $p_{d,1}$ by maximizing $p_\theta(x, z)/p_{n,1}(x, z)$ for real data (and its code), whereas the second term favors mis-alignment between p_θ and $p_{n,0}$ for noise samples (by minimizing the likelihood ratio $p_\theta(x, z)/p_{n,0}(x, z)$).

Eq. 20 immediately allows to establish one pair (f_1, f_0) satisfying the double ELBO property: we choose $f_1(r) = \log r$, which yields $f_1'(r) = 1/r$ and therefore $f_0(r) = -r$. Both f_1 and f_0 are concave. Further,

$$S_1(\mu) = \log \frac{\mu}{1-\mu} \quad S_0(1-\mu) = -\frac{\mu}{1-\mu} \quad (27)$$

and therefore

$$G(\mu) = \mu S_1(\mu) + (1-\mu) S_0(1-\mu) = \mu \left(\log \frac{\mu}{1-\mu} - 1 \right), \quad (28)$$

which is convex in $(0, 1)$. Thus, we have established the existence of one PSR allowing the ELBO being applied on both terms as in Eq. 16. This example can be generalized to the following parametrized family of PSRs:

Lemma 2. *A family of PSRs satisfying the double ELBO property is given by*

$$f_1(r) = \log(r + \beta) \quad f_0(r) = \beta \log(r + \beta) - r \quad (29)$$

for any $\beta \geq 0$.

Proof. This follows from

$$f_0'(r) = -r f_1'(r) = -\frac{r}{r+\beta} = -\frac{r+\beta-\beta}{r+\beta} = -1 + \frac{\beta}{r+\beta} \implies f_0(r) = \beta \log(r + \beta) - r.$$

Further, G'' can be calculated as

$$G''(\mu) = -\frac{1}{(1-\mu)^2(\beta\mu - \mu - \beta)} = \frac{1}{(1-\mu)^2(\mu + \beta(1-\mu))} > 0, \quad (30)$$

which establishes the convexity of G (due to $(1-\mu)^2 > 0$ and $\mu + \beta(1-\mu) > 0$ for $\mu \in (0, 1)$ and $\beta \geq 0$). \square

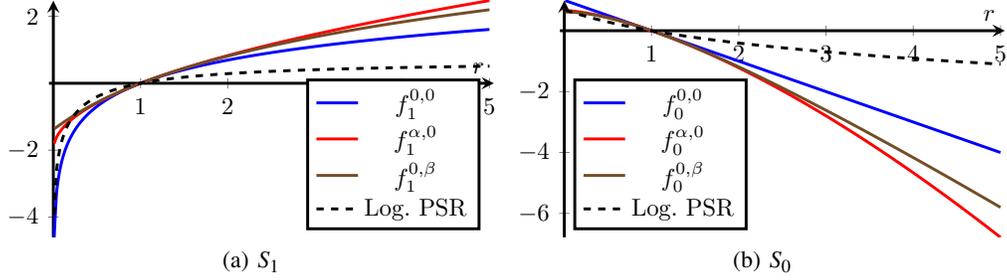


Figure 1: Several pairs (f_1, f_0) , in particular $(f_1^{0,0}, f_0^{0,0})$, $(f_1^{\alpha,0}, f_0^{\alpha,0})$ and $(f_1^{0,\beta}, f_0^{0,\beta})$ for $\alpha = 1/2$ and $\beta = 1$ (solid curves). Both f_1 and f_0 are concave functions. The pair (f_1, f_0) induced by the logarithmic PSR is shown for reference (dashed curve, which is concave in (a), but convex in (b)).

A 2-parameter family of PSRs is given next.

Lemma 3. For $\alpha \in (0, 1]$ and $\beta \geq 0$ we choose

$$f_1(r) = \frac{1}{\alpha}(r + \beta)^\alpha \quad f_0(r) = -\frac{1}{\alpha+1}(r + \beta)^{\alpha+1}.$$

This pair induces a strictly PSR satisfying the double ELBO property.

