Statistical Efficiency of Score Matching: The View from Isoperimetry

Frederic Koehler* Alexander Heckett[†] Andrej Risteski[‡]

December 26, 2022

Abstract

Deep generative models parametrized up to a normalizing constant (e.g. energy-based models) are difficult to train by maximizing the likelihood of the data because the likelihood and/or gradients thereof cannot be explicitly or efficiently written down. Score matching is a training method, whereby instead of fitting the likelihood $\log p(x)$ for the training data, we instead fit the score function $r \times \log p(x)$ — obviating the need to evaluate the partition function. Though this estimator is known to be consistent, its unclear whether (and when) its statistical efficiency is comparable to that of maximum likelihood — which is known to be (asymptotically) optimal. We initiate this line of inquiry in this paper, and show a tight connection between statistical efficiency of score matching and the isoperi-metric properties of the distribution being estimated — i.e. the Poincare, log-Sobolev and isoperimetric constant — quantities which govern the mixing time of Markov processes like Langevin dynamics. Roughly, we show that the score matching estimator is statistically comparable to the maximum likelihood when the distribution has a small isoperimetric constant. Conversely, if the distribution has a large isoperimetric constant — even for simple families of distributions like exponential families with rich enough sufficient statistics — score matching will be substantially less efficient than maximum likelihood. We suitably formalize these results both in the finite sample regime, and in the asymptotic regime. Finally, we identify a direct parallel in the discrete setting, where we connect the statistical properties of pseudolikelihood estimation with approximate tensorization of entropy and the Glauber dynamics.

1 Introduction

Energy-based models (EBMs) are deep generative models parametrized up to a constant of parametrization, namely p(x) / exp(f(x)). The primary training challenge is the fact that evaluating the likelihood (and gradients thereof) requires evaluating the partition function of the model, which is generally computationally intractable — even when using relatively sophisticated MCMC techniques. Recent works, including the seminal paper of Song and Ermon [2019], circumvent this difficulty by instead fitting the score function of the model, that is $r_x \log p(x)$. Though not obvious how to evaluate this loss from training samples only, Hyvarinen [2005] showed this can be done via integration by parts, and the estimator is consistent (that is, converges to the correct value in the limit of infinite samples).

The maximum likelihood estimator is the de-facto choice for model-fitting for its well-known property of being statistically optimal in the limit where the number of samples goes to infinity [Van der Vaart, 2000]. It is unclear how much worse score matching can be — thus, it's unclear how much statistical efficiency we sacrifice for the algorithmic convenience of avoiding partition functions. In the seminal paper [Song and Ermon, 2019], it was conjectured that multimodality, as well as a low-dimensional manifold structure may cause difficulties for score matching — which was the reason the authors proposed annealing by convolving the input samples with a sequence of Gaussians with different variance. Though the intuition for this is natural: having poor estimates for the score in "low probability" regions of

^{*}fkoehler@stanford.edu, Stanford University. Supported in part by NSF award CCF-1704417, NSF award IIS-1908774, and N. Anari's Sloan Research Fellowship.

[†]aheckett@andrew.cmu.edu, Carnegie Mellon University.

[‡] a ristesk@andrew.cmu.edu, Carnegie Mellon University. Supported in part by NSF award IIS-2211907 and an Amazon Research Award on "Causal + Deep Out-of-Distribution Learning".

the distribution can "propagate" into bad estimates for the likelihood once the score vector field is "integrated" — making this formal seems challenging.

