
Provable benefits of annealing for estimating normalizing constants

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recent research has developed several Monte Carlo methods for estimating the
2 normalization constant (partition function) based on the idea of annealing. This
3 means sampling successively from a path of distributions which interpolate be-
4 tween a tractable “proposal” distribution and the unnormalized “target” distribution.
5 Prominent estimators in this family include annealed importance sampling and
6 annealed noise-contrastive estimation (NCE). Such methods hinge on a number
7 of design choices: which estimator to use, which path of distributions to use and
8 whether to use a path at all; so far, there is no definitive theory on which choices are
9 efficient. Here, we evaluate each design choice by the asymptotic estimation error it
10 produces. First, we show that using NCE is more efficient than the importance sam-
11 pling estimator, but in the limit of infinitesimal path steps, the difference vanishes.
12 Second, we find that using the geometric path brings down the estimation error
13 from an exponential to a polynomial function of the parameter distance between
14 the target and proposal distributions. Third, we find that the arithmetic path, while
15 rarely used, can offer optimality properties over the universally-used geometric
16 path. In fact, in a particular limit, the optimal path is arithmetic. Based on this
17 theory, we finally propose a two-step estimator to approximate the optimal path in
18 an efficient way.

19 1 Introduction

20 Recent progress in generative modeling has sparked renewed interest in models of data that are defined
21 by an unnormalized distribution. One prominent example is energy-based (or score-based) models,
22 which are increasingly used in deep learning [1], and for which there are a variety of parameter
23 estimation procedures [2–5]. Another example comes from Bayesian statistics, where the posterior
24 model of parameters given data is frequently known only up to a proportionality constant. Such
25 models can be evaluated and compared by the probability they assign to a dataset, yet this requires
26 computing their normalization constants (partition functions) which are typically high-dimensional,
27 intractable integrals.

28 Monte-Carlo methods have been successful at computing these integrals using sampling methods [6].
29 The most common is importance sampling [6] which draws a sample from a tractable, “proposal”
30 distribution to integrate the unnormalized “target” density. Noise-contrastive estimation (NCE) [3]
31 uses a sample from *both* the proposal and the target, to compute the integral. Yet such methods
32 suffer from high variance, especially when the “gap” between the proposal and target densities is
33 large [7–9]. This has motivated various approaches to gradually bridge the gap with intermediate
34 distributions, which is loosely referred to as “annealing”. Among them, annealed importance sampling
35 (AIS) [10–12] is widely adopted: it has been used to compute the normalization constants of deep
36 stochastic models [13, 14] or to motivate a lower-bound for learning objectives [15, 16]. To integrate

Name	Loss identified by $\phi(\mathbf{x})$	Estimator \hat{Z}_1	MSE
IS	$x \log x$	$\mathbb{E}_{p_0} \frac{f_1}{f_0}$	$\frac{1+\nu}{\nu N} \mathcal{D}_{\chi^2}(p_1, p_0)$
RevIS	$-\log x$	$(\mathbb{E}_{p_1} \frac{f_0}{f_1})^{-1}$	$\frac{1+\nu}{N} \mathcal{D}_{\chi^2}(p_0, p_1)$
NCE	$x \log x - (1+x) \log(\frac{1+x}{2})$	implicit	$\frac{(1+\nu)^2}{\nu N} \frac{\mathcal{D}_{\text{HM}}(p_1, p_0)}{1 - \mathcal{D}_{\text{HM}}(p_1, p_0)}$

Table 1: Some estimators of the normalization obtained by minimizing a classification loss, and their estimation error in terms of well-known divergences [22]. For details and definitions, see Appendix A.

37 the unnormalized "target" density, it draws a sample from an entire path of distributions between
38 the proposal and the target. While annealed importance sampling has been shown to be effective
39 empirically, its theoretical understanding remains limited [17, 18]: it is yet unclear when annealing is
40 effective, for which annealing paths, and whether AIS is a statistically efficient way to do it.

41 In this paper, we define a family of annealed Bregman estimators (ABE) for the normalization constant.
42 We show that it is general enough to recover many classical estimators as a special case, including
43 importance sampling, noise-contrastive estimation, umbrella sampling [19], bridge sampling [20]
44 and annealed importance sampling. We provide a statistical analysis of its hyperparameters such
45 as the choice of paths, and show the following. First, we show that using NCE is more efficient
46 than the importance sampling estimator, but in the limit of infinitesimal path steps, the difference
47 vanishes. Second, we find that the near-universally used geometric path brings down the estimation
48 error from an exponential to a polynomial function of the parameter distance between the target
49 and proposal distributions. Third, we find that using the recently introduced arithmetic path [21] is
50 exponentially inefficient in its basic form, yet it can be reparameterized to be in some sense optimal.
51 Based on this optimality result, we finally propose a two-stage estimation procedure which first finds
52 an approximation of the optimal (arithmetic) path, then uses it to estimate the normalization constant.

53 2 Background

54 **Importance sampling and NCE** The problem considered here is how to compute the normalization
55 constant ¹, *i.e.* the integral of some unnormalized density $f_1(\mathbf{x})$ called "target". Importance sampling
56 and noise-contrastive estimation are two common estimators which integrate the unnormalized target
57 over a random sample drawn from a tractable density $p_0(\mathbf{x})$ called "proposal" (Table 1, column 3). In
58 fact, they are part of a larger family of Monte-Carlo estimators which can be interpreted as solving a
59 binary classification task, aiming to distinguish between a sample drawn from the proposal and another
60 from the target [22]. Each estimator is obtained by minimizing a specific binary classification loss
61 that is identified by a convex function $\phi(\mathbf{x})$. For example, minimizing the classification loss identified
62 by $\phi_{\text{IS}}(x) = x \log x$ yields the importance sampling estimator. Similarly, $\phi_{\text{RevIS}}(x) = -\log x$ leads
63 to the reverse importance sampling estimator [23], and $\phi_{\text{NCE}}(x) = x \log x - (1+x) \log((1+x)/2)$
64 to the noise-contrastive estimator. These estimators are summarized in Table 1.

65 **Annealed estimators** Annealing extends the above "binary" setup, by introducing a sequence of
66 $K + 1$ distributions from the proposal to the target (included). The idea will be to draw a sample from
67 *all* these distributions to integrate the target $f_1(\mathbf{x})$. These intermediate distributions are obtained
68 from a path $(f_t)_{t \in [0,1]}$, defined by interpolating between the proposal p_0 and unnormalized target f_1 :
69 this path is therefore unnormalized. Different interpolation schemes can be chosen. A general one,
70 explained in Figure 1, is to take the q -mean of the proposal and target [21]. Two values of q are of
71 particular interest: $q \rightarrow 0$ defines a near-universal path [18], obtained by taking the geometric mean
72 of the target and proposal, while $q = 1$ defines a path obtained by the arithmetic mean.

73 Once a path is chosen, it can be uniformly² discretized into a sequence of $K + 1$ unnormalized
74 densities, denoted by $(f_{k/K})_{k \in [0,K]}$ with corresponding normalizations $(Z_{k/K})_{k \in [0,K]}$. In practice,
75 samples are drawn from the corresponding normalized densities $(p_{k/K})_{k \in [0,K]}$ using Markov Chain

¹in this paper we also say we "estimate" the the normalization constant, though this terminology is unconventional as estimation traditionally refers to the *parameters* of a statistical model

²other discretization schemes can be equivalently achieved by re-parameterizing the path [17]

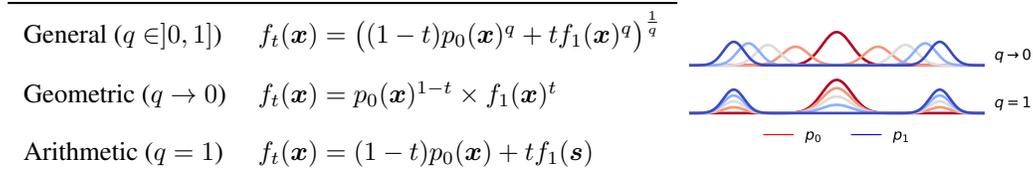


Figure 1: Geometric and arithmetic path between the proposal and target distributions. Here, the proposal (red) is a standard gaussian. The target (blue) is a gaussian mixture with two modes, and same first and second moments as the proposal.

76 Monte Carlo (MCMC). This sampling step incurs a computational cost, which is paid in the hope
77 of reducing the variance of the estimation. It is common in the literature [17, 18] to assume *perfect*
78 *sampling*, meaning the MCMC has converged and produced exact and independent samples from the
79 distributions along the path, which simplifies the analysis.

80 **Estimation error** A measure of "quality" is required to compare different estimation choices,
81 such as whether to anneal and which path to use. Such a measure is given by the Mean Squared
82 Error (MSE), which is generally tractable when written at the first order in the asymptotic limit of a
83 large sample size [24, Eq. 5.20]. These expressions have been derived for estimators obtained by
84 minimizing a classification loss [22] and are included in table 1. They measure the "gap" between
85 the proposal and target distributions using statistical divergences. Note also that the estimation error
86 depends on the *normalized* target density (column 4), while the estimators are computed using the
87 *unnormalized* target density (column 3). Further details are available in Appendix A.

88 3 Annealed Bregman Estimators of the normalization constant

89 The question that we try to answer in this paper is: *How should we choose the $K + 1$ distributions in*
90 *annealing, and how are their samples best used?* To answer this, we will study the error produced by
91 different estimation choices. But first we define the set of estimators for which the analysis is done.

92 **Definition of Annealed Bregman Estimators** We now define a new family of estimators, which
93 we call annealed Bregman estimators (ABE); this terminology is explained later. We will show that
94 this is a general class of estimators for computing the normalization using a sample drawn from the
95 sequence of $K + 1$ distributions. For ABE, the log normalization $\log Z_1$ is estimated additively along
96 the sequence of distributions

$$\widehat{\log Z_1} = \sum_{k=0}^{K-1} \log \left(\frac{Z_{(k+1)/K}}{Z_{k/K}} \right) + \log Z_0 . \quad (1)$$

97 Defining the estimation in log-space is analytically convenient, as it is easier to analyze a sum of
98 estimators than a product. Exponentiating the result leads to an estimator of Z_1 . We naturally extend
99 the binary setup ($K = 1$) of Chehab et al. [22] and propose to compute each of the intermediate
100 log-ratios, by solving a classification task between samples drawn from their corresponding densities
101 $p_{k/K}$ and $p_{(k+1)/K}$. Each binary classification loss is now identified by a convex function $\phi_k(\mathbf{x})$ and
102 defined as

$$\mathcal{L}_k(\beta) := \mathbb{E}_{\mathbf{x} \sim p_{k/K}} [\phi'_k(r(\mathbf{x}; \beta)) \times r(\mathbf{x}; \beta) - \phi_k(r(\mathbf{x}))] - \mathbb{E}_{\mathbf{x} \sim p_{(k+1)/K}} [\phi'_k(r(\mathbf{x}; \beta))] \quad (2)$$

103 where $r(\mathbf{x}; \beta)$ is parameterized by the unknown log-ratio $\beta^* = \log(Z_{(k+1)/K}/Z_{k/K})$

$$r(\mathbf{x}; \beta) = \exp(-\beta) \times f_{(k+1)/K}(\mathbf{x})/f_{k/K}(\mathbf{x}) . \quad (3)$$

104 The convex functions $(\phi_k)_{k \in [0, K-1]}$ which identify the classification losses are called "Bregman"
105 generators, hence ABE. As mentioned above, we assume perfect sampling and allocate the total
106 sample size N equally among the K estimators in the sum.

107 **Hyperparameters** The annealed Bregman estimator depends on the following hyperparameters:
 108 (1) the choice of path q ; (2) the number of distributions along that path $K + 1$ (including the proposal
 109 and the target); (3) the classification losses identified by the convex functions $(\phi_k)_{k \in [0, K-1]}$.