Proof. Both f_1 and f_0 are clearly concave. We deduce

$$f_1'(r) = (r + \beta)^{\alpha-1} \quad f_0'(r) = -r(r + \beta)^{\alpha-1} = -rf_1'(r), \quad (31)$$

hence (f_1, f_0) satisfy the condition in Eq. 20. $G''(\mu)$ can be calculated as

$$G''(\mu) = \left(\frac{\mu + \beta(1 - \mu)}{1 - \mu} \right)^\alpha \cdot \frac{\mu + \beta(2 - \alpha)(1 - \mu)}{(1 - \mu)^2(\mu + \beta(1 - \mu))^2}. \quad (32)$$

The first factor is positive for $\alpha \in (0, 1]$, $\beta \geq 0$ and $\mu \in (0, 1)$. Analogously, the second factor is positive since the numerator is positive for the allowed values of (μ, α, β) , and the denominator is a product of squares. \square

Since

$$\lim_{\alpha \rightarrow 0^+} f_0'(r; \alpha, \beta) = -1 \implies \lim_{\alpha \rightarrow 0^+} f_1'(r; \alpha, \beta) = (r + \beta)^{-1}, \quad (33)$$

we deduce that the limit $\alpha \rightarrow 0^+$ yields the pair (f_1, f_0) from Lemma 2 (up to constants independent of r).

For visualization purposes it is convenient to normalize f_1 and f_0 such that $f_1(1) = f_0(1) = 0$ and $f_1'(1) = 1$ (and therefore $f_0'(1) = -1$). With such normalization the above pairs are given by

$$\begin{aligned} f_1(r; \alpha, \beta) &= \frac{(1+\beta)^{1-\alpha}}{\alpha} ((r + \beta)^\alpha - (1 + \beta)^\alpha) \\ f_0(r; \alpha, \beta) &= -\frac{(1+\beta)^{1-\alpha}}{\alpha(\alpha+1)} ((\alpha r - \beta)(r + \beta)^\alpha - (\alpha - \beta)(1 + \beta)^\alpha). \end{aligned} \quad (34)$$

Few instances of $(f_1^{\alpha,\beta}, f_0^{\alpha,\beta})$ are depicted in Fig. 1. We further introduce the fully variational NCE loss parametrized by (α, β) ,

$$J_{\text{fVNCE}}^{\alpha,\beta}(\theta, q_1, q_0) \stackrel{\text{def}}{=} \mathbb{E}_{(x,z) \sim p_{d,1}} \left[f_1 \left(\frac{p_\theta(x,z)}{p_{n,1}(x,z)}; \alpha, \beta \right) \right] + \mathbb{E}_{(x,z) \sim p_{n,0}} \left[f_0 \left(\frac{p_\theta(x,z)}{p_{n,0}(x,z)}; \alpha, \beta \right) \right]. \quad (35)$$

We would like to get a better understanding of these PSRs in terms of losses used for binary classification. Recall that

$$r = \frac{p}{q} \quad \mu = \frac{p}{p+q} = \frac{1}{1+1/r} = \frac{r}{1+r} = \sigma(\Delta) \quad r = \frac{\mu}{1-\mu} = \frac{\sigma(\Delta)}{1-\sigma(\Delta)} = \frac{\sigma(\Delta)}{\sigma(-\Delta)} = \exp(\Delta).$$

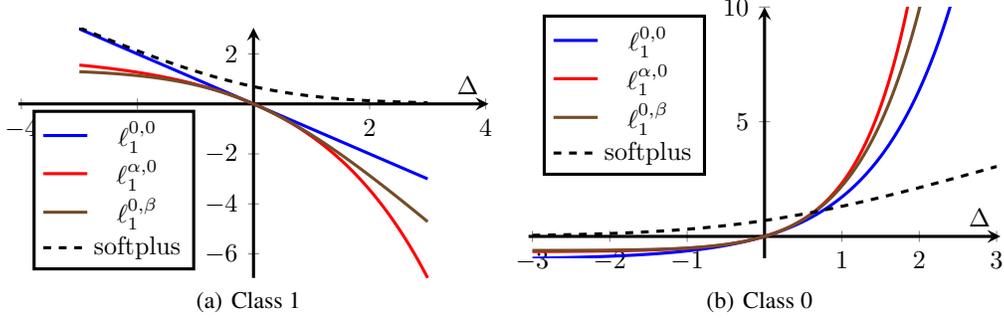


Figure 2: The PSRs from Fig. 1 reinterpreted as binary classification losses in terms of log-ratios $\Delta = \log r$. The soft-plus loss corresponds to the logarithmic PSR.