We show that the right mathematical tools to formalize, and substantially generalize such intuitions are functional analytic tools that characterize isoperimetric properties of the distribution in question. Namely, we show three quantities, the Poincaré, log-Sobolev and isoperimetric constants (which are all in turn very closely related, see Section 2), tightly characterize how much worse the efficiency of score matching is compared to maximum likelihood. These quantities can be (equivalently) viewed as: (1) characterizing the mixing time of Langevin dynamics — a stochastic differential equation used to sample from a distribution p(x) / exp(f(x)), given access to a gradient oracle for f; (2) characterizing "sparse cuts" in the distribution: that is sets S, for which the surface area of the set S can be much smaller than the volume of S. Notably, multimodal distributions, with well-separated, deep modes have very big log-Sobolev/Poincaré/isoperimetric constants [Gayrard et al., 2004, 2005], as do distributions supported over manifold with negative curvature [Hsu, 2002] (like hyperbolic manifolds). Since it is commonly thought that complex, high dimensional distribution deep generative models are trained to learn do in fact exhibit multimodal and low-dimensional manifold structure, our paper can be interpreted as showing that in many of these settings, score matching may be substantially less statistically efficient than maximum likelihood. Thus, our results can be thought of as a formal justification of the conjectured challenges for score matching in Song and Ermon [2019], as well as a vast generalization of the set of "problem cases" for score matching. This also shows that surprisingly, the same obstructions for efficient inference (i.e. drawing samples from a trained model, which is usual done using Langevin dynamics for EBMs) are also an obstacle for efficient learning using score matching.

We roughly show the following results:

- 1. For finite number of samples n, we show that if we are trying to estimate a distribution from a class with Rademacher complexity bounded by R_n , as well as a log-Sobolev constant bounded by C_{LS} , achieving score matching loss at most implies that we have learned a distribution that's no more than $C_{LS}R_n$ away from the data distribution in KL divergence. The main tool for this is showing that the score matching objective is at most a multiplicative factor of C_{LS} away from the KL divergence to the data distribution.
- 2. In the asymptotic limit (i.e. as the number of samples n ! 1), we focus on the special case of estimating the parameters of a probability distribution of an exponential family fp(x) / exp(h; F(x)i) for some sufficient statistics F using score matching. If the distribution p we are estimating has Poincare constant bounded by C_P have asymptotic efficiency that differs by at most a factor of C_P . Conversely, we show that if the family of sufficient statistics is sufficiently rich, and the distribution p we are estimating has isoperimetric constant lower bounded by C_{LS} , then the score matching loss is less efficient than the MLE estimator by at least a factor of C_{LS} .
- 3. Based on our new conceptual framework, we identify a precise analogy between score matching in the continuous setting and pseudolikelihood methods in the discrete (and continuous) setting. This connection is made by replacing the Langevin dynamics with its natural analogue the Glauber dynamics (Gibbs sampler). We show that the approximation tensorization of entropy inequality [Marton, 2013, Caputo et al., 2015], which guarantees rapid mixing of the Glauber dynamics, allows us to obtain finite-sample bounds for learning distributions in KL via pseudolikelihood in an identical way to the log-Sobolev inequality for score matching. A variant of this connection is also made for the related ratio matching estimator of Hyvariñen [2007b].
- 4. In Section 7, we perform several simulations which illustrate the close connection between isoperimetry and the performance of score matching. We give examples both when fitting the parameters of an exponential family and when the score function is fit using a neural network.

2 Preliminaries

Definition 1 (Score matching). Given a ground truth distribution p with sufficient decay at infinity and a smooth distribution q, the score matching loss (at the population level) is defined to be

$$J_{p}(q) := \frac{1}{2} E_{Xp}[kr \log p(X) \quad r \log q(X)k^{2}] + K_{p} = E_{Xp} \quad Tr \, r^{2} \log q + \frac{1}{2} k_{P}^{1} \log qk^{2}$$
 (1)

where K_p is a constant independent of q. The last equality is due to Hyvärinen [2005]. Given samples from p, the training loss $J_0^{\Lambda}(q)$ is defined by replacing the rightmost expectation with the average over data.

Functional and Isoperimetric Inequalities. Let q(x) be a smooth probability density over R^d . A key role in this work is played by the log-Sobolev, Poincaré, and isoperimetric constants of q — closely related geometric quantities, connected to the mixing of the Langevin dynamics, which have been deeply studied in probability theory and geometric and functional analysis (see e.g. [Gross, 1975, Ledoux, 2000, Bakry et al., 2014]).

Definition 2. The log-Sobolev constant $C_{LS}(q)$ 0 is the smallest constant so that for any probability density p(x)

$$KL(p;q) C_{LS}(q)I(pjq)$$
 (2)

where $KL(p;q) = E_{Xp}[log(p(X)=q(X))]$ is the Kullback-Leibler divergence or relative entropy and the relative Fisher information I(p j q) is defined ¹ as $I(p j q) := E_q r log \frac{p}{a}$; $r_a \frac{p}{a}$.