110 Different combinations of these hyperparameters recover several common estimators of the log-
 111 partition function. In binary case of $K = 1$ this includes importance sampling, reverse importance
 112 sampling, and noise-contrastive estimation, each obtained for a different choice of the classification
 113 loss [22]. To build intuition, consider $K = 2$ so that we add a single intermediate distribution $p_{1/2}$
 114 to the sequence. Using the importance sampling loss ($\phi_0 = x \log x$) for the first ratio, and reverse
 115 importance sampling ($\phi_1 = -\log x$) for the second ratio, recovers the *bridge sampling estimator* as
 116 a special case [20]

$$\widehat{\log Z_1} = -\log \mathbb{E}_{p_1} \frac{f_{1/2}}{f_1} + \log \mathbb{E}_{p_0} \frac{f_{1/2}}{f_0} \log Z_0 . \quad (4)$$

117 Alternatively, we can use these classification losses in reverse order: reverse importance sampling
 118 ($\phi_0 = -\log x$) for the first ratio, and importance sampling ($\phi_1 = x \log x$) for the second ratio, and
 119 recover the *umbrella sampling estimator* [19] also known as the *ratio sampling estimator* [25]

$$\widehat{\log Z_1} = \log \mathbb{E}_{p_{1/2}} \frac{f_1}{f_{1/2}} - \log \mathbb{E}_{p_{1/2}} \frac{f_0}{f_{1/2}} \log Z_0 . \quad (5)$$

120 Another option yet, is to use the same classification loss for all ratios. With importance sampling
 121 ($\phi_k = x \log x, \forall k \in [0, K-1]$), we recover the *annealed importance sampling estimator* [10–12]

$$\widehat{\log Z_1} = \sum_{k=1}^K \log \mathbb{E}_{\mathbf{x} \sim p_{k-1}} \left[\frac{f_k}{f_{k-1}}(\mathbf{x}) \right] + \log Z_0 . \quad (6)$$

122 The family of annealed Bregman estimators is visibly large enough to include many existing esti-
 123 mators, obtained for different hyperparameter choices. This raises the fundamental question of how
 124 these hyperparameters should be chosen, in particular in the challenging case where the *target and*
 125 *proposal have little overlap and the data is high dimensional*. To answer this question, we will study
 126 the estimation error produced by different hyperparameter choices.

127 4 Statistical analysis of the hyperparameters

128 We consider a fixed data budget N and investigate how the remaining hyperparameters are best
 129 chosen for statistical efficiency. The starting point for the analysis is that as ABE estimates the
 130 normalization in log-space, the estimator is obtained by a sum of independent and asymptotically
 131 unbiased estimators [26] given in Eq. 1 and thus the mean squared errors written in table 1 are
 132 additive. Each individual error actually measures an overlap between two consecutive distributions
 133 along the path, and annealing integrates these overlaps.

134 4.1 Classification losses, ϕ_k

135 Given the popularity of annealed importance sampling, we should first ask if the importance sampling
 136 loss is really an acceptable default. We recall an important limitation of annealed importance
 137 sampling [27]: its estimation error is notoriously sensitive to distribution tails. Without annealing, it
 138 is infinite when the target p_1 has a heavier tail than the proposal p_0 . When annealing with a geometric
 139 path, for example between two gaussians with different covariances $p_0 = \mathcal{N}(\mathbf{0}, \mathbf{Id})$ and $p_1 =$
 140 $\mathcal{N}(\mathbf{0}, 2\mathbf{Id})$, the geometric path produces gaussians with increasing variances $\Sigma_t = (1 - t/2)^{-1} \mathbf{Id}$
 141 and therefore increasing tails. Hence, the same tail mismatch holds along the path. Note that this
 142 concern is a realistic one for natural image data, as the target distribution over images is typically
 143 super-gaussian [28] while the proposal is usually chosen as gaussian.

144 This warrants a better choice for the loss: In the binary setup ($K = 1$), the NCE loss is optimal [20, 22]
 145 and its error can be orders of magnitude less than importance sampling [20]. We extend this optimality
 146 result over a sequence of distributions $K > 1$ and also show that the gap between annealed IS and
 147 annealed NCE is closed in the limit of a continuous path:

148 **Theorem 1** (Estimation error and the Fisher-Rao path length) *For a finite value of K , the optimal*
 149 *loss is NCE*

$$\text{MSE}(p_0, p_1; q, K, N, \phi_{\text{NCE}}) \leq \text{MSE}(p_0, p_1; q, K, N, \phi), \quad \forall q, K, N, \phi. \quad (7)$$

150 *In the limit of $K \rightarrow \infty$, NCE, IS, and revIS converge to the same estimation error, given by the*
 151 *Fisher-Rao path length from the proposal to the target*

$$\text{MSE}(p_0, p_1; q, K, N, \phi) \rightarrow \frac{1}{N} \int_0^1 I(t) dt, \quad \text{for } K \rightarrow \infty, \phi \in \{\phi_{\text{NCE}}, \phi_{\text{IS}}, \phi_{\text{RevIS}}\} \quad (8)$$

152 *where the Fisher-Rao metric $I(t) := \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}, t)} [(\frac{d}{dt} \log p_t(\mathbf{x}))^2]$ defined as the Fisher information*
 153 *over the path, using time t as the parameter.*

154 This is proven in Appendix B. While the NCE estimator requires solving a (potentially non-convex)
 155 scalar optimization problem in Eq. 2 and IS does not, this is the price to pay for statistical optimality.
 156 In the following we will keep the optimal NCE loss and will indicate the dependency of the estimation
 157 error on ϕ_{NCE} with a subscript, instead. We highlight that our theorems in this paper apply to the
 158 MSE in the limit of $K \rightarrow \infty$: their results hold the same for the IS and RevIS losses by virtue of
 159 theorem 1. Just as in the binary case, while the estimator is computed with the *unnormalized* path of
 160 densities (Eq. 2), the estimation error depends on the *normalized* path of densities (Eq. 8).

161 4.2 Number of distributions, $K + 1$

162 It is known that estimating the normalization constant using plain importance sampling ($K = 1$) can
 163 produce a statistical error than is exponential in the dimension [6, Example 9.1]. We show that in the
 164 binary case, NCE also suffers from an estimation error that scales exponentially with the dimension.

165 In the following, we consider a proposal p_0 and target p_1 that are in an exponential family; note that
 166 certain exponential families have universal approximation capabilities [29]. The exponential family
 167 is defined as

$$p(\mathbf{x}; \boldsymbol{\theta}) := \exp(\langle \boldsymbol{\theta}, \mathbf{t}(\mathbf{x}) \rangle - \log Z(\boldsymbol{\theta})) \quad (9)$$

168 where $Z(\boldsymbol{\theta}) = \int \exp(\langle \boldsymbol{\theta}_1, \mathbf{t}(\mathbf{x}) \rangle)$. We will consider that the (unnormalized) target density f_1 is what
 169 we call a *simply unnormalized model* defined as

$$f_1(\mathbf{x}) = \exp(\langle \boldsymbol{\theta}_1, \mathbf{t}(\mathbf{x}) \rangle) \quad (10)$$

170 Note that in general, a pdf can be unnormalized in many ways: one can multiply an unnormalized
 171 density by any positive function of $\boldsymbol{\theta}$ and it will still be unnormalized. However, the simple and
 172 intuitive case defined above is what we encounter in the analysis below.

173 For exponential families, the log-normalization $\log Z(\boldsymbol{\theta})$ is a convex function (“log-sum-exp”) of the
 174 parameter $\boldsymbol{\theta}$ [30], which implies $0 \preceq \nabla_{\boldsymbol{\theta}}^2 \log Z(\boldsymbol{\theta})$. In our theorems we use the further assumptions
 175 of strong convexity with constant M , and/or smoothness with constant L (gradient is L -Lipschitz):

$$\nabla_{\boldsymbol{\theta}}^2 \log Z(\boldsymbol{\theta}) \succeq M \mathbf{Id} \quad (11)$$

176

$$\nabla_{\boldsymbol{\theta}}^2 \log Z(\boldsymbol{\theta}) \preceq L \mathbf{Id} \quad (12)$$

177 For exponential families, the derivatives of the log partition function yield moments of the sufficient
 178 statistics, so we are effectively assuming that the eigenvalues of $\nabla_{\boldsymbol{\theta}}^2 \log Z(\boldsymbol{\theta}) = \text{Cov}_{\mathbf{x} \sim p}[\mathbf{t}(\mathbf{x})]$ are
 179 bounded, which will be the case for parameters in a bounded domain $\boldsymbol{\theta} \in \Theta$. An example along with
 180 the proofs of the following theorems 2 and 3, are provided in Appendix B.

181 **Theorem 2** (Exponential error of binary NCE) *Assume the proposal p_0 is from the normalized*
 182 *exponential family, while the (unnormalized) target f_1 is from the simply unnormalized exponential*
 183 *family (Eq. 10). The log-partition function $\log Z(\boldsymbol{\theta})$ is assumed to be strongly convex (Eq. 11).*
 184 *Then in the binary case $K = 1$, the estimation error of NCE is (at least) exponential in the parameter-*
 185 *distance between the proposal and the target*

$$\text{MSE}_{\text{NCE}}(p_0, p_1; q, K, N) \geq \frac{4}{N} \exp\left(\frac{1}{8} M \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\|^2\right) - 1, \quad \text{when } K = 1 \quad (13)$$

186 *where M is the strong convexity constant of $\log Z(\boldsymbol{\theta})$.*

187 Annealing the importance sampling estimator (increasing K) was proposed in hope that we can trade
 188 the statistical cost in the dimension for a computational cost (number of classification tasks) which is
 189 more acceptable. Yet, there is no definitive theory on the ability of annealing to reduce the statistical
 190 cost in a general setup [18, 21]. For both importance sampling and noise-contrastive estimation, we
 191 prove that annealing with the near-universal geometric path brings down the estimation error, from
 192 exponential to polynomial in the parameter-distance between the proposal and target. Given that
 193 the parameter-distance itself scales as \sqrt{D} with the dimension, using these paths effectively makes
 194 annealed estimation amenable to high-dimensional problems. This corroborates empirical [31] and
 195 theoretical [17] results which suggested in simple cases that annealing with an appropriate path can
 196 reduce the estimator error up to several orders of magnitude.

197 **Theorem 3** (Polynomial error of annealed NCE with a geometric path) *Assume the proposal p_0*
 198 *is from the normalized exponential family, while the (unnormalized) target f_1 is from the simply*
 199 *unnormalized exponential family (Eq. 10). The log-partition function $\log Z(\theta)$ is assumed to be*
 200 *strongly convex and smooth (Eq. 11, Eq. 12).*
 201 *Then in the annealing limit of a continuous path $K \rightarrow \infty$, the estimation error of annealed NCE with*
 202 *the geometric path is (at most) polynomial in the parameter-distance between the proposal and the*
 203 *target*

$$\text{MSE}_{\text{NCE}}(p_0, p_1; q, K, N) \leq \frac{L^2}{MN} \|\theta_1 - \theta_0\|^2, \quad \text{when } K \rightarrow \infty, q = 0 \quad (14)$$

204 where M and L are respectively the strong convexity and smoothness constants of $\log Z(\theta)$.

205 To our knowledge, this is the first result building on Gelman and Meng [17, Table 1] and Grosse et al.
 206 [18] which showcases the benefits of annealed estimation for a general target distribution.

207 We conclude that annealing with the near-universal geometric path provably benefits noise-contrastive
 208 estimation, as well as importance sampling and reverse importance sampling, when the proposal and
 209 target distributions have little overlap.

210 4.3 Path parameter, q — geometric vs. arithmetic

211 Despite the near-universal popularity of the geometric path ($q \rightarrow 0$), it is worth asking if there are
 212 other simple paths that are more optimal. Interpolating moments of exponential families was shown
 213 to outperform the geometric path by Grosse et al. [18], yet building such a path requires knowing
 214 the exponential family of the target. Other paths based on the arithmetic mean (and generalizations)
 215 of the target and proposal, were proposed in Masrani et al. [21], without a definitive theory of the
 216 estimation error.