Here Δ is the logit of the binary classifier. We minimize a classification loss, hence we consider the negated PSRs. Thus, we obtain for the logarithmic PSR,

$$\begin{aligned} -\log(\mu) &= -\log(\sigma(\Delta)) = \log(1 + \exp(-\Delta)) = \text{soft-plus}(-\Delta) \\ -\log(1 - \mu) &= -\log(1 - \sigma(\Delta)) = -\log(\sigma(-\Delta)) = \text{soft-plus}(\Delta), \end{aligned}$$

where $\text{soft-plus}(u) \stackrel{\text{def}}{=} \log(1 + e^u)$. Inserting $f_1(r) = \log(r + \beta)$ and $f_0(r) = \beta \log(r + \beta) - r$ yields

$$\begin{aligned} -f_1(r) &= -\log(r + \beta) = -\log(e^\Delta + \beta) \doteq -\text{soft-max}(\Delta, \log \beta) = \text{soft-min}(-\Delta, -\log \beta) \\ -f_0(r) &= r - \beta \log(r + \beta) = e^\Delta + \beta \text{soft-min}(-\Delta, -\log \beta) \end{aligned}$$

Finally, $f_1(r) = r^\alpha / \alpha$, $f_0(r) = -r^{\alpha+1} / (\alpha + 1)$ results in

$$-f_1(r) = -\frac{1}{\alpha} r^\alpha = -\frac{1}{\alpha} e^{\alpha\Delta} \quad -f_0(r) = \frac{1}{\alpha+1} r^{\alpha+1} = \frac{1}{\alpha+1} e^{(\alpha+1)\Delta}.$$

Graphically, the difference between the logistic classification loss and the double-ELBO losses is, that the logistic loss solely penalizes incorrect predictions and the double ELBO losses strongly favor true positives instead (as shown in Fig. 2).

We conclude this section by noting that non-negative linear combinations of double ELBO pairs have the double ELBO property as well:

Corollary 1. *The set of pairs with the double ELBO property is a convex cone.*

This follows from the linearity of the relations Eq. 20 and Eq. 8.

5 Instances of Fully Variational NCE

In this section we discuss several instances of $J_{\text{fvNCE}}^{\alpha,\beta}$ for specific choices of α and β . For easier identification of known frameworks we focus on normalized model distributions p_θ , but the extension to unnormalized models is straightforward.

5.1 Variational auto-encoders: $(\alpha, \beta) = (0, 0)$

We choose $(\alpha, \beta) = (0, 0)$ in the 2-parameter family given in Lemma 3, i.e. $f_1(r) = \log r$ and $f_0(r) = -r$. The resulting fully variational NCE objective therefore is given by

$$J_{\text{fvNCE}}^{0,0}(\theta, q_1, q_0) = \mathbb{E}_{(x,z) \sim p_{d,1}} \left[\log \left(\frac{p_\theta(x,z)}{p_{n,1}(x,z)} \right) \right] - \mathbb{E}_{(x,z) \sim p_{n,0}} \left[\frac{p_\theta(x,z)}{p_{n,0}(x,z)} \right]. \quad (36)$$

We first focus on the second term:

$$\mathbb{E}_{(x,z) \sim p_{n,0}} \left[\frac{p_\theta(x,z)}{p_{n,0}(x,z)} \right] = \sum_{x,z: p_{n,0}(x,z) > 0} p_\theta(x,z) \leq 1. \quad (37)$$

Now if $\text{supp}(p_\theta) \subseteq \text{supp}(p_{n,0})$, then the r.h.s. of Eq. 37 is exactly 1, otherwise it is bounded by 1 from above.³ We assume that $\text{supp}(p_\theta) \subseteq \text{supp}(p_{n,0})$, then the last term in Eq. 36 is 1, and since $\mathbb{E}_{x \sim p_d} [\log p_n(x)]$ is constant, we obtain

$$J_{\text{fvNCE}}^{0,0}(\theta, q_1, q_0) \doteq \mathbb{E}_{x \sim p_d, z \sim q_1(Z|x)} \left[\log \left(\frac{p_\theta(x, z)}{q_1(z|x)} \right) \right]. \quad (38)$$