The log-Sobolev inequality is equivalent to exponential ergodicity of the Langevin dynamics for q, a canonical Markov process which preserves and is used for sampling q, described by the Stochastic Differential Equation $dX_t = r \log q(X_t) dt + \sqrt{2} dB_t$. Precisely, if p_t is the distribution of the continuous-time Langevin Dynamics² for q started from X_0 p, then $I(p j q) = \frac{dt}{dt} L(p_t; q) j_{t=0}$ and so by integrating

$$KL(p_t;q) e^{t=C_{LS}} KL(p;q)$$
: (3)

This holding for any p and t is an equivalent characterization of the log-Sobolev constant (Theorem 3.20 of Van Handel [2014]). For a class of distributions P, we can also define the restricted log-Sobolev constant $C_{LS}(q; P)$ to be the smallest constant such that (2) holds under the additional restriction that p 2 P — see e.g. Anari et al. [2021b]. For P an infinitesimal neighborhood of p, the restricted log-Sobolev constant of q becomes half of the Poincare constant or inverse spectral gap $C_P(q)$:

Definition 3. The Poincaré constant C_P (q) 0 is the smallest constant so that for any smooth function f,

$$Var_q(f) C_P(q)E_qkrfk^2:$$
 (4)

It is related to the log-Sobolev constant by C_P 2 C_{LS} (Lemma 3.28 of Van Handel [2014]).

Similarly, the Poincaré inequality implies exponential ergodicity for the ²-divergence:

$$^{2}(p_{t};q) e^{2t=C_{p} 2}(p;q)$$
:

This holding for every p and t is an equivalent characterization of the Poincaré constant (Theorem 2.18 of Van Handel [2014]). We can equivalently view the Langevin dynamics in a functional-analytic way through its definition as a Markov semigroup, which is equivalent to the SDE definition via the Fokker-Planck equation [Van Handel, 2014, Bakry et al., 2014]. From this perspective, we can write $p_t = qH_{tq}$ where H_t is the Langevin semigroup for q, so $H_t = e^{tL}$ with generator

$$Lf = hrlogq; rfi + f:$$

In this case, the Poincaré constant has a direct interpretation in terms of the inverse spectral gap of L, i.e. the inverse of the gap between its two largest eigenvalues.

Both Poincaré and log-Sobolev inequalities measure the isoperimetric properties of q from the perspective of functions; they are closely related to the isoperimetric constant:

¹There are several alternatives formulas for I (p j q), see Remark 3.26 of Van Handel [2014].

²See e.g. Vempala and Wibisono [2019] for more background and the connection to the discrete time dynamics.

Definition 4. The isoperimetric constant $C_{1S}(q)$ is the smallest constant, s.t. for every set S,

where S = fx : d(x; S) g and d(x; S) denotes the (Euclidean) distance of x from the set S. The isoperimetric constant is related to the Poincare constant by $C_P = 4C_{1|S}^2$ (Proposition 8.5.2 of Bakry et al. [2014]). Assuming S is chosen so $_{S}^{R} q(x)dx < 1=2$, the left hand side can be interpreted as the volume and the right hand side as the surface area of S with respect to q.

A strengthened isoperimetric inequality (Bobkov inequality) upper bounds the log-Sobolev constant, see Ledoux [2000], Bobkov [1997].

Mollifiers We recall the definition of one of the standard mollifiers/bump functions, as used in e.g. Hormander [2015]. Mollifiers are smooth functions useful for approximating non-smooth functions: convolving a function with a mollifier makes it "smoother", in the sense of the existence and size of the derivatives. Precisely, define the (infinitely differentiable) function : R^d ! R as

$$(y) = \int_{1^{d}}^{1} e^{\int_{1^{d}}^{1} e$$

where $I_d := {R \atop e} e^{-1=(1-jyj^2)} dy$.