217 Next, we analyze the error of the arithmetic path. We prove that the arithmetic path ($q = 1$) does *not*
 218 exhibit the same benefits as the geometric path: in general, its estimation error grows exponentially in
 219 the parameter-distance between the target and proposal distributions. However, in the case where an
 220 oracle gives us the normalization Z_1 to be used only in the construction of the path (we will discuss
 221 what this means in practice below), the arithmetic path can be reparameterized so as to bring down
 222 the estimation error to polynomial, even constant, in the parameter-distance. We start by the negative
 223 result.

224 **Theorem 4** (Exponential error of annealed NCE with an arithmetic path) *Assume the proposal p_0*
 225 *is from the normalized exponential family, while the (unnormalized) target f_1 is from the simply*
 226 *unnormalized exponential family (Eq. 10). The log-partition function $\log Z(\theta)$ is assumed to be*
 227 *strongly convex (Eq. 11).*
 228 *Consider the annealing limit of a continuous path $K \rightarrow \infty$ path, and of a far-away target $\|\theta_1 - \theta_0\| \rightarrow$*
 229 *∞ . For estimating the log normalization of the (unnormalized) target density f_1 , the estimation error*
 230 *of annealed NCE with the arithmetic path is (at least) exponential in the parameter-distance between*
 231 *the proposal and the target.*

$$\text{MSE}_{\text{NCE}}(p_0, p_1; q, K, N) = O\left(f\left(\frac{1}{N} \exp\left(\frac{M}{2} \|\theta_1 - \theta_0\|^2\right)\right)\right), \quad \text{when } K \rightarrow \infty, q = 1 \quad (15)$$

232 where f is an increasing function defined in Appendix B.3.

Path name	Unnormalized density	Normalized density	Error
Geometric	$f_t(\mathbf{x}) = p_0(\mathbf{x})^{1-t} f_1(\mathbf{x})^t$	$p_t(\mathbf{x}) \propto p_0(\mathbf{x})^{1-t} p_1(\mathbf{x})^t$	poly
Arithmetic	$f_t(\mathbf{x}) = (1 - w_t)p_0(\mathbf{x}) + w_t f_1(\mathbf{x})$	$p_t(\mathbf{x}) = (1 - \tilde{w}_t)p_0(\mathbf{x}) + \tilde{w}_t p_1(\mathbf{x})$	
vanilla	$w_t = t$	$\tilde{w}_t = \frac{t Z_1}{(1-t) + t Z_1}$	exp
oracle	$w_t = \frac{t}{t + Z_1(1-t)}$	$\tilde{w}_t = t$	poly
oracle-trig	$w_t = \frac{\sin^2\left(\frac{\pi t}{2}\right)}{\sin^2\left(\frac{\pi t}{2}\right) + Z_1 \cos^2\left(\frac{\pi t}{2}\right)}$	$\tilde{w}_t = \sin^2\left(\frac{\pi t}{2}\right)$	const

Table 2: Geometric and arithmetic paths, defined in the space of unnormalized densities (second column); ‘oracle’ and ‘oracle-trig’ are reparameterizations of the arithmetic path which depend on the true normalization Z_1 . The corresponding normalized densities (third column) produce an estimation error (fourth column) which we quantify.

233 We suggest an intuitive explanation for this negative result. We begin with the observation that the
234 estimation error (Eq. 8) depends on the *normalized* path of densities. Suppose the target model
235 is rescaled by a constant 100, so that the new unnormalized target density is $f_1(\mathbf{x}) \times 100$ and its
236 new normalization is $Z_1 \times 100$. Looking at table 2, this rescaling does not modify the geometric
237 path of normalized densities, while it does the arithmetic path of normalized densities. Because the
238 estimation error depends on path of normalized densities, this makes the arithmetic choice sensitive
239 to target normalization, even more so as the parameter distance grows and the log-normalization
240 with it, as a strongly convex function of it (Appendix, Eq. 85). This suggests making the arithmetic
241 path of normalized distributions ‘robust’ to the choice of Z_1 . We will show this can be achieved by
242 re-parameterizing the path in terms of Z_1 .

243 We next prove that certain reparameterizations can bring down the error to a polynomial and even
244 constant function of the parameter-distance between the target and proposal. The following theorems
245 may seem purely theoretical, as if necessitating an oracle for Z_1 , but they will actually lead to an
246 efficient estimation algorithm later.

247 **Theorem 5** (Polynomial error of annealed NCE with an arithmetic path and oracle) *Assume the same*
248 *as in Theorem 4, replacing the strong convexity of the log-partition by smoothness (Eq. 12). Addition-*
249 *ally, suppose an oracle gives the normalization constant Z_1 to be used only in the reparameterization*
250 *of the arithmetic path with $t \rightarrow \frac{t}{t + Z_1(1-t)}$ (see Table 2). This brings down the estimation error of*
251 *annealed NCE to (at most) polynomial in the parameter-distance*

$$\text{MSE}_{\text{NCE}}(p_0, p_1; q, K, N) \leq \frac{1}{N} (2 + L \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\|^2), \quad \text{when } K \rightarrow \infty, q = 1 \quad (16)$$

252 where L is the smoothness constant of $\log Z(\boldsymbol{\theta})$.

253 In fact, supposing we have (oracle) access to the normalizing constant Z_1 , the arithmetic path can
254 even be reparameterized such that it is the optimal path in a certain limit. We next prove such
255 optimality in the limits of a continuous path $K \rightarrow \infty$ and ‘far-away’ target and proposal:

256 **Theorem 6** (Constant error of annealed NCE with an arithmetic path and oracle) *Suppose we*
257 *can successively take the limit of a continuous annealing path $K \rightarrow \infty$, then the limit of a target*
258 *distribution that tends toward no overlap with the proposal $p_1(\mathbf{x})p_0(\mathbf{x}) \rightarrow 0$ pointwise (and assuming*
259 *domination by an integrable function)³. Then the optimal annealing path convergences pointwise to*
260 *an arithmetic path reparameterized trigonometrically with $t \rightarrow \frac{t}{\sin^2\left(\frac{\pi t}{2}\right) + Z_1(1 - \sin^2\left(\frac{\pi t}{2}\right))}$. In that limit,*
261 *the estimation error is constant with respect to the parameter-distance*

$$\text{MSE}_{\text{NCE}}(p_0, p_1; q, K, N) \sim \frac{2\pi^2}{N}, \quad \text{when } K \rightarrow \infty, q = 1 \quad (17)$$

262 **Two-step estimation** Thus, we see that, perhaps unsurprisingly, the ‘optimal’ mixture weights in
263 the space of unnormalized densities depends on the true Z_1 : however, this dependency is simple. We

³this effectively assumes that $K \rightarrow \infty$ faster than $p_1(\mathbf{x})p_0(\mathbf{x}) \rightarrow 0$, so that the error in the first limit is dominated by the error in the second limit.

264 propose a two-step estimation method: first, Z_1 is pre-estimated, for example using the geometric
 265 path; second, the estimate of Z_1 is plugged into the "oracle" or "oracle-trig" weight of the arithmetic
 266 path (table 2, column 2), and which is used to obtain a second estimation of Z_1 . Note that pre-
 267 estimating a problematic (hyper)parameter, here Z_1 , has proved beneficial to reduce the estimation
 268 error of NCE in a related context [32].

269 5 Numerical results

270 We now present numerical evidence for our theory and validate our two-step estimators. Importantly,
 271 we do *not* claim to achieve state of the art in terms of practical evaluation of the normalization
 272 constants; our goal is to support our theoretical analysis. We follow the evaluation methods of
 273 importance sampling literature [18] and evaluate our methods on synthetic gaussians. This setup
 274 is specially convenient for validating our theory: the optimal estimation error can conveniently be
 275 computed in closed-form, so too can the geometric and arithmetic paths which avoids a sampling error
 276 from MCMC algorithms. These derivations are included in the Appendix B. We specifically consider
 277 the high-dimensional setting, where the computation of the determinant of a high-dimensional
 278 (covariance) matrix which appears in the normalization of a gaussian, can in fact be challenging [33].

279 **Numerical Methods** The proposal distribution is always a standard gaussian, while the target
 280 differs by the second moment: $p_1 = \mathcal{N}(\mathbf{0}, 2\mathbf{Id})$ in Figure 2, $p_1 = \mathcal{N}(\mathbf{0}, 0.25\mathbf{Id})$ in Figure 3b
 281 and $p_1 = \mathcal{N}(\mathbf{0}, \sigma^2\mathbf{Id})$ in Figure 3a, where the target variance decreases as $\sigma(i) = i^{-1}$ so that the
 282 (natural) parameter distance grows linearly [30, Part II-4]. We use a sample size of $N = 50000$
 283 points, and, unless otherwise mentioned, $K + 1 = 10$ distributions from the annealing paths and
 284 the dimensionality is 50. To compute an estimator of the normalization constant using the non-
 285 convex NCE loss, we used a non-linear conjugate gradient scheme implemented in Scipy [34]. We
 286 chose the conjugate-gradient algorithm as it is deterministic (no residual variance like in SGD). The
 287 empirical Mean-Squared Error was computed over 100 random seeds, parallelized over 100 CPUs.
 288 The longest experiment took 7 wall-clock hours to complete. For the two-step estimators ("two-step"
 289 and "two-step (trig)"), a pre-estimate of the normalization was first computed using the geometric
 290 path with 10 distributions. Then, this estimate was used to re-parameterize an arithmetic path with
 291 10 distributions which produced the second estimate.

292 **Results** Figure 2 numerically supports the optimality of the NCE loss for a finite K (here, $K = 2$
 293 so three distributions are used) proven in Theorem 1. Figure 3 validates our main results for annealing
 294 paths. It shows how the estimation error scales with the proposal and target distributions growing
 295 apart, either with the parameter-distance in Figure 3a or with the dimensionality in Figure 3b.
 296 Using no annealing path ($K = 1$) produces an estimation error which grows linearly in log space;
 297 this numerically supports the exponential growth predicted by Theorem 2. Meanwhile, annealing
 298 ($K \rightarrow \infty$) sets the estimation error on different trends, depending on the choice of path. Choosing
 299 the geometric path brings the growth down to sub-exponential, as predicted by Theorem 3, while
 300 choosing the (basic) arithmetic path does not as in Theorem 4. To alleviate this, our two-step
 301 estimation methods consist in reparameterizing the arithmetic path so that it actually does bring down
 302 the estimation error. In fact, our two-step estimators in table 2 empirically approach the optimal
 303 estimation error in Figure 3. While this requires more computation, it has the appeal of making the
 304 estimation error *constant* with respect to the parameter-distance between the target and proposal
 305 distributions. Practically, this means that in Figure 3a, regular Noise-Contrastive Estimation (black,
 306 full line) fails when the parameter-distance between the target and proposal distributions is higher
 307 than 20, while our two-step estimators remain optimal.

308 We next explain interesting observations in Figure 3 which are actually coherent with our theory.
 309 First, in Figure 3a, the "two-step (trig)" estimator is only optimal when the parameter-distance
 310 between the target and proposal distributions is larger than 10. This is because the optimality of
 311 this two-step estimator was derived in Theorem 6 conditionally on non-overlapping distributions,
 312 here achieved by a large parameter-distance. Second, in both Figures 3a and 3b, the "two-step"
 313 estimator empirically achieves the optimal estimation error that was predicted for the "two-step (trig)"
 314 estimator. This suggests our polynomial upper bound from Theorem 5 may be loose in certain cases.
 315 This further explains why, in Figure 3a, the arithmetic path is near-optimal for a single value of the
 316 parameter-distance. At this value of 20, the partition function happens to be equal to one $Z(\theta_1) = 1$,
 317 so that the arithmetic path is effectively the same as the "two-step" estimator.

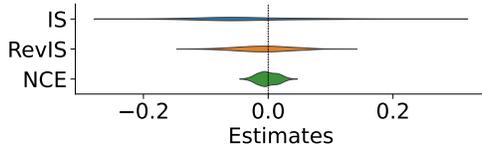


Figure 2: Optimality of the NCE loss, using the geometric path with $K = 2$.