After factorizing $p_\theta(x, z) = p_\theta(x|z)p_Z(z)$ this can be identified as the variational autoencoder loss (up to constants independent of θ and q_1),

$$J_{\text{fvNCE}}^{0,0}(\theta, q_1) \doteq \underbrace{\mathbb{E}_{x \sim p_d} \left[\mathbb{E}_{z \sim q_1(Z|x)} [\log p_\theta(x|z)] - D_{KL}(q_1(Z|x} \| p_Z) \right]}_{\stackrel{\text{def}}{=} J_{\text{VAE}}(\theta, q_1)}. \quad (39)$$

Thus, in this setting standard VAE training can be understood as variance-reduced implementation of $J_{\text{fvNCE}}^{0,0}$ (since the stochastic second term becomes a closed-form constant). If $\text{supp}(p_\theta) \not\subseteq \text{supp}(p_{n,0})$, then

$$-\mathbb{E}_{(x,z) \sim p_{n,0}} \left[\frac{p_\theta(x, z)}{p_{n,0}(x, z)} \right] \geq -1 \quad (40)$$

and optimizing the VAE loss J_{VAE} is maximizing a lower bound of $J_{\text{fvNCE}}^{0,0}$. Now let $q_0(Z|X)$ be a deterministic encoder, i.e. $q_0(z|x) = \mathbf{1}[z = g_0(x)]$. In this setting

$$\begin{aligned} J_{\text{fvNCE}}^{0,0}(\theta, q_1, q_0) &\doteq J_{\text{VAE}}(\theta, q_1) - \mathbb{E}_{x \sim p_n} \left[\frac{p_\theta(x, g_0(x))}{p_n(x)} \right] \\ &= J_{\text{VAE}}(\theta, q_1) - \sum_x p_\theta(x, g_0(x)). \end{aligned} \quad (41)$$

Intuitively, $J_{\text{fvNCE}}^{0,0}$ aims to autoencode real data well, but at the same time prefers poor reconstructions for arbitrary inputs. $J_{\text{fvNCE}}^{0,0}$ uses importance weighting to estimate $\sum_x p_\theta(x, g_0(x))$. This term only becomes relevant in the objective if the two encoders q_1 and q_0 are tied in some way (otherwise g_0 may map the input to a constant code that is unlikely to be sampled from q_1).

It is interesting to note that deterministic (and tied) encoders yield somewhat different objectives when comparing classical autoencoders, VAEs and the fully variational NCE:

$$J_{\text{AE}}(\theta, g) = \mathbb{E}_{x \sim p_d} [\log p_\theta(x|g(x))] \quad (42)$$

$$J_{\text{VAE}}(\theta, g) = J_{\text{AE}}(\theta, g) + \mathbb{E}_{x \sim p_d} [\log p_Z(g(x))] - \gamma \quad (43)$$

$$J_{\text{fvNCE}}^{0,0}(\theta, g) = J_{\text{VAE}}(\theta, g) - \sum_x p_\theta(x, g(x)), \quad (44)$$

where $\gamma := \max_z \log p_Z(z)$ is introduced to ensure $\log p_Z(z) - \gamma \leq 0$,⁴ which allows us to obtain the following chain of inequalities,

$$J_{\text{AE}}(\theta, g) \geq J_{\text{VAE}}(\theta, g) \geq J_{\text{fvNCE}}^{0,0}(\theta, g). \quad (45)$$

$J_{\text{fvNCE}}^{0,0}$ can be also interpreted as a well-justified instance of regularized autoencoders [6]. When using tied stochastic encoders $q_0 = q_1$ satisfying $\text{supp}(p_\theta) \subseteq \text{supp}(p_{n,0})$, using the empirical version the 2nd expectation in Eq. 36 (instead of dropping it due to being a constant) can be beneficial in scenarios explicitly requiring poor reconstruction of certain inputs. The downside is a higher variance in the empirical loss and its gradients. Overall, a variational autoencoder can be generally understood as variance-reduced instance of fully variational NCE.