We will use the basic estimate 8 d B $_d$ < I $_d$ < B $_d$ where B $_d$ is the volume of the unit ball in R d , which follows from the fact that e $^{1=(1\ jyj\)}$ $^{1=2}4$ for kyk $^{1=2}$ and e $^{1=(1\ jyj\)}$ $^{1=4}$ for kyk $^{1=2}$ for kyk $^{1=2}$ and e $^{1=(1\ jyj\)}$ $^{1=4}$ for kyk $^{1=2}$ for kyk $^{1=2}$ and e $^{1=(1\ jyj\)}$ $^{1=4}$ for kyk $^{1=2}$ for kyk $^{$

$$r_y$$
 $(y) = (2=I_d)e^{-1=(1-kyk^2)} \frac{y}{(1-kyk^2)^2} = \frac{2y}{(1-kyk^2)^2}$ (y)

It is straightforward to check that $\sup_{y} k r_y$ (y)k < 1=I_d. For > 0, we'll also define a "sharpening" of , namely (y) = $\binom{d}{d}$ (y=) so that = 1 and (by chain rule)

$$r_y$$
 $(y) = {}^{d-1}(r)(y=) = \frac{2y={}^{2}(y=)}{(y=k^{2})}$

so in particular $k r_v k_2$ d $^{1}=I_d$.

Glauber dynamics. The Glauber dynamics will become important in Section 5 as the natural analogue of the Langevin dynamics. The Glauber dynamics or Gibbs sampler for a distribution p is the standard sampler for discrete spin systems — it repeatedly selects a random coordinate and then resamples the spin X_i there according to the distribution p conditional on all of the other ones (i.e. conditional on X_i). See e.g. Levin and Peres [2017]. This is the standard sampler for discrete systems, but it also applies and has been extensively studied for continuous ones (see e.g. Marton [2013]).

Reach and Condition Number of a Manifold. For a smooth submanifold M of Euclidean space, the reach $_{\rm M}$ is the smallest radius r so that every point with distance at most r to the manifold M has a unique nearest point on M [Federer, 1959]; the reach is guaranteed to be positive for compact manifolds. The reach has a few equivalent characterizations (see e.g. Niyogi et al. [2008]); a common terminology is that the condition number of a manifold is $1=_{\rm M}$.

Notation. For a random vector X , $\chi := E[X X^T] E[X]E[X]^T$ denotes its covariance matrix.

3 Learning Distributions from Scores: Nonasymptotic Theory

Though consistency of the score matching estimator was proven in Hyvärinen [2005], it is unclear what one can conclude about the proximity of the learned distribution from a finite number of samples. Precisely, we would like a guarantee that shows that if the training loss (i.e. empirical estimate of (1)) is small, the learned distribution is close to the ground truth distribution (e.g. in the KL divergence sense). However, this is not always true! We will see an illustrative example where this is not true in Section 7 and also establish a general negative result in Section 4.

In this section, we prove (Theorem 1) that minimizing the training loss does learn the true distribution, assuming that the class of distributions we are learning have bounded complexity and small log-Sobolev constant. First, we formalize the connection to the log-Sobolev constant:

Proposition 1. The log-Sobolev inequality for q is equivalent to the following inequality over all smooth probability densities p:

$$KL(p;q) 2C_{LS}(q)(J_p(q) J_p(p)):$$
(6)

More generally, for a class of distribution p 2 P the restricted log-Sobolev constant is the smallest constant such that KL(p;q) $C_{LS}(q;P)(J_p(q))$ for all distributions p.

Proof. This follows from the following equivalent form for the relative Fisher information (e.g. Shao et al. [2019], Vempala and Wibisono [2019])

$$I(p j q) = E_q hr \frac{p}{q}; rlog \frac{p}{q}i$$

$$= E_p h \frac{q}{p} r \frac{p}{q}; rlog \frac{p}{q}i = E_p hrlog \frac{p}{q}; rlog \frac{p}{q}i = E_p krlog p rlog qk^2:$$
 (7)

Using this and (1) the log-Sobolev inequality can be rewritten as KL(p;q) $C_{LS}(J_p(q))$ which proves the first claim, and the same argument shows the second claim.