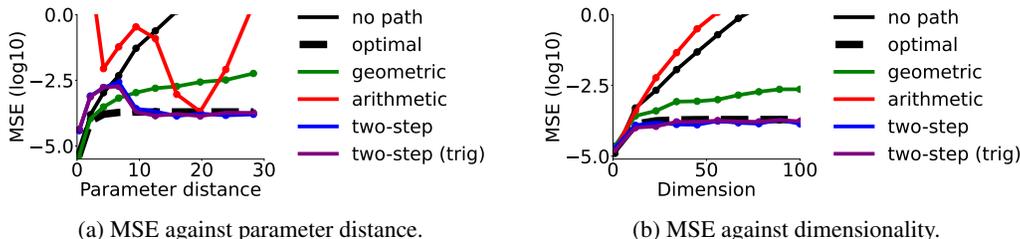


Figure 3: Estimation error as the target and proposal distributions grow apart. Without annealing, the error is exponential in the parameter distance (diagonal in log-scale). Annealing with the geometric path and our two-step methods brings down the error to slower growth, as predicted by our theorems.

318 6 Related work and Limitations

319 Previous work has mainly focused on annealed importance sampling [18, 35], which is a special case
 320 of our annealed Bregman estimator. They have evaluated the merits of different paths empirically,
 321 using an approximation of the estimation error called Effective Sample Size (ESS) and the consistency-
 322 gap. In our analysis, we consider consistent estimators and derive and optimize the exact estimation
 323 error of the optimal Noise-Contrastive Estimation. Liu et al. [27] considered the NCE estimate for Z
 324 (not $\log Z$) with the name “discriminance sampling”, and annealed the estimator using an extended
 325 state-space construction similar to Neal [10]. Their analysis of the estimation error is relevant but
 326 does not deal with hyperparameters other than the classification loss.

327 We made the common assumption of perfect sampling [18] in order to study the estimation error
 328 only and obtain practical guidelines to reduce it. We note however, that this leaves a gap to bridge
 329 with a practical setup where the sampling error cannot be ignored; in fact, annealed importance
 330 sampling [10] was originally proposed such that the samples can be obtained from a Markov Chain
 331 that has not converged. It might also be argued that the limit of almost no overlap between proposal
 332 and target, which we use a lot, is unrealistic. To see why it can be realistic, consider the case of
 333 natural image data. A typical proposal is Gaussian, since nothing much more sophisticated is tractable
 334 in high dimensions. However, there is almost no overlap between Gaussian data and natural images,
 335 which is seen in the fact that a human observer can effortlessly discriminate between the two.

336 7 Conclusion

337 We defined a class of estimators of the normalization constant, annealed Bregman estimation, which
 338 relies on a sampling phase from a path of distributions, and an estimation phase where these samples
 339 are used to estimate the log-normalization of the target distribution. Our results suggest a number of
 340 simple recommendations regarding hyperparameter choices in annealing. First, if the path has very
 341 few intermediate distributions, it is better to choose NCE due to its statistical optimality (Theorem
 342 1). If however, the path has many intermediate distributions and approaches the annealing limit, it
 343 is better to use IS due to its computational simplicity (since the statistical efficiency is the same
 344 as NCE). Annealing can always provide substantial benefits (Theorem 2). Moreover, if we have
 345 a reasonable a priori estimate of Z_1 , the arithmetic path achieves very low error (Theorem 5) —
 346 sometimes even approaching optimal (Theorem 6). On the other hand, even absent an initial estimate
 347 of Z_1 , the geometric path can exponentially reduce the estimation error compared with no annealing
 348 (Theorems 2 and 3).

349 **References**

- 350 [1] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and
351 Ben Poole. Score-based generative modeling through stochastic differential equations. *ArXiv*,
352 abs/2011.13456, 2021.
- 353 [2] A. Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of*
354 *Machine Learning Research*, 6(24):695–709, 2005.
- 355 [3] M. Gutmann and A. Hyvärinen. Noise-contrastive estimation of unnormalized statistical models,
356 with applications to natural image statistics. *Journal of Machine Learning Research*, 13(11):
357 307–361, 2012.
- 358 [4] G.E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural*
359 *Computation*, 14(8):1771–1800, 2002.
- 360 [5] R. Gao, E. Nijkamp, D.P. Kingma, Z. Xu, A.M. Dai, and Y. Nian Wu. Flow contrastive
361 estimation of energy-based models. *2020 IEEE/CVF Conference on Computer Vision and*
362 *Pattern Recognition (CVPR)*, pages 7515–7525, 2020.
- 363 [6] A.B. Owen. *Monte Carlo theory, methods and examples*. 2013.
- 364 [7] B. Liu, E. Rosenfeld, P. Ravikumar, and A. Risteski. Analyzing and improving the opti-
365 mization landscape of noise-contrastive estimation. In *International Conference on Learning*
366 *Representations (ICLR)*, 2022.
- 367 [8] H. Lee, C. Pabbaraju, A.P. Sevekari, and A. Risteski. Pitfalls of gaussians as a noise distribution
368 in NCE. In *International Conference on Learning Representations (ICLR)*, 2023.
- 369 [9] O. Chehab, A. Gramfort, and A. Hyvärinen. The optimal noise in noise-contrastive learning is
370 not what you think. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, volume 180,
371 pages 307–316. PMLR, 2022.
- 372 [10] R.M. Neal. Annealed importance sampling. *Statistics and Computing*, 11:125–139, 1998.
- 373 [11] C. Jarzynski. Nonequilibrium equality for free energy differences. *Physical Review Letters*, 78:
374 2690–2693, 1996.
- 375 [12] C. Jarzynski. Equilibrium free-energy differences from nonequilibrium measurements: A
376 master-equation approach. *Physical Review E*, 56:5018–5035, 1997.
- 377 [13] R. Salakhutdinov and G. Hinton. Replicated softmax: an undirected topic model. In *Neural*
378 *Information Processing Systems (NIPS)*, 2009.
- 379 [14] Y. Dauphin and Y. Bengio. Stochastic ratio matching of rbms for sparse high-dimensional
380 inputs. In *Neural Information Processing Systems (NIPS)*, volume 26. Curran Associates, Inc.,
381 2013.
- 382 [15] V. Masrani, T.A. Le, and F.D. Wood. The thermodynamic variational objective. In *Neural*
383 *Information Processing Systems (NeurIPS)*, 2019.
- 384 [16] R. Brekelmans, S. Huang, M. Ghassemi, G. Ver Steeg, R.B. Grosse, and A. Makhzani. Improv-
385 ing mutual information estimation with annealed and energy-based bounds. In *International*
386 *Conference on Learning Representations (ICLR)*, 2022.
- 387 [17] A. Gelman and X.-L. Meng. Simulating normalizing constants: From importance sampling to
388 bridge sampling to path sampling. *Statistical Science*, 13:163–185, 1998.
- 389 [18] R. Grosse, C. Maddison, and R. Salakhutdinov. Annealing between distributions by averaging
390 moments. In *Advances in Neural Information Processing Systems (NIPS)*, volume 26. Curran
391 Associates, Inc., 2013.
- 392 [19] G.M. Torrie and J.P. Valleau. Nonphysical sampling distributions in monte carlo free-energy
393 estimation: Umbrella sampling. *Journal of Computational Physics*, 23(2):187–199, 1977.

- 394 [20] X.-L. Meng and W.H. Wong. Simulating ratios of normalizing constants via a simple identity: a
395 theoretical exploration. *Statistica Sinica*, pages 831–860, 1996.
- 396 [21] V. Masrani, R. Brekelmans, T. Bui, F. Nielsen, A. Galstyan, G. Ver Steeg, and F. Wood. q-paths:
397 Generalizing the geometric annealing path using power means. In *Conference on Uncertainty*
398 *in Artificial Intelligence (UAI)*, volume 161, pages 1938–1947. PMLR, 27–30 Jul 2021.
- 399 [22] O. Chehab, A. Gramfort, and A. Hyvärinen. Optimizing the noise in self-supervised learning:
400 from importance sampling to noise-contrastive estimation, 2023.
- 401 [23] M.A. Newton. Approximate bayesian-inference with the weighted likelihood bootstrap. *Journal*
402 *of the royal statistical society series b-methodological*, 1994.
- 403 [24] A.W. van der Vaart. *Asymptotic Statistics*. Asymptotic Statistics. Cambridge University Press,
404 2000.
- 405 [25] M.-H. Chen and Q.-M. Shao. On monte carlo methods for estimating ratios of normalizing
406 constants. *The Annals of Statistics*, 25(4):1563–1594, 1997.
- 407 [26] M. Uehara, T. Matsuda, and F. Komaki. Analysis of noise contrastive estimation from the
408 perspective of asymptotic variance. *ArXiv*, 2018. doi: 10.48550/ARXIV.1808.07983.
- 409 [27] Q. Liu, J. Peng, A.T. Ihler, and J.W. Fisher III. Estimating the partition function by discriminace
410 sampling. In *Uncertainty in Artificial Intelligence (UAI)*, 2015.
- 411 [28] A. Hyvärinen, J. Hurri, and P.O. Hoyer. Natural image statistics - a probabilistic approach to
412 early computational vision. In *Computational Imaging and Vision*, 2009.
- 413 [29] B. Sriperumbudur, K. Fukumizu, A. Gretton, A. Hyvärinen, and R. Kumar. Density estimation
414 in infinite dimensional exponential families. *Journal of Machine Learning Research*, 18(57):
415 1–59, 2017.
- 416 [30] F. Nielsen and V. Garcia. Statistical exponential families: A digest with flash cards, 2011.
- 417 [31] X.-L. Meng and S. Schilling. Fitting full-information item factor models and an empirical
418 investigation of bridge sampling. *Journal of the American Statistical Association*, 91(435):
419 1254–1267, 1996.
- 420 [32] M. Uehara, T. Kanamori, T. Takenouchi, and T. Matsuda. A unified statistically efficient
421 estimation framework for unnormalized models. In *International Conference on Artificial*
422 *Intelligence and Statistics (AISTATS)*, volume 108, pages 809–819. PMLR, 2020.
- 423 [33] L. Ellam, H. Strathmann, M.A. Girolami, and I. Murray. A determinant-free method to simulate
424 the parameters of large gaussian fields. *Stat*, 6:271–281, 2017.
- 425 [34] P. Virtanen, R. Gommers, T.E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski,
426 P. Peterson, W. Weckesser, J. Bright, S.J. van der Walt, M. Brett, J. Wilson, J.K. Millman,
427 N. Mayorov, A.R.J. Nelson, E. Jones, R. Kern, E. Larson, C.J. Carey, I. Polat, Y. Feng, E.W.
428 Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E.A. Quintero, C.R.
429 Harris, A.M. Archibald, A.H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors.
430 SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:
431 261–272, 2020.
- 432 [35] T. Kiwaki. Variational optimization of annealing schedules, 2015.
- 433 [36] M. Pihlaja, M. Gutmann, and A. Hyvärinen. A family of computationally efficient and simple
434 estimators for unnormalized statistical models. In *Uncertainty in Artificial Intelligence (UAI)*,
435 2010.
- 436 [37] L. Wang, D.E. Jones, and X.-L. Meng. Warp bridge sampling: The next generation. *Journal of*
437 *the American Statistical Association*, 117(538):835–851, 2022.
- 438 [38] H. Xing. Improving bridge estimators via f-GAN. *Statistics and Computing*, 32, 2022.
- 439 [39] Y. Polyanskiy and Y. Wu. *Information Theory: From Coding to Learning*. Cambridge University
440 Press, 2022.
- 441 [40] S.-I. Amari and H. Nagaoka. *Methods of information geometry*. AMS, 2000.

442 In the following, we will study the estimation error of of annealed Bregman estimation (ABE) in two
 443 importants setups: the log-normalization is computed using two distributions ($K = 1$), the proposal
 444 and the target, or else using a path of distributions ($K \rightarrow \infty$).

445 The anonymized code used for the experiments is available at <https://anonymous.4open.science/r/annealed-nce-submission-B8A0>.