5.2 “Robustified” VAEs: $(\alpha, \beta) = (0, 1)$

Now we consider the pair $f_1(r) = \log(1+r)$ and $f_0(r) = \log(1+r) - r$. We read

$$\begin{aligned} J_{\text{fvNCE}}^{0,1}(\theta, q_1, q_0) &= \mathbb{E}_{(x,z) \sim p_{d,1}} \left[\log \left(1 + \frac{p_\theta(x, z)}{p_{n,1}(x, z)} \right) \right] \\ &\quad + \mathbb{E}_{(x,z) \sim p_{n,0}} \left[\log \left(1 + \frac{p_\theta(x, z)}{p_{n,0}(x, z)} \right) - \frac{p_\theta(x, z)}{p_{n,0}(x, z)} \right]. \end{aligned} \quad (46)$$

³If we use unnormalized models p_θ^0 , then Eq. 37 is bounded by $Z(\theta)$.

⁴This is only necessary for continuous latent variables as pmf’s are always in $[0, 1]$.

We assume $\text{supp}(p_\theta) \subseteq \text{supp}(p_{n,0})$, then the 3rd term can be dropped (see Sec. 5.1). With tied encoders $q = q_1 = q_0$ we arrive at a near-symmetric cost

$$\begin{aligned} J_{\text{fvNCE}}^{0,1} &\doteq \mathbb{E}_{(x,z) \sim p_{d,1}} \left[\log \left(1 + \frac{p_\theta(x,z)}{p_{n,1}(x,z)} \right) \right] + \mathbb{E}_{(x,z) \sim p_{n,0}} \left[\log \left(1 + \frac{p_\theta(x,z)}{p_{n,0}(x,z)} \right) \right] \\ &= \mathbb{E}_{(x,z) \sim p_{d,1}} \left[\log \left(1 + \frac{p_\theta(x,z)}{p_n(x)q(z|x)} \right) \right] + \mathbb{E}_{(x,z) \sim p_{n,1}} \left[\log \left(1 + \frac{p_\theta(x,z)}{p_n(x)q(z|x)} \right) \right] \\ &= \mathbb{E}_{(x,z) \sim p_{d,1}} [\text{soft-plus}(\Delta(x,z))] + \mathbb{E}_{(x,z) \sim p_{n,1}} [\text{soft-plus}(\Delta(x,z))], \end{aligned} \quad (47)$$

where we introduced the shorthand notation $\Delta(x,z) = \log p_\theta(x,z) - \log p_n(x) - \log q(z|x)$. This lower bound is tight if $q(z|x) = p_\theta(z|x) = p_\theta(x,z)/p_\theta(x)$. In this case the ratio inside the log simplifies to

$$\frac{p_\theta(x,z)}{p_n(x)q(z|x)} = \frac{p_\theta(x,z)p_\theta(x)}{p_n(x)p_\theta(x,z)} = \frac{p_\theta(x)}{p_n(x)} \quad (48)$$

and $\Delta(x,z) = \log p_\theta(x) - \log p_n(x)$. Note that $\log p_n(x)$ is expected to be small for real samples x and large for noise samples. $J_{\text{fvNCE}}^{0,1}$ can be interpreted as a version of VAEs aiming to reconstruct both real and noise samples well, but is based on a robustified reconstruction error (but with different and sample dependent truncation values for real and noise samples). In practice this cost appears to behave similar to AEs and VAEs (see Sec. 6.2 and Table 1).

5.3 Weighted squared distance: $(\alpha, \beta) = (1, 0)$

As a last example we consider $f_1(r) = r$ and $f_0(r) = -r^2/2$:

$$J_{\text{fvNCE}}^{1,0}(\theta, q_1, q_0) = \mathbb{E}_{(x,z) \sim p_{d,1}} \left[\frac{p_\theta(x,z)}{p_{n,1}(x,z)} \right] - \frac{1}{2} \mathbb{E}_{(x,z) \sim p_{n,0}} \left[\left(\frac{p_\theta(x,z)}{p_{n,0}(x,z)} \right)^2 \right] \quad (49)$$

Note that the encoder q_1 cancels in the first term, as

$$\mathbb{E}_{(x,z) \sim p_{d,1}} \left[\frac{p_\theta(x,z)}{p_{n,1}(x,z)} \right] = \sum_{x,z} \frac{p_d(x)q_1(z|x)p_\theta(x,z)}{p_n(x)q_1(z|x)} = \mathbb{E}_{\substack{x \sim p_d \\ z \sim p_Z}} \left[\frac{p_\theta(x|z)}{p_n(x)} \right]. \quad (50)$$