Remark 1 (Interpretation of Score Matching). The left hand side of (6) is $KL(p;q) = E_p[log p]$ $E_p[log q]$. The first term is independent of q and the second term is the likelihood, the objective for Maximum Likelihood Estimation. So (6) shows that the score matching objective is a relaxation (within a multiplicative factor of $C_{LS}(q)$) of maximum-likelihood via the log-Sobolev inequality. We discuss connections to other proposed interpretations in Appendix A.

Remark 2. Interestingly, the log-Sobolev constant which appears in the bound is that of q and not p the ground truth distribution. This is useful because q is known to the learner whereas p is only indirectly observed. If q is actually close to p, the log-Sobolev constants are comparable due to the Holley-Stroock perturbation principle (Proposition 5.1.6 of Bakry et al. [2014]).

The connection between the score matching loss and the relative Fisher information used in (7) is not new to this work—see the Related Work section for more discussion and references. The useful statistical implications which we discuss next are new to the best of our knowledge. Combining Proposition 1, bounds on log-Sobolev constants from the literature, and fundamental tools from generalization theory allows us to derive finite-sample guarantees for learning distributions in K L divergence via score matching. ³

Theorem 1. Suppose that P is a class of probability distributions containing p and define

$$C_{LS}(P; P) := \sup_{q \ge P} C_{LS}(q; P) \sup_{q \ge P} C_{LS}(q)$$

to be the worst-case (restricted) log-Sobolev constant in the class of distributions. (For example, if every distribution in P is -strongly log concave then C_{LS} 1=2 by Bakry-Emery theory [Bakry et al., 2014].) Let

$$R_n := E_{X_1; ...; X_n; 1; ...; n} \sup_{i=1}^{n} \frac{1}{i} \operatorname{Tr} r^2 \log q(X_i) + \frac{1}{i} \ker_{i=1}^{n} \operatorname{og} q(X_i) k^2 = 1$$

³We use the simplest version of Rademacher complexity bounds to illustrate our techniques. ²Standard literature, e.g. Shalev-Shwartz and Ben-David [2014], Bartlett et al. [2005] contains more sophisticated versions, and our techniques readily generalize.

be the expected Rademacher complexity of the class given n samples $X_1; :::; X_n$ pi.i.d. and independent $_1; :::;_n$ U nif1g i.i.d. Rademacher random variables. Let p be the score matching estimator from n samples, i.e. p = arg min_{q2P} J_p(q). Then

$$EKL(p; p) 4C_{LS}(P; P)R_n$$
:

In particular, if $C_{LS}(P; P) < 1$ then $\lim_{n \geq 1} E_{KL}(p; p^n) = 0$ as long as $\lim_{n \geq 1} R_n = 0$.

Proof. By the standard symmetrization argument (Theorem 26.3 of Shalev-Shwartz and Ben-David [2014]) we have $EJ_p(p)$ $J_p(p)$ $2R_n$, so by Proposition 1 we have EKL(p;p) $EC_{LS}(P)(J_p(p))$ $J_p(p)$ $2C_{LS}(P)R_n$.

Example 1. Suppose we are fitting an isotropic Gaussian in d dimensions with unknown mean satisfying kk R. The class of distributions P is q with kk R of the form $q(x) / exp kx k^2=2$ so the expected Rademacher complexity can be upper bounded as so:

where the inequality is Jensen's inequality and in the last step we expanded the square and used that $E_{ij} = 1(i = j)$ and $E_i k X_i k^2 = R^2 + d$. Recall that the standard Gaussian distribution is 1-strongly log concave so $C_{LS} = 1=2$. Hence we have the concrete bound $E_i k L(p; p^s) = R^2 + d$.

4 Statistical cost of score matching: asymptotic results

In this section, we compare the asymptotic efficiency of the score matching estimator in exponential families to the effiency of the maximum likelihood estimator. Because we are considering asymptotics, we might expect (recall the discussion in Section 2) that the relevant functional inequality will be the local version of the log-Sobolev inequal-ity around the true distribution p, which is the Poincare inequality for p. Our results will show precisely how this occurs and characterize the situations where score matching is substantially less statistically efficient than maximum likelihood.