447 A No annealing, $K = 1$

448 We use [22, Eq.21] for the estimation error of any suitably parameterized ⁴ classifier $F(\mathbf{x}; \beta)$ between
 449 two distributions p_1 and p_0 . The estimation error is measured by the Mean-Squared Error (MSE)

$$\text{MSE}_{\hat{\beta}}(p_n, \nu, \phi, N) = \frac{\nu + 1}{N} \text{tr}(\Sigma) \quad (18)$$

450 which depends on the sample sizes $N = N_1 + N_0$, their ratio $\nu = N_1/N_0$, the Bregman classification
 451 loss indexed by the convex function $\phi(x)$, and the asymptotic variance matrix

$$\Sigma = \mathbf{I}_w^{-1} \left(\mathbf{I}_v - \left(1 + \frac{1}{\nu}\right) \mathbf{m}_w \mathbf{m}_w^\top \right) \mathbf{I}_w^{-1} . \quad (19)$$

452 Here, $\mathbf{m}_w(\beta^*)$, $\mathbf{I}_w(\beta^*)$ and $\mathbf{I}_v(\beta^*)$ are the reweighted mean and covariances of the paramete-
 453 gradient of the classifier, also known as the ‘‘relative’’ Fisher score $\nabla_{\beta} F(\mathbf{x}; \beta^*)$,

$$\mathbf{m}_w(\beta^*) = \mathbb{E}_{\mathbf{x} \sim p_d} [w(\mathbf{x}) \nabla_{\beta} F(\mathbf{x}; \beta^*)] \quad (20)$$

$$\mathbf{I}_w(\beta^*) = \mathbb{E}_{\mathbf{x} \sim p_d} [w(\mathbf{x}) \nabla_{\beta} F(\mathbf{x}; \beta^*) \nabla_{\beta} F(\mathbf{x}; \beta^*)^\top] \quad (21)$$

$$\mathbf{I}_v(\beta^*) = \mathbb{E}_{\mathbf{x} \sim p_d} [v(\mathbf{x}) \nabla_{\beta} F(\mathbf{x}; \beta^*) \nabla_{\beta} F(\mathbf{x}; \beta^*)^\top] \quad (22)$$

454 where the reweighting of data points is by $w(\mathbf{x}) := \frac{p_1(\mathbf{x}) \phi''(\frac{p_1}{\nu p_0}(\mathbf{x}))}{\nu p_0(\mathbf{x})}$ and by $v(\mathbf{x}) =$
 455 $w(\mathbf{x})^2 \frac{\nu p_0(\mathbf{x}) + p_1(\mathbf{x})}{\nu p_0(\mathbf{x})}$, which are all evaluated at the true parameter value β^* .

456 **Scalar parameterization** We now consider a specific parameterization of the classifier:

$$F(\mathbf{x}; \beta) = \log \left(\frac{f_1(\mathbf{x})}{\nu f_0(\mathbf{x})} \right) - \beta \quad (23)$$

457 where the optimal parameter is the log-ratio of normalizations $\beta^* = \log(Z_1/Z_0)$. Consequently, we
 458 have $\nabla_{\beta} F(\mathbf{x}; \beta^*) = -1$ and plugging this into the above quantities yields

$$\text{MSE} = \frac{1 + \nu}{T} \left(\frac{\mathbb{E}_{\mathbf{x} \sim p_1} [w^2(\mathbf{x}) \frac{\nu p_0(\mathbf{x}) + p_1(\mathbf{x})}{\nu p_0(\mathbf{x})}]}{\mathbb{E}_{\mathbf{x} \sim p_1} [w(\mathbf{x})]^2} - \left(1 + \frac{1}{\nu}\right) \right)$$

459 which matches the formula found in [20, Eq 3.2]. For different choices of the Bregman classification
 460 loss, the estimation error is written using a divergence between the two distributions

Name	Loss identified by $\phi(x)$	Estimator	MSE
IS	$x \log x$	$\log \mathbb{E}_{p_0} \frac{f_1}{f_0}$	$\frac{1+\nu}{\nu N} \mathcal{D}_{\chi^2}(p_1, p_0)$
461 RevIS	$-\log x$	$-\log \mathbb{E}_{p_1} \frac{f_0}{f_1}$	$\frac{1+\nu}{N} \mathcal{D}_{\chi^2}(p_0, p_1)$
NCE	$x \log x - (1+x) \log(\frac{1+x}{2})$	implicit	$\frac{(1+\nu)^2}{\nu N} \frac{\mathcal{D}_{\text{HM}}(p_1, p_0)}{1 - \mathcal{D}_{\text{HM}}(p_1, p_0)}$
IS-RevIS	$(1 - \sqrt{x})^2$	$\log \mathbb{E}_{p_0} \frac{f_1}{f_0} - \log \mathbb{E}_{p_1} \frac{f_0}{f_1}$	$\frac{(1+\nu)^2}{\nu N} \frac{1 - (1 - \mathcal{D}_{H^2}(p_d, p_n))^2}{(1 - \mathcal{D}_{H^2}(p_d, p_n))^2}$

462 where

463 $\mathcal{D}_{\chi^2}(p_1, p_0) := \left(\int \frac{p_1^2}{p_0} \right) - 1$ is the chi-squared divergence

⁴technically, the formula was derived in [36, 3] assuming the classifier was parameterized as $F(\mathbf{x}; \beta) = \log p_1(\mathbf{x}; \beta) / \nu p_0(\mathbf{x})$ but the proof seems to generalize to any well-defined parameterization $F(\mathbf{x}; \beta)$.

464 $D_{H^2}(p_1, p_0) := 1 - \left(\int \sqrt{p_1 p_0} \right) \in [0, 1]$ is the squared Hellinger distance
465 $\mathcal{D}_{\text{HM}}(p_1, p_0) := 1 - \int (\pi p_1^{-1} + (1 - \pi) p_0^{-1})^{-1} = 1 - \frac{1}{\pi} \mathbb{E}_{p_1} \frac{\pi p_0}{(1 - \pi) p_1 + \pi p_0} \in [0, 1]$
466 is the harmonic divergence with weight $\pi \in [0, 1]$.
467 Here, the weight $\pi = P(Y = 0) = \frac{T_0}{T} = \frac{\nu}{1 + \nu}$.

468 **Proof of Theorem 2** *Exponential error of binary NCE*

469 The estimation error of binary NCE is expressed in terms of the harmonic divergence

$$\text{MSE} = \frac{4}{N} \frac{\mathcal{D}_{\text{HM}}(p_1, p_0)}{1 - \mathcal{D}_{\text{HM}}(p_1, p_0)} \quad (24)$$

470 which is intractable for general exponential families. Instead, we can lower-bound the estimation
471 error. To do so, we lower-bound the harmonic divergence using the inequality of means (harmonic vs.
472 geometric)

$$\mathcal{D}_{\text{HM}}(p_1, p_0) = 1 - \int \frac{2p_0 p_1}{p_0 + p_1} \geq 1 - \int \sqrt{p_0 p_1} = \mathcal{D}_{H^2}(p_0, p_1) \quad (25)$$

473 and therefore

$$\text{MSE}_{\text{LB}} = \frac{4}{N} \frac{\mathcal{D}_{H^2}(p_1, p_0)}{1 - \mathcal{D}_{H^2}(p_1, p_0)} . \quad (26)$$

474 This lower bound is expressed in terms of the squared Hellinger distance, that is tractable for
475 exponential families:

$$\mathcal{D}_{H^2}(p_1, p_0) := 1 - \int_{\mathbf{x} \in \mathbb{R}^D} \sqrt{p_1 p_0} d\mathbf{x} \quad (27)$$

$$= 1 - \int_{\mathbf{x} \in \mathbb{R}^D} \frac{1}{Z(\boldsymbol{\theta}_1)^{\frac{1}{2}} Z(\boldsymbol{\theta}_0)^{\frac{1}{2}}} \exp\left(\frac{1}{2}(\boldsymbol{\theta}_1 + \boldsymbol{\theta}_0)^\top \mathbf{t}(\mathbf{x})\right) d\mathbf{x} \quad (28)$$

$$= 1 - \frac{Z(\frac{1}{2}\boldsymbol{\theta}_1 + \frac{1}{2}\boldsymbol{\theta}_0)}{Z(\boldsymbol{\theta}_1)^{\frac{1}{2}} Z(\boldsymbol{\theta}_0)^{\frac{1}{2}}} \quad (29)$$

$$= 1 - \exp\left(\log Z\left(\frac{1}{2}\boldsymbol{\theta}_1 + \frac{1}{2}\boldsymbol{\theta}_0\right) - \frac{1}{2}\log Z(\boldsymbol{\theta}_1) - \frac{1}{2}\log Z(\boldsymbol{\theta}_0)\right) . \quad (30)$$

476 We now wish to lower bound MSE_{LB} , and therefore $\mathcal{D}_{H^2}(p_1, p_0)$, by an expression which is expo-
477 nential in the parameter distance $\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\|$. To do so, we note that for exponential families, the
478 log-normalization is convex in the parameters. Here, we further assume strong convexity, so that

$$\log Z\left(\frac{1}{2}\boldsymbol{\theta}_1 + \frac{1}{2}\boldsymbol{\theta}_0\right) \leq \frac{1}{2}\log Z(\boldsymbol{\theta}_1) + \frac{1}{2}\log Z(\boldsymbol{\theta}_0) - \frac{1}{8}M\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\|^2 \quad (31)$$

479 where M is the strong convexity constant. Plugging this back into the squared Hellinger distance, we
480 obtain

$$\mathcal{D}_{H^2}(p_1, p_0) \geq 1 - \exp\left(-\frac{1}{8}M\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\|^2\right) \quad (32)$$

481 so that the MSE

$$\text{MSE} \geq \frac{4}{N} \frac{\mathcal{D}_{H^2}(p_1, p_0)}{1 - \mathcal{D}_{H^2}(p_1, p_0)} \geq \frac{4}{N} \exp\left(\frac{1}{8}M\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\|^2\right) - 1 \quad (33)$$

482 grows exponentially with the euclidean distance between the parameters.

483 **B Annealing limit, $K \rightarrow \infty$**

484 We now consider annealing paths $(p_t)_{t \in [0,1]}$ that interpolate between between the proposal p_0 and
 485 the target p_1 .

486 **B.1 Estimation error**

487 We first show the optimality of the NCE loss within the family of Annealed Bregman Estimators, in
 488 the sense that it produces the smallest estimation error. We then study the estimation error of different
 489 annealed Bregman estimators in the annealing limit of a continuous path ($K \rightarrow \infty$).

490 **Proof of Theorem 1** *Optimality of the NCE loss and the estimation error in the annealing limit*
 491 $K \rightarrow \infty$

- Optimality of the NCE loss

492 Because the annealed Bregman estimator is built by adding independent estimators

$$\widehat{\log Z_1} = \sum_{k=0}^{K-1} \log \left(\frac{\widehat{Z_{(k+1)/K}}}{Z_{k/K}} \right) + \log Z_0 . \quad (34)$$

494 the total Mean Squared Error (MSE) is the sum of each MSEs for each estimator (indexed
 495 by $k \in \llbracket 0, K-1 \rrbracket$)

$$\text{MSE}((\phi_k)_{k \in \llbracket 0, K \rrbracket}) = \sum_{k=0}^{K-1} \text{MSE}_k(\phi_k) \quad (35)$$

496 where we highlighted the dependency on the classification losses identified by $(\phi_k)_{k \in \llbracket 0, K \rrbracket}$.
 497 The MSEs follow Eq. 24. It was shown by Meng and Wong [20] that for any of these MSEs,
 498 the optimal loss is identified by $\phi_k(x) = x \log x - (1+x) \log(\frac{1+x}{2})$ and is in fact the NCE
 499 loss [22]. Thus the sum of MSEs is minimized for the same loss.