Therefore q_1 does not appear in the r.h.s. of Eq. 49 and can be omitted. Further, the last term in $J_{\text{fvNCE}}^{1,0}$ is the (Neyman) χ^2 -divergence between $p_\theta(X,Z)$ and $p_n(X)q_0(Z|X)$. After some algebraic manipulations it can be shown that $J_{\text{fvNCE}}^{1,0}$ is (up to constants) a weighted squared distance,

$$J_{\text{fvNCE}}^{1,0}(\theta, q_0) \doteq \frac{1}{2} \sum_{x,z} \frac{(p_\theta(x,z) - p_{d,0}(x,z))^2}{p_{n,0}(x,z)}. \quad (51)$$

Overall the aim is to minimize the weighted squared distance between the generative joint model $p_\theta(X,Z)$ and the data-encoder induced one $p_d(X)q(Z|X)$. In contrary to the setting where $\alpha = 0$ (or α is at least small) and therefore it is natural to model $\log p_\theta$, it seems more natural to model p_θ directly (instead of the log-likelihood) in Eq. 49. Hence, the choice $\alpha = 1$ is connected to density ratio estimation [30, 29], that typically uses shallow mixture models to represent the density ratio p_θ/p_n . In fact, $J_{\text{fvNCE}}^{1,0}$ in Eq. 49 is closely related to least-squares importance fitting [13] when $q_0 = q_1$.

6 Numerical Experiments

In this section we illustrate the difference in the behavior of several instances of $J_{\text{fvNCE}}^{\alpha,\beta}$ —in particular in comparison with classical autoencoders and VAEs—on toy examples.

6.1 Noise-penalized variational autoencoders

First, we demonstrate the capability to steer the behavior of an 784-256-784 autoencoder (with deterministic encoder) by using $J_{\text{fvNCE}}^{0,0}$ (Eq. 44). The noise distribution p_n is a kernel density estimate

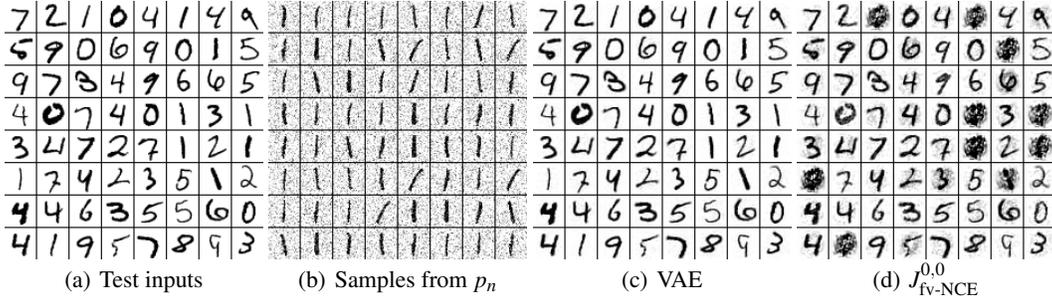


Figure 3: The impact of the 2nd term in Eq. 44 on the reconstruction of test inputs (a). p_n is chosen as a kernel-density estimator of several digits in a validation set showing “1”, with samples shown in (b). Reconstructions of the inputs using a VAE-trained encoder-decoder are given in (c), and (d) shows the corresponding reconstruction for an encoder-decoder trained using $J_{\text{fv-NCE}}^{0,0}$ (Eq. 44). Input patches showing a “1” are poorly reconstructed (as intended).

of inputs depicting the digit “1” from a validation set. Since the cost for false positives induced by $-f_0^{0,0} = r$ is higher than the cost for false negatives ($-f_1^{0,0}(r) = -\log r$), anything resembling a digit “1” is expected to be poorly reconstructed—even when those digits appear frequently in the training data. Fig. 3 visually verifies this on test inputs. This feature of Eq. 44 is useful when training data for OOD detection is contaminated by outliers, but a collection of outliers is available; or when an autoencoder-based OOD detector is required to identify certain patterns as OOD.

6.2 Stronger noise penalization using $J_{\text{fv-NCE}}^{\alpha,0}$

Since $f_0^{\alpha,0}$ penalizes false positives stronger than $f_1^{\alpha,0}$ does for false negatives, we expect different solutions for different choices of α . With infinite data and correctly specified models $\log p_\theta$, all PSRs will return the same solution (up to the issue of local maxima), but we only have finite training data and clearly underspecified models.