Setup. In this section, we will focus on distributions from exponential families. We will consider estimating the parameters of an exponential family using two estimators, the classical maximum likelihood estimator (MLE), and the score matching estimator; we will use that the score matching estimator arg min $_0$ J $_p$ (p_0) admits a closed-form formula in this setting.

Definition 5 (Exponential family). For sufficient statistics $F:R^d!R^m$, the exponential family of distributions associated with F is $fp(x) / exp(h; F(x)i) j 2 R^mg$:

Definition 6 (MLE, Van der Vaart [2000]). Given i.i.d. samples $x_1; ::: ; x_n$ p, the maximum likelihood estimator is $M_L^c = \arg \max_{0.2} E[\log p_0(X)]$, where E denotes the expectation over the samples. As n ! 1 and under appropriate regularity conditions, we have $p_{\overline{n}} = N(0; M_LE)$, where $M_LE := 1$ and M_LE

Proposition 2 (Score matching estimator, Equation (34) of Hyvarinen [2007b]). Given i.i.d. samples $x_1; ::: ; x_n = p$, the score matching estimator equals $\int_{SM} e^{-\frac{p}{2}} [(JF)_X (JF)_X^T]^{-\frac{1}{p}} f$, where $(JF)_X : m = d$ is the Jacobian of F at the point X , $f = \int_{-1}^{1} e^{-\frac{p}{2}} f$ is the Laplacian and it is applied coordinate wise to the vector-valued function F .

4.1 Asymptotic normality

Next, we recall the asymptotic normality of the score matching estimator and give a formula for the limiting renormalized covariance matrix $_{S\,M}$ established by Forbes and Lauritzen [2015] (see also Theorem 6 of Barp et al. [2019] and Corollary 1 of Song et al. [2020]). Since the MLE also satisfies asymptotic normality with an explicit covariance matrix, we can then proceed in the next sections to compare their relative efficiency (as in e.g. Section 8.2 of Van der Vaart [2000]) by comparing the asymptotic covariances $_{S\,M}$ and $_{M\,LE}$.

Proposition 3 (Asymptotic normality, Forbes and Lauritzen [2015]). As n ! 1, and assuming sufficient smoothness and decay conditions so that score matching is consistent (see Hyvarinen [2005]) we have the following convergence in distribution: $P = \overline{n}(\hat{S}_M)$ | N (0; SM), where

$$SM := E[(JF)_X (JF)_X^T]^{-1}_{X} (JF)_X (JF)_X^{TF} E[(JF)_X (JF)^T]^{-1}_{X}.$$
 (8)

Proof. We include the proof for the reader's convenience. From Hyvärinen [2005], we have consistency of score matching (Theorem 2) and in particular the formula

$$= E[(JF)_X (JF)^T_{\vee}]^{-1}EF:$$
 (9)

We now compute the limiting distribution of the estimator as the number of samples $n \ ! \ 1$. We will need to use some standard results from probability theory such as Slutsky's theorem and the central limit theorem, see e.g. Van der Vaart [2000] or Durrett [2019] for references. To minimize ambiguity, let E_n denote the empirical expectation over n i.i.d. samples samples and let n denote the score matching estimator n from n samples. Define n; and n; by the equations

$$\mathbf{E}_{n}[(JF)_{X}(JF)_{X}^{T}] = E[(JF)_{X}(JF)_{X}^{T}] + {}_{n;1} = {}^{p}_{n}$$

and

$$f_n F = Ef + {}_{n;2} = {}^p n$$
:

By the central limit theorem, n = (n;1;n;2) converges in distribution to a multivariate Gaussian (with a covariance matrix that we won't need explicitly) as $n \mid 1$. From the definition

$$\hat{n} = E_{n}^{\wedge}[(JF)_{X}(JF)_{X}^{\top}]^{-1}E^{\uparrow}$$

$$= [E[(JF)_{X}(JF)_{X}^{\top}]^{-1}E_{n}[(JF)_{X}(JF)_{X}^{\top}]]^{-1}E[(JF)_{X}(JF)_{X}^{\top}]^{-1}E^{\uparrow}F$$

and we now simplify the expression on the right hand side. By applying (9) we have