- Annealed Noise-Contrastive Estimation (NCE)

501 We are interested in the estimation error (asymptotic MSE) obtained for the NCE loss. Based
 502 off table 1, it is written as

$$\text{MSE} = \frac{4K}{N} \sum_{k=0}^{K-1} \frac{\mathcal{D}_{\text{HM}}(p_{k/K}, p_{(k+1)/K})}{1 - \mathcal{D}_{\text{HM}}(p_{k/K}, p_{(k+1)/K})} . \quad (36)$$

503 The estimation error of balanced ($\nu = 1$) NCE-JS between two ‘‘close’’ distributions p_t and
 504 p_{t+h} , is

$$\text{MSE}(p_t, p_{t+h}) \propto \frac{\mathcal{D}_{\text{HM}}(p_t, p_{t+h})}{1 - \mathcal{D}_{\text{HM}}(p_t, p_{t+h})} \quad (37)$$

505 The estimation error can be simplified using a Taylor expansion. To do so, we recall that
 506 \mathcal{D}_{HM} is an f-divergence generated by $\phi(x) = 1 - \frac{x}{\pi + (1-\pi)x}$ [37, 38] ($\pi = \frac{1}{2}$ here) and its
 507 expansion is therefore [39, Eq.7.64]

$$\mathcal{D}_{\text{HM}}(p_t, p_{t+h}) = \frac{1}{2} h^2 \nabla_t^2 \mathcal{D}_{\text{HM}}(p_t, p_{t+h}) + o(h^2) \quad (38)$$

$$= \frac{1}{2} \phi''(1) h^2 I(t) + o(h^2) = \frac{1}{4} h^2 I(t) + o(h^2) . \quad (39)$$

508 It follows that

$$\frac{\mathcal{D}_{\text{HM}}(p_t, p_{t+h})}{1 - \mathcal{D}_{\text{HM}}(p_t, p_{t+h})} = \frac{1}{4} I(t) h^2 + o(h^2) . \quad (40)$$

509 Summing these estimation errors along the path of distributions with $h = 1/K$,

$$\text{MSE} = \frac{4K}{N} \sum_{k=0}^{K-1} \left(\frac{1}{4} I(t) \frac{1}{K^2} + o\left(\frac{1}{K^2}\right) \right) \quad (41)$$

$$= \left(\frac{1}{NK} \sum_{k=0}^{K-1} I(t) \right) + o(1) \underset{K \rightarrow \infty}{\sim} \frac{1}{N} \int_0^1 I(t) dt . \quad (42)$$

510
511

In the case of a parametric path $p(\mathbf{x}|\boldsymbol{\theta}(t))_{t \in [0,1]}$, the proof is the same. Simply, the second-order term in the Taylor expansion of Eq. 39 is computed using the chain rule

$$\nabla_t^2 \text{MSE}(p_{\boldsymbol{\theta}(t)}, p_{\boldsymbol{\theta}(t+h)}) \quad (43)$$

$$= \dot{\boldsymbol{\theta}}(t)^\top \nabla_{\boldsymbol{\theta}}^2 \text{MSE}(p_{\boldsymbol{\theta}(t)}, p_{\boldsymbol{\theta}(t+h)}) \dot{\boldsymbol{\theta}}(t) + \ddot{\boldsymbol{\theta}}(t)^\top \nabla_{\boldsymbol{\theta}} \text{MSE}(p_{\boldsymbol{\theta}(t)}, p_{\boldsymbol{\theta}(t+h)}) \quad (44)$$

$$= \dot{\boldsymbol{\theta}}(t)^\top \nabla_{\boldsymbol{\theta}}^2 \text{MSE}(p_{\boldsymbol{\theta}(t)}, p_{\boldsymbol{\theta}(t+h)}) \dot{\boldsymbol{\theta}}(t)^\top + 0 \quad (45)$$

$$= \dot{\boldsymbol{\theta}}(t)^\top \mathbf{I}(\boldsymbol{\theta}(t)) \dot{\boldsymbol{\theta}}(t) \quad (46)$$

512
513

- Annealed importance sampling (IS)

Similarly, for the choice of the importance sampling base estimator,

$$\text{MSE} = \frac{2K}{N} \sum_{k=0}^{K-1} \mathcal{D}_{\chi^2}(p_{(k+1)/K}, p_{k/K}) \quad (47)$$

$$= \frac{2K}{N} \sum_{k=0}^{K-1} \mathcal{D}_{\text{rev}\chi^2}(p_{k/K}, p_{(k+1)/K}) = \frac{2K}{N} \sum_{k=0}^{K-1} \left(\frac{1}{2} \phi''(1) I(t) \frac{1}{K^2} + o\left(\frac{1}{K^2}\right) \right) \quad (48)$$

$$= \left(\frac{1}{NK} \sum_{k=0}^{K-1} I(t) \right) + o(1) \underset{K \rightarrow \infty}{\sim} \frac{1}{N} \int_0^1 I(t) dt . \quad (49)$$

514
515
516

given that $\phi(x) = -\log(x)$ and therefore $\phi''(1) = 1$ for the reverse χ^2 divergence.

- Annealed reverse importance sampling (RevIS)

Similarly, for the choice of the reverse importance sampling base estimator,

$$\text{MSE} = \frac{2K}{N} \sum_{k=0}^{K-1} \mathcal{D}_{\chi^2}(p_{k/K}, p_{(k+1)/K}) \quad (50)$$

$$= \frac{2K}{N} \sum_{k=0}^{K-1} \left(\frac{1}{2} \phi''(1) I(t) \frac{1}{K^2} + o\left(\frac{1}{K^2}\right) \right) \quad (51)$$

$$= \left(\frac{1}{NK} \sum_{k=0}^{K-1} I(t) \right) + o(1) \underset{K \rightarrow \infty}{\sim} \frac{1}{N} \int_0^1 I(t) dt . \quad (52)$$

517

given that $\phi(x) = x \log(x)$ and therefore $\phi''(1) = 1$ for the χ^2 divergence.

518 B.2 Examples of paths

519 **Geometric path** The geometric path is defined in the space of unnormalized densities by

$$f_t(\mathbf{x}) := p_0(\mathbf{x})^{1-t} f_1(\mathbf{x})^t = p_0(\mathbf{x})^{1-t} p_1(\mathbf{x})^t Z_1^t \propto p_0(\mathbf{x})^{1-t} p_1(\mathbf{x})^t \quad (53)$$

520 so in the space of normalized densities, the path is

$$p_t := \frac{p_0(\mathbf{x})^{1-t} p_1(\mathbf{x})^t}{Z_t} \quad (54)$$

521 where the normalization is

$$Z_t := \int_{\mathbf{x} \in \mathbb{R}^d} p_0(\mathbf{x})^{1-t} p_1(\mathbf{x})^t d\mathbf{x} = \mathbb{E}_{\mathbf{x} \sim p_1} \left[\left(\frac{p_0(\mathbf{x})}{p_1(\mathbf{x})} \right)^t \right] = \mathbb{E}_{\mathbf{x} \sim p_0} \left[\left(\frac{p_1(\mathbf{x})}{p_0(\mathbf{x})} \right)^{1-t} \right] . \quad (55)$$

522 **Arithmetic path** The arithmetic path is defined in the space of unnormalized densities by

$$f_t(\mathbf{x}) := (1-t)p_0(\mathbf{x}) + t f_1(\mathbf{x}) = (1-t)p_0 + t Z_1 p_1 \quad (56)$$

$$\propto \frac{(1-t)}{(1-t) + t Z_1} p_0 + \frac{t Z_1}{(1-t) + t Z_1} p_1 \quad (57)$$

523 so in the space of normalized densities, the path is actually a mixture between the target and the
524 proposal, where the weight of the mixture is a nonlinear function of the target normalization

$$p_t := (1 - \tilde{w}_t) p_0 + \tilde{w}_t p_1, \quad \tilde{w}_t = \frac{t Z_1}{(1-t) + t Z_1} . \quad (58)$$

525 **Optimal path** We know (e.g. from Gelman and Meng [17, Eq. 49]) that the optimal path is

$$p_t(\mathbf{x}) = (a(t)\sqrt{p_0(\mathbf{x})} + b(t)\sqrt{p_1(\mathbf{x})})^2 \quad (59)$$

526 where the coefficients $a(t)$ and $b(t)$

$$a(t) = \frac{\cos((2t-1)\alpha_H)}{2\cos(\alpha_H)} - \frac{\sin((2t-1)\alpha_H)}{2\sin(\alpha_H)} \quad (60)$$

$$b(t) = \frac{\cos((2t-1)\alpha_H)}{2\cos(\alpha_H)} + \frac{\sin((2t-1)\alpha_H)}{2\sin(\alpha_H)} \quad (61)$$

527 are simple functions of the squared Hellinger distance $\mathcal{D}_{H^2}(p_0, p_1)$ between the proposal and the
528 target⁵

$$\alpha_H = \arctan\left(\sqrt{\frac{\mathcal{D}_{H^2}(p_0, p_1)}{2 - \mathcal{D}_{H^2}(p_0, p_1)}}\right) \in [0, \frac{\pi}{4}] . \quad (62)$$

529 The estimation error produced by that optimal path is [17, Eq. 48]

$$\text{MSE} = \frac{1}{N} \int_0^1 I(t) dt = \frac{1}{N} 16\alpha_H^2 . \quad (63)$$

530 For two gaussians

$$p_0 := \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \quad (64)$$

$$p_1 := \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \quad (65)$$

531 the squared Hellinger distance can be written in closed-form

$$\mathcal{D}_{H^2}(p_0, p_1) = 1 - \frac{|\boldsymbol{\Sigma}_0|^{\frac{1}{4}} |\boldsymbol{\Sigma}_1|^{\frac{1}{4}}}{|\frac{1}{2}\boldsymbol{\Sigma}_0 + \frac{1}{2}\boldsymbol{\Sigma}_1|^{\frac{1}{2}}} \exp\left(-\frac{1}{8}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \left(\frac{1}{2}\boldsymbol{\Sigma}_0 + \frac{1}{2}\boldsymbol{\Sigma}_1\right)^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)\right) \quad (66)$$

532 and plugs into the optimal path formula, which is also obtained in closed-form.

533 B.3 Estimation error from taking different paths

534 **Proof of Theorem 3** *Polynomial error of annealed NCE with the geometric path*

535 We next study the estimation error produced by the geometric path (Figure 1). In the annealing limit
536 $K \rightarrow \infty$, the MSE is written as

$$\text{MSE} = \frac{1}{N} \int_0^1 I(t) dt . \quad (67)$$

537 We recall from Grosse et al. [18] that the geometric path is closed for distributions in the exponential
538 family: all distributions along the path remain in the exponential family. Furthermore, their Fisher
539 information can be written in terms of the terms parameters [18, Eq. 17]; this is based off a result
540 of exponential families from [40, Section 3.3]

$$I(t) = \dot{\boldsymbol{\theta}}(t)^\top \mathbf{I}(\boldsymbol{\theta}(t)) \dot{\boldsymbol{\theta}}(t) = \dot{\boldsymbol{\theta}}(t)^\top \dot{\boldsymbol{\mu}}(t) \quad (68)$$

541 where $\boldsymbol{\mu}(t)$ are the generalized moments, defined as $\boldsymbol{\mu}(t) = \mathbb{E}_{\mathbf{x} \sim p_t(\mathbf{x})}[\mathbf{t}(\mathbf{x})] = \nabla_{\boldsymbol{\theta}} \log Z_t(\boldsymbol{\theta})$. It
542 follows,

$$\text{MSE} = \frac{1}{N} \int_0^1 \dot{\boldsymbol{\theta}}(t)^\top \dot{\boldsymbol{\mu}}(t) dt . \quad (69)$$

543 The geometric path is defined in parameter space by $\boldsymbol{\theta}_t = t\boldsymbol{\theta}_1 + (1-t)\boldsymbol{\theta}_0$, therefore

$$\text{MSE} = \frac{1}{N} (\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0)^\top \int_0^1 \dot{\boldsymbol{\mu}}(t) dt = \frac{1}{N} (\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0)^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) \quad (70)$$

⁵In Gelman and Meng [17, Eq. 49], the Hellinger distance is defined such that it is in $[0, \sqrt{2}]$. We here instead use the conventional definition of the squared Hellinger distance which is normalized so that it is in $[0, 1]$.