We fix the decoder variance to $\sigma_{\text{dec}}^2 = 1/8^2$ and use a kernel density estimate with bandwidth $\sigma_{\text{kde}} = 2\sigma_{\text{dec}}$ as noise distribution p_n . By setting $\alpha > 0$, noise samples (which are near the training data in this setting) force the model p_θ to explicitly concentrate on the training data. Samples $x \sim p_n$ have a larger reconstruction error as compared to the VAE setting ($\alpha = 0$). Table 1 lists average decoding log-likelihoods for several values of α . VAEs reconstruct noise samples worse than standard autoencoders (AEs) due to their latent code regularization. This behavior is generally amplified for increasing α , as the difference between the average reconstruction error grows with α . We also include $J_{\text{fvNCE}}^{0,1}$ (Sec. 5.2) for reference, which behaves in practice similar to VAEs. Fig. 4 visualizes the decreasing reconstruction quality of samples drawn from p_n .

In order to avoid vanishing gradients when $\alpha > 0$ in the initial training phase, in view of Cor. 1 we use actually a linear combination of $J_{\text{fvNCE}}^{\alpha,0}$ (with weight 0.9) and $J_{\text{fvNCE}}^{0,0}$ (with weight 0.1) as training loss. Table 1 lists the values for two ReLU-based MLP networks (trained from the same random initial weights) obtained after 100 epochs. Since the log-ratios such as $\log r = \log p_\theta(x, z) - \log p_{n,1}(x, z)$ can attain large magnitudes, expressions such as r^α and $r^{\alpha+1}$ are evaluated using a “clipped” exponential function: we use the first-order approximation $e^T(u - T + 1)$ when $u > T$ for a threshold value T , which is chosen as $T = 10$ in our implementation.

7 Conclusion

In this work we propose fully variational noise-contrastive estimation as a tractable method to apply noise-contrastive estimation on latent variable models. As with most variational inference methods, the resulting empirical loss only needs samples from the data, noise and encoder distributions. We are largely interested in the existence and basic properties of such framework and unravel a connection with variational autoencoders. In light of this connection, VAEs are now justified to be steered explicitly towards poorly reconstructing samples from a user-specified noise distribution.

(a) 784-128-784				(b) 784-256-128-256-784			
Method/ (α, β)	$x \sim p_d$	$x \sim p_n$	Difference	Method/ (α, β)	$x \sim p_d$	$x \sim p_n$	Difference
AE	766	-1138	1904	AE	756	-919	1675
VAE	765	-1609	2374	VAE	769	-1373	2142
(1/256, 0)	749	-1665	2414	(1/256, 0)	774	-1402	2176
(1/64, 0)	698	-1863	2561	(1/64, 0)	748	-1520	2268
(1/16, 0)	753	-1818	2571	(1/16, 0)	777	-1463	2240
(0, 1)	736	-1656	2392	(0, 1)	775	-1372	2147

Table 1: Average log-likelihood $\log p_\theta(x|g(x))$ in nats. Higher values indicate lower reconstruction error.

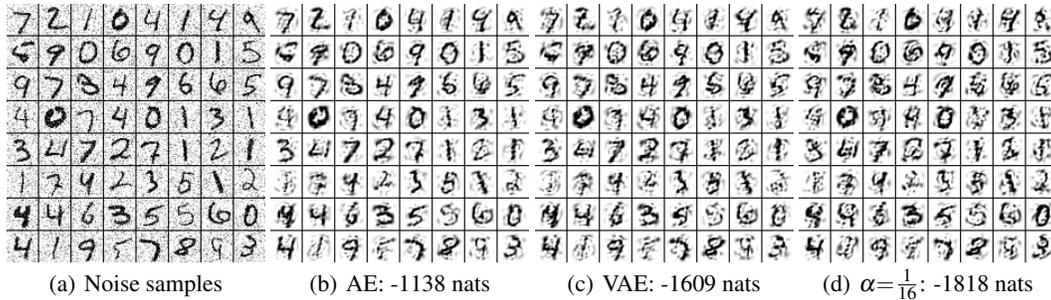


Figure 4: Reconstruction of samples $x \sim p_n$ (a). VAEs (c) and $J_{\text{fvNCE}}^{\alpha, 0}$ (d) increasingly force such samples to be poorly reconstructed compared to AEs (b), while maintaining a similar reconstruction error for training data $x \sim p_d$ (see Table 1).