$$E[(JF)_{X}(JF)_{X}^{T}]^{-1}\hat{E}_{n}F = E[(JF)_{X}(JF)^{T}]_{X}^{-1}(EF + _{n;2} = ^{p}n)$$

$$= + E[(JF)_{X}(JF)^{T}]_{n:2}^{-1} = p - _{n}$$

Since

$$E[(JF)_{x}(JF)_{x}^{T}]^{-1}\hat{E}_{n}[(JF)_{x}(JF)_{x}^{T}] = I + E[(JF)_{x}(JF)_{x}^{T}]^{-1}_{n;1} = p$$

and $(I + X)^{-1} = I + X + X^{2}$ we have by applying Slutsky's theorem that

$$E[(JF)_{X}(JF)_{X}^{T}] \stackrel{1}{\leftarrow} e_{n}[(JF)_{X}(JF)_{X}^{T}]] \stackrel{1}{\rightarrow} = I \quad E[(JF)_{X}(JF)_{X}^{T}] \stackrel{1}{\rightarrow}_{n;1} = \stackrel{p}{\rightarrow} n + O_{P}(1=n)$$

where we use the standard notation $Y_n = O_P(1=n)$ to indicate that $nY_n = f(n) ! 0$ in probability for any function f(n) ! 1. Hence

$$\hat{n} = [E[(JF)_{X}(JF)^{T}_{X}]^{-1}E^{\hat{n}}_{n}[(JF)_{X}(JF)^{T}_{X}]^{-1}E[(JF)_{X}(JF)^{T}_{X}]^{-1}E^{\hat{n}}_{n}F$$

$$= I E[(JF)_{X}(JF)^{T}_{X}]^{-1}_{n;1} = p^{\hat{n}}_{n} + O_{P}(1=n) (+ E[(JF)_{X}(JF)^{T}]_{X}^{-1}_{n;2} = n^{\hat{p}}_{n}$$

and applying Slutsky's theorem again, we find

From the definition, we know

$$\frac{1}{P_n}(n_{;1} \quad n_{;2}) = E_n[^{\wedge}(JF)_X(JF)_X \quad ^TF] \quad E[(JF)_X(JF)_X \quad F]$$

so altogether by the central limit theorem, we have

$$p$$
 $n(^{\land}_{n})$! N $0; E[(JF)_{X}(JF)^{T}]_{X}^{1}_{(JF)_{X}(JF)_{Y}+F^{T}}E[(JF)_{X}(JF)^{T}]^{1}_{X}$

as claimed.

4.2 Statistical efficiency of score matching under a Poincaré inequality

Our first result will show that if we are estimating a distribution with a small Poincaré constant (and some relatively mild smoothness assumptions), the statistical efficiency of the score matching estimator is not much worse than the maximum likelihood estimator.

Theorem 2 (Efficiency under a Poincaré inequality). Suppose the distribution p satisfies a Poincare inequality with constant C_P . Then we have

More generally, the same bound holds assuming only the following restricted version of the Poincaré inequality: for any w, Var(hw; F(x)i) C_P $Ekrhw; F(x)ik_2$.

Remark 3. To interpret the terms in the bound, the quantities $E_p \ k(JF)_X k_{OP}^4$ and $EkF k^2$ can be seen as a measure of the smoothness of the sufficient statistics F, and kk as a bound on the radius of parameters for the exponential family. In Section 7 we will give an example to show bounded smoothness is indeed necessary for score matching to be efficient.

Remark 4. A direct consequence of this result is that with 99% probability and for sufficiently large n,

nk
$$_{SM} k^2 \text{ nEk} \quad _{MLE} k^2 \stackrel{?}{_{\Lambda}} O C_P m kk^2 E_2 k (JF)_X k_{OP} + EkF k_4^2$$
 (10)

. So if the the distribution is smooth and Poincaré, score matching achieves small '2 error provided MLE does. To show this, since $p = 10^{10} \, \text{m} \, \text{$

$$nk = s_M^{\Lambda} k^2 = O(E_{ZN(0; s_M)} kZ k^2) = O(Tr_{SM}) = O(mk_{SM} k_{OP})$$
:

On the other hand, by Fatou's lemma we have that

$$\lim \inf_{n \mid 1} nEk \quad M_{LE}^{\Lambda} k^2 \quad E_{ZN(0; MLE)} kZk^2 = Tr(MLE) \quad k \quad MLE k_{OP}$$

where in the first expression $\bigwedge_{L \in I}$ implicitly depends on n, the number of samples. Combining these two observations with Theorem 2 and gives inequality 10.