544 as in [18, Eq. 17]. For exponential families, $\log Z(\boldsymbol{\theta})$ is convex in $\boldsymbol{\theta}$. Here, we further assume strong
 545 convexity (with constant M) and smoothness (with constant L) so that

$$(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0)^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) = (\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0)^\top (\nabla_{\boldsymbol{\theta}} \log Z_t(\boldsymbol{\theta}_1) - \nabla_{\boldsymbol{\theta}} \log Z_t(\boldsymbol{\theta}_0)) \quad (71)$$

$$\leq \frac{1}{M} \|\nabla_{\boldsymbol{\theta}} \log Z_t(\boldsymbol{\theta}_1) - \nabla_{\boldsymbol{\theta}} \log Z_t(\boldsymbol{\theta}_0)\|^2 \leq \frac{L^2}{M} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\|^2 \quad (72)$$

546 so that the MSE

$$\text{MSE} \leq \frac{L^2}{MN} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\|^2 \quad (73)$$

547 is polynomial in the euclidean distance between the parameters.

548 **Proof of Theorem 4** *Exponential error of annealed NCE with the arithmetic path and "vanilla"*
 549 *schedule*

550 We now study the estimation error produced by the arithmetic path with "vanilla" schedule (table 2,
 551 line 3). Similarly, we start with the formula of the estimation error of NCE in the limit of a continuous
 552 path

$$\text{MSE} = \frac{1}{N} \int_0^1 I(t) dt \quad (74)$$

553 where $I(t) = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}, t)} [(\frac{d}{dt} \log p(\mathbf{x}, t))^2]$ is the Fisher information over the path, using time t as the
 554 parameter. The arithmetic path is a gaussian mixture (see table 2) so we will conveniently use the
 555 parametric form of the path to compute the Fisher information

$$I(t) = \dot{\tilde{w}}_t^\top I(\tilde{w}_t) \dot{\tilde{w}}_t \quad (75)$$

556 where the parameter here is the weight of the Gaussian mixture $\tilde{w}_t = tZ_1/(tZ_1 + 1 - t)$. We will
 557 need to compute two quantities: the Fisher information to that mixture parameter (not the time), and
 558 the parameter speed $\dot{\tilde{w}}_t$.

$$I(\tilde{w}_t) := \mathbb{E}_{\mathbf{x} \sim p_{\tilde{w}_t}} \left[\left(\frac{\partial \log p_{\tilde{w}_t}(\mathbf{x})}{\partial \tilde{w}_t} \right)^2 \right] = \mathbb{E}_{\mathbf{x} \sim p_{\tilde{w}_t}} \left[\left(\frac{1}{p_{\tilde{w}_t}(\mathbf{x})} \frac{\partial p_{\tilde{w}_t}(\mathbf{x})}{\partial \tilde{w}_t} \right)^2 \right] \quad (76)$$

$$= \int_{\mathbf{x} \in \mathbb{R}^D} \frac{(p_1(\mathbf{x}) - p_0(\mathbf{x}))^2}{p_{\tilde{w}_t}(\mathbf{x})} d\mathbf{x} = \int_{\mathbf{x} \in \mathbb{R}^D} \frac{(p_1(\mathbf{x}) - p_0(\mathbf{x}))^2}{(1 - \tilde{w}_t)p_0(\mathbf{x}) + \tilde{w}_t p_1(\mathbf{x})} d\mathbf{x} \quad (77)$$

$$\geq \int_{\mathbf{x} \in \mathbb{R}^D} \frac{(p_1(\mathbf{x}) - p_0(\mathbf{x}))^2}{p_0(\mathbf{x}) + p_1(\mathbf{x})} d\mathbf{x} = \int p_0(\mathbf{x}) \frac{\left(1 - \frac{p_1(\mathbf{x})}{p_0(\mathbf{x})}\right)^2}{1 + \frac{p_1(\mathbf{x})}{p_0(\mathbf{x})}} = \mathcal{D}_\phi(p_1, p_0) \quad (78)$$

559 which is an f-divergence with generator $\phi(x) = (1 - x)^2/(1 + x)$ that provides a t -independent
 560 lower bound. This will allow us to factor this quantity out of the integral defining the MSE, and
 561 simplify computations. We also have

$$\dot{\tilde{w}}_t := \frac{\partial}{\partial t} \tilde{w}_t = \frac{1}{t(1-t)} \times \sigma \left(\log \frac{tZ_1}{1-t} \right) \times \left(1 - \sigma \left(\log \frac{tZ_1}{1-t} \right) \right) \quad (79)$$

$$= \frac{1}{t(1-t)} \times \frac{tZ_1}{(1-t) + tZ_1} \times \frac{(1-t)}{(1-t) + tZ_1} = \frac{Z_1}{((1-t) + tZ_1)^2} \quad (80)$$

562 where we choose to keep the dependency on t . The intuition is that integrating this quantity will yield
 563 a function of Z_1 , which will drive the MSE toward high values. We next show this rigorously and
 564 finally compute the estimation error.

$$\text{MSE} = \frac{1}{N} \int_0^1 I(t) dt = \frac{1}{N} \int_0^1 \dot{\tilde{w}}(t) I(\tilde{w}(t)) \dot{\tilde{w}}(t) dt \quad (81)$$

$$\geq \frac{1}{N} \times \mathcal{D}_\phi(p_1, p_0) \times \int_0^1 \dot{\tilde{w}}(t)^2 dt = \frac{1}{N} \times \mathcal{D}_\phi(p_1, p_0) \times Z_1^2 \times \int_0^1 \frac{1}{(t(Z_1 - 1) + 1)^4} dt \quad (82)$$

$$= \frac{1}{N} \times \mathcal{D}_\phi(p_1, p_0) \times Z_1^2 \times \frac{Z_1^2 + Z_1 + 1}{3Z_1^3} = \frac{1}{3N} \times \mathcal{D}_\phi(p_1, p_0) \times (Z_1^{-1} + 1 + Z_1) \quad (83)$$

565 We would like to write Z_1 in terms of the parameters. To do so, we now suppose the unnormalized
 566 target is in a simply unnormalized exponential family. Consequently,

$$Z_1 := \exp(\log Z(\boldsymbol{\theta}_1) - \log Z(\boldsymbol{\theta}_0) + \log Z(\boldsymbol{\theta}_0)) \quad (84)$$

$$\geq \exp\left(\nabla \log Z(\boldsymbol{\theta}_0)(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0) + \frac{M}{2} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\|^2 + \log Z(\boldsymbol{\theta}_0)\right) . \quad (85)$$

567 using the strong convexity of the log-partition function. It follows that in the limit of parameter-
 568 distance $\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\| \rightarrow \infty$, the MSE grows (at least) exponentially with the parameter-distance

$$\text{MSE} = O\left(\frac{1}{3N} \times \mathcal{D}_\phi(p_1, p_0) \times \exp\left(\frac{M}{2} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\|^2\right)\right), \quad \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\|^2 \rightarrow \infty \quad (86)$$

569 **Proof of Theorem 5** *Polynomial error of annealed NCE with the arithmetic path and "oracle"*
 570 *schedule*

571 We now study the estimation error produced by the arithmetic path with "oracle" schedule (table 2,
 572 line 4). Similarly, we start with the formula of the estimation error of NCE annealed over a continuous
 573 path

$$\text{MSE} = \frac{1}{N} \int_0^1 I(t) dt . \quad (87)$$

574 where $I(t) = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}, t)} \left[\left(\frac{d}{dt} \log p(\mathbf{x}, t) \right)^2 \right]$ is the Fisher information over the path, using time t as the
 575 parameter. The arithmetic path is the gaussian mixture $p_t(\mathbf{x}) = tp_1(\mathbf{x}) + (1-t)p_0(\mathbf{x})$ (see table 2).
 576 The Fisher information is therefore

$$I(t) := \mathbb{E}_{\mathbf{x} \sim p_t} \left[\left(\frac{\partial \log p_t(\mathbf{x})}{\partial t} \right)^2 \right] = \mathbb{E}_{\mathbf{x} \sim p_t} \left[\left(\frac{1}{p_t(\mathbf{x})} \frac{\partial p_t(\mathbf{x})}{\partial t} \right)^2 \right] \quad (88)$$

$$= \int_{\mathbf{x} \in \mathbb{R}^D} \frac{(p_1(\mathbf{x}) - p_0(\mathbf{x}))^2}{p_t(\mathbf{x})} d\mathbf{x} = \int_{\mathbf{x} \in \mathbb{R}^D} \frac{(p_1(\mathbf{x}) - p_0(\mathbf{x}))^2}{(1-t)p_0(\mathbf{x}) + tp_1(\mathbf{x})} d\mathbf{x} \quad (89)$$

$$\leq \int_{\mathbf{x} \in \mathbb{R}^D} \frac{p_1(\mathbf{x})^2 + p_0(\mathbf{x})^2}{(1-t)p_0(\mathbf{x}) + tp_1(\mathbf{x})} d\mathbf{x} \quad (90)$$

577 where we choose to keep the dependency on t in the bound.

578 We briefly justify this choice. We had first tried a t -independent bound, which led to an upper bound
 579 of the MSE that was too loose. We share insight as to why: first, recognize that the fraction can be
 580 broken in two terms, each of them a chi-square divergence between an endpoint of the path (p_0 or
 581 p_1) and the mixture p_t . Each of them admits a t -independent upper bound given by the chi-square
 582 divergence between the endpoints p_0 and p_1 , using lemma 1. However, the chi-square divergence
 583 between two gaussians, for example, is exponential (not polynomial) in the natural parameters [39,
 584 eq 7.41]. In fact, plotting $I(t)$ for a univariate gaussian model revealed that it took high values at the
 585 endpoints $t = 0$ and $t = 1$, and was near zero almost everywhere else in the interval $t \in [0, 1]$, which
 586 again suggested that dropping the dependency on t was unreasonable.

587 Now we can compute the estimation error, as

$$\text{MSE} = \frac{1}{N} \int_0^1 I(t) dt \leq \frac{1}{N} \int_{\mathbb{R}^d} \int_0^1 \frac{p_1(\mathbf{x})^2 + p_0(\mathbf{x})^2}{(1-t)p_0(\mathbf{x}) + tp_1(\mathbf{x})} dt d\mathbf{x} = \frac{1}{N} (J_1 + J_2) . \quad (91)$$

588 Let us try to solve one of these integrals, say J_1 .

$$J_1 = \int_{\mathbb{R}^d} \int_0^1 \frac{p_1(\mathbf{x})^2}{(1-t)p_0(\mathbf{x}) + tp_1(\mathbf{x})} dt d\mathbf{x} = \int_{\mathbb{R}^d} \frac{p_1(\mathbf{x})^2}{p_0(\mathbf{x})} \left(\int_0^1 \frac{1}{1 + t \left(\frac{p_1(\mathbf{x})}{p_0(\mathbf{x})} - 1 \right)} dt \right) d\mathbf{x} \quad (92)$$

$$= \int_{\mathbb{R}^d} \frac{p_1(\mathbf{x})^2}{p_0(\mathbf{x})} \left(\frac{1}{\frac{p_1}{p_0} - 1} \log \frac{p_1}{p_0} \right) d\mathbf{x} = 1 + \mathbb{E}_{p_1} \left[\frac{1}{1 - \frac{p_0}{p_1}} \log \frac{p_1}{p_0} - 1 \right] = 1 + \mathcal{D}_\phi(p_0, p_1) . \quad (93)$$

589 which we rewrote using an f-divergence defined by $\phi(x) = \frac{-\log(x)}{1-x} - 1$. Similarly, we obtain

$$J_2 = \int_{\mathbb{R}^d} \int_0^1 \frac{p_0(\mathbf{x})^2}{(1-t)p_0(\mathbf{x}) + tp_1(\mathbf{x})} dt d\mathbf{x} = \int_{\mathbb{R}^d} \frac{p_0(\mathbf{x})^2}{p_1(\mathbf{x})} \left(\int_0^1 \frac{1}{\frac{p_0(\mathbf{x})}{p_1(\mathbf{x})} + t(1 - \frac{p_0(\mathbf{x})}{p_1(\mathbf{x})})} dt \right) d\mathbf{x} \quad (94)$$

$$= \int_{\mathbb{R}^d} \frac{p_0(\mathbf{x})^2}{p_1(\mathbf{x})} \left(\frac{1}{\frac{p_0}{p_1} - 1} \log \frac{p_0}{p_1} \right) d\mathbf{x} = 1 + \mathbb{E}_{p_0} \left[\frac{1}{1 - \frac{p_1}{p_0}} \log \frac{p_0}{p_1} - 1 \right] = 1 + \mathcal{D}_\phi(p_1, p_0) . \quad (95)$$