The utility of our framework for improved OOD detection and enabling general energy-based decoder models is left as future work. Further, the highly asymmetric nature of the classification loss suggests a potential but yet-to-explore connection with one-class SVMs [27] and support vector data description [31].

Acknowledgement This work was partially supported by the Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by the Knut and Alice Wallenberg Foundation.

References

- [1] Ciwan Ceylan and Michael U Gutmann. Conditional noise-contrastive estimation of unnormalised models. In *International Conference on Machine Learning*, pages 726–734. PMLR, 2018.
- [2] Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. *Advances in neural information processing systems*, 29, 2016.
- [3] Alexander Philip Dawid and Monica Musio. Theory and applications of proper scoring rules. *Metron*, 72(2):169–183, 2014.
- [4] Peter Dayan, Geoffrey E Hinton, Radford M Neal, and Richard S Zemel. The helmholtz machine. *Neural computation*, 7(5):889–904, 1995.
- [5] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021.
- [6] Partha Ghosh, Medhi SM Sajjadi, Antonio Vergari, and Michael Black. From variational to deterministic autoencoders. In *8th International Conference on Learning Representations*, 2020.
- [7] Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American statistical Association*, 102(477):359–378, 2007.

- [8] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, 2014.
- [9] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 297–304, 2010.
- [10] Michael U Gutmann and Aapo Hyvärinen. Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics. *Journal of Machine Learning Research*, 13(Feb):307–361, 2012.
- [11] Aapo Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4), 2005.
- [12] Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233, 1999.
- [13] Takafumi Kanamori, Shohei Hido, and Masashi Sugiyama. A least-squares approach to direct importance estimation. *The Journal of Machine Learning Research*, 10:1391–1445, 2009.
- [14] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *Proceedings of the 2nd International Conference on Learning Representations (ICLR)*, 2014.
- [15] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *Advances in neural information processing systems*, 31, 2018.
- [16] Polina Kirichenko, Pavel Izmailov, and Andrew G Wilson. Why normalizing flows fail to detect out-of-distribution data. *Advances in neural information processing systems*, 33:20578–20589, 2020.
- [17] Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. *Advances in Neural Information Processing Systems*, 33:21464–21475, 2020.
- [18] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [19] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [20] Matthew Parry, A Philip Dawid, Steffen Lauritzen, et al. Proper local scoring rules. *Annals of Statistics*, 40(1):561–592, 2012.
- [21] Miika Pihlaja, Michael Gutmann, and Aapo Hyvärinen. A family of computationally efficient and simple estimators for unnormalized statistical models. In *Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*, pages 442–449, 2010.
- [22] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [23] Jie Ren, Peter J Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark Depristo, Joshua Dillon, and Balaji Lakshminarayanan. Likelihood ratios for out-of-distribution detection. *Advances in neural information processing systems*, 32, 2019.
- [24] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538. PMLR, 2015.
- [25] Danilo J Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 1278–1286, 2014.
- [26] Benjamin Rhodes and Michael U Gutmann. Variational noise-contrastive estimation. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2741–2750. PMLR, 2019.
- [27] Bernhard Schölkopf, John C Platt, John Shawe-Taylor, Alex J Smola, and Robert C Williamson. Estimating the support of a high-dimensional distribution. *Neural computation*, 13(7):1443–1471, 2001.
- [28] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2020.

- [29] Masashi Sugiyama, Taiji Suzuki, and Takafumi Kanamori. *Density ratio estimation in machine learning*. Cambridge University Press, 2012.
- [30] Masashi Sugiyama, Taiji Suzuki, and Takafumi Kanamori. Density-ratio matching under the bregman divergence: a unified framework of density-ratio estimation. *Annals of the Institute of Statistical Mathematics*, 64(5):1009–1044, 2012.
- [31] David MJ Tax and Robert PW Duin. Support vector data description. *Machine learning*, 54(1):45–66, 2004.
- [32] Houssam Zenati, Manon Romain, Chuan-Sheng Foo, Bruno Lecouat, and Vijay Chandrasekhar. Adversarially learned anomaly detection. In *2018 IEEE International conference on data mining (ICDM)*, pages 727–736. IEEE, 2018.