590 Putting this together, we get

$$\text{MSE} \leq \frac{1}{N} (2 + \mathcal{D}_\phi(p_0, p_1) + \mathcal{D}_\phi(p_1, p_0)) . \quad (96)$$

591 How does this divergence depend on the parameter-distance $\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\|$? Does it bring down the
592 dependency from exponential to something lower? We next analyze this:

$$\mathcal{D}_\phi(p_0, p_1) + 1 = \mathbb{E}_{p_1} \frac{1}{1 - \frac{p_0}{p_1}} \log \frac{p_1}{p_0}$$

593 which looks like a Kullback-Leibler divergence, where the integrand is reweighted by $\frac{1}{1 - \frac{p_0(\mathbf{x})}{p_1(\mathbf{x})}}$. Note

594 that

$$\begin{cases} \frac{1}{1 - \frac{p_0(\mathbf{x})}{p_1(\mathbf{x})}} \geq 1 & p_0(\mathbf{x}) \leq p_1(\mathbf{x}) \\ \frac{1}{1 - \frac{p_0(\mathbf{x})}{p_1(\mathbf{x})}} < 1 & p_0(\mathbf{x}) > p_1(\mathbf{x}) \end{cases} \quad (97)$$

595 which motivates separating the integral over both domains

$$1 + \mathcal{D}_\phi(p_0, p_1) = \int_{\{\mathbf{x} \in \mathbb{R}^D | p_0(\mathbf{x}) \leq p_1(\mathbf{x})\}} p_1(\mathbf{x}) \log \frac{p_1(\mathbf{x})}{p_0(\mathbf{x})} \frac{1}{1 - \frac{p_0(\mathbf{x})}{p_1(\mathbf{x})}} \quad (98)$$

$$+ \int_{\{\mathbf{x} \in \mathbb{R}^D | p_0(\mathbf{x}) > p_1(\mathbf{x})\}} p_1(\mathbf{x}) \log \frac{p_1(\mathbf{x})}{p_0(\mathbf{x})} \frac{1}{1 - \frac{p_0(\mathbf{x})}{p_1(\mathbf{x})}} \quad (99)$$

$$\leq \int_{\{\mathbf{x} \in \mathbb{R}^D | p_0(\mathbf{x}) \leq p_1(\mathbf{x})\}} p_1(\mathbf{x}) + \int_{\{\mathbf{x} \in \mathbb{R}^D | p_0(\mathbf{x}) > p_1(\mathbf{x})\}} p_1(\mathbf{x}) \log \frac{p_1(\mathbf{x})}{p_0(\mathbf{x})} \quad (100)$$

$$\leq 1 + D_{\text{KL}}(p_1, p_0) \quad (101)$$

596 Hence we get

$$\text{MSE} \leq \frac{1}{N} \times (2 + D_{\text{KL}}(p_0, p_1) + D_{\text{KL}}(p_1, p_0)) . \quad (102)$$

597 We now suppose the proposal and target are distributions in an exponential family. The KL divergence
598 between exponential distributions with parameters $\boldsymbol{\theta}_0$ and $\boldsymbol{\theta}_1$, is given by the Bregman divergence of
599 the log-partition on the swapped parameters [30, Eq. 29]

$$D_{\text{KL}}(p_0, p_1) = \mathcal{D}_{\log Z}^{\text{Bregman}}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_0) := \log Z(\boldsymbol{\theta}_1) - \log Z(\boldsymbol{\theta}_0) - \nabla \log Z(\boldsymbol{\theta}_0)(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0) \quad (103)$$

$$\leq \frac{L}{2} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\|^2 \quad (104)$$

600 Hence

$$\text{MSE} \leq \frac{1}{N} \times (2 + L \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\|^2) \quad (105)$$

601 using the L -smoothness of the log-partition function $\log Z(\boldsymbol{\theta})$.

602 **Discussion on the assumptions for theorems 2, 3, 4, 5** For these theorems, we have supposed that
603 the target and proposal distributions are in an exponential family with a log partition that verifies

$$M \mathbf{Id} \preceq \nabla_{\boldsymbol{\theta}}^2 \log Z(\boldsymbol{\theta}) \preceq L \mathbf{Id} . \quad (106)$$

604 We now look at the validity of this assumption for a simple example: the univariate gaussian, which
605 is in an exponential family. The canonical parameters are its mean and variance (μ, v) . Written as an
606 exponential family,

$$p(\mathbf{x}) := \exp(\langle \boldsymbol{\theta}, \mathbf{t}(\mathbf{x}) \rangle - \log Z(\boldsymbol{\theta})) \quad (107)$$

607 the natural parameters are $\boldsymbol{\theta} = (\mu/v, -1/(2v))$, associated with the sufficient statistics $\mathbf{t}(x) =$
608 (x, x^2) [30]. The log-partition function and its derivatives are

$$\log Z(\boldsymbol{\theta}) = -\frac{\theta_1^2}{4\theta_2} - \frac{1}{2} \log(-2\theta_2) \quad (108)$$

$$\nabla \log Z(\boldsymbol{\theta}) = \mathbb{E}_{x \sim p}[t(x)] = \left(-\frac{\theta_1}{2\theta_2}, -\frac{1}{2\theta_2} + \frac{\theta_1^2}{4\theta_2^2} \right) \quad (109)$$

$$\nabla^2 \log Z(\boldsymbol{\theta}) = \text{Var}_{x \sim p}[t(x)] = \frac{1}{2\theta_2} \begin{pmatrix} -1 & \frac{\theta_1}{\theta_2} \\ \frac{\theta_1}{\theta_2} & \frac{1}{2} \frac{\theta_1^2}{\theta_2} \end{pmatrix} = \begin{pmatrix} v & 2\mu v \\ 2\mu v & 2v^2 - \mu^2 \end{pmatrix} \quad (110)$$

609 When the mean is zero, the eigenvalues of the hessian are in fact the diagonal values $(v, 2v^2)$, and
 610 they are bounded if and only if the variance v is bounded.

611 **Proof of Theorem 6** *Constant error of annealed NCE with the arithmetic path and "oracle-trig"*
 612 *schedule*

613 We now study the estimation error produced by the arithmetic path with "oracle-trig" schedule
 614 (table 2, line 5). We write the optimal path of Eq. 59 in the limit where the distributions do not
 615 overlap: $p_0(\mathbf{x})p_1(\mathbf{x}) \rightarrow 0$ pointwise and is bounded by an integrable function. In that limit, many
 616 quantities involved in the optimal distribution simplify

$$\mathcal{D}_{H^2}(p_0, p_1) = 1 - \int \sqrt{p_0 p_1} \rightarrow 1 \quad (111)$$

$$\alpha_H = \arctan \left(\sqrt{\frac{\mathcal{D}_{H^2}(p_0, p_1)}{2 - \mathcal{D}_{H^2}(p_0, p_1)}} \right) \rightarrow \frac{\pi}{4} \quad (112)$$

$$a(t) = \frac{\cos((2t-1)\alpha_H)}{2 \cos(\alpha_H)} - \frac{\sin((2t-1)\alpha_H)}{2 \sin(\alpha_H)} \rightarrow \cos\left(\frac{\pi t}{2}\right) \quad (113)$$

$$b(t) = \frac{\cos((2t-1)\alpha_H)}{2 \cos(\alpha_H)} + \frac{\sin((2t-1)\alpha_H)}{2 \sin(\alpha_H)} \rightarrow \sin\left(\frac{\pi t}{2}\right) . \quad (114)$$

617 All these limits are pointwise: for the first line, the dominated convergence theorem is used to justify
 618 the pointwise convergence of the integral $\int p_0 p_1 \rightarrow 0$ (L2 convergence of $\sqrt{p_0 p_1}$ and consequently
 619 the pointwise convergence of the integral $\int \sqrt{p_0 p_1} \rightarrow 0$ (L1 convergence of $\sqrt{p_0 p_1}$). This leads to
 620 the following simplification of the optimal path

$$p_t(\mathbf{x}) = (a(t)\sqrt{p_0(\mathbf{x})} + b(t)\sqrt{p_1(\mathbf{x})})^2 = a(t)^2 p_0(\mathbf{x}) + b(t)^2 p_1(\mathbf{x}) + 2a(t)b(t)\sqrt{p_0 p_1} \quad (115)$$

$$\rightarrow \cos^2\left(\frac{\pi t}{2}\right) p_0(\mathbf{x}) + \sin^2\left(\frac{\pi t}{2}\right) p_1(\mathbf{x}) \quad (116)$$

621 which is the arithmetic path with "oracle-trig" schedule defined in table 2 (line 5). The trigonometric
 622 weights evolve slowly at the end points $t = 0$ and $t = 1$. The estimation error in Eq. 63 produced by
 623 this path converges to

$$\text{MSE} = \frac{1}{N} \int_0^1 I(t) dt = \frac{1}{N} 16\alpha_H^2 \sim \frac{1}{N} 16 \frac{\pi^2}{8} = \frac{1}{N} 2\pi^2 . \quad (117)$$

624 which is a constant function of the parameter-distance.

625 **C Useful Lemma**

626 **Lemma 1** *Chi-square divergence of between a density and a mixture* We wish to upper bound the
 627 *chi-square divergence between a distribution $p(\mathbf{x})$ and a mixture $wp(\mathbf{x}) + (1 - w)q(\mathbf{x})$, where*
 628 $0 < w < 1$.

$$\mathcal{D}_{\chi^2}(p, wp + (1 - w)q) = \int_{\mathbf{x} \in \mathbb{R}^D} \frac{p(\mathbf{x})^2}{wp(\mathbf{x}) + (1 - w)q(\mathbf{x})} d\mathbf{x} - 1 \quad (118)$$

$$= \int_{\{\mathbf{x} \in \mathbb{R}^D | p(\mathbf{x}) < q(\mathbf{x})\}} \frac{p(\mathbf{x})^2}{wp(\mathbf{x}) + (1 - w)q(\mathbf{x})} d\mathbf{x} \quad (119)$$

$$+ \int_{\{\mathbf{x} \in \mathbb{R}^D | p(\mathbf{x}) > q(\mathbf{x})\}} \frac{p(\mathbf{x})^2}{wp(\mathbf{x}) + (1 - w)q(\mathbf{x})} d\mathbf{x} - 1 \quad (120)$$

$$\leq \int_{\{\mathbf{x} \in \mathbb{R}^D | p(\mathbf{x}) < q(\mathbf{x})\}} \frac{p(\mathbf{x})^2}{wp(\mathbf{x}) + (1 - w)p(\mathbf{x})} d\mathbf{x} \quad (121)$$

$$+ \int_{\{\mathbf{x} \in \mathbb{R}^D | p(\mathbf{x}) > q(\mathbf{x})\}} \frac{p(\mathbf{x})^2}{wq(\mathbf{x}) + (1 - w)q(\mathbf{x})} d\mathbf{x} - 1 \quad (122)$$

$$\leq \int_{\mathbf{x} \in \mathbb{R}^D} p(\mathbf{x}) d\mathbf{x} + \int_{\mathbf{x} \in \mathbb{R}^D} \frac{p(\mathbf{x})^2}{q(\mathbf{x})} d\mathbf{x} - 1 \quad (123)$$

$$= \int_{\mathbf{x} \in \mathbb{R}^D} \frac{p(\mathbf{x})^2}{q(\mathbf{x})} d\mathbf{x} = \mathcal{D}_{\chi^2}(p, q) + 1 \quad (124)$$