The main lemma to prove the theorem is the following:

Lemma 1. $E[(JF)_X(JF)_X^T]^{-1}$ $C_{P_F}^{-1}$ where C_P is the Poincare constant of p.

Proof. For any vector w 2 R^m, we have by the Poincaré inequality that

$$C_P hw; E[(JF)_X (JF)_X^T]wi = C_P Ekr_x hw; F(x)ij_X k_2^2 Var(hw; F(x)i) = hw; Fwi$$

This shows $C_P E[(JF)_X (JF)_X^T]_F$ and inverting both sides, using the well-known fact that the matrix inverse is operator monotone [Toda, 2011], gives the result.

We will also need the following helper lemma:

Lemma 2. For any random vectors A; B we have A + B = 2A + 2B.

Proof. For any vector w we have

$$Var(hw; A + Bi) = Var(hw; Ai) + 2Cov(hw; Aihw; Bi) + Var(hw; Bi)$$

$$Var(hw; Ai) + 2 \frac{D}{Var(hw; Ai)Var(hw; Bi) + Var(hw; Bi)}$$

$$2Var(hw; Ai) + 2Var(hw; Bi)$$

where the first inequality is Cauchy-Schwarz for variance and the second is ab $a^2=2+b^2=2$. We proved for this for every vector which proves the PSD inequality.

With this in mind, we can proceed to the proof of Theorem 2:

Proof of Theorem 2. Recall from Proposition 3 that

$$_{SM} := E[(JF)_X(JF)_X^T]^{-1}_{(JF)_X(JF)_X^{\mp}F} E[(JF)_X(JF)^T]^{-1}_X$$

By Lemma 1 and submultiplicativity of the operator norm, we have

We will finally bound the two operator norms on the right hand side. By Lemma 2, we have

$$(JF)_x(JF)_x + 2_{(JF)_x(JF)_x} + 2_{FT}$$

Furthermore, we have

$$k_{(JF)x(JF)^T}k_{OP}k_{E[(JF)x(JF)^T}(J_xF)x(JF)^T]k_{OP}x_Ek(JF)xk_{OP}kk^2$$

and

$$k_F k_{OP} kE(F)(F)^T k_{OP} Tr E(F)(F)^T EkFk^2$$

which implies the statement of the theorem.

4.3 Statistical efficiency lower bounds from sparse cuts

In this section, we prove a converse to Theorem 2: whereas a small (restricted) Poincare constant upper bounds the variance of the score matching estimator, if the Poincare constant of our target distribution is large and we have sufficiently rich sufficient statistics, score matching will be extremely inefficient compared to the MLE. In fact, we will be able to do so by taking an arbitrary family of sufficient statistics, and adding a single sufficient statistic! Informally, we'll show the following:

Consider estimating a distribution p in an exponential family with isoperimetric constant C_{1S} . Then, p can be viewed as a member of an enlarged exponential family with one more $(O_{@S}(1)\text{-Lipschitz})$ sufficient statistic, such that score matching has asymptotic relative efficiency $_{@S}(C_{1S})$ compared to the MLE, where $_{@S}$ denotes the boundary of the isoperimetric cut of p and $_{@S}$ indicates a constant depending only on the geometry of the manifold $_{@S}$.

As noted in Section 2, a large Poincare constant implies a large isoperimetric constant — so we focus on showing that the score matching estimator is inefficient when there is a set S which is a "sparse cut". Our proof uses differential geometry, so our final result will depend on standard geometric properties of the boundary @S — e.g., we use the concept of the reach $_{\rm M}$ of a manifold which was defined in the preliminaries (Section 2). The full proofs are in Appendix B. We now give the formal statement.

Theorem 3 (Inefficiency of score matching in the presence of sparse cuts). There exists an absolute constant c > 0such that the following is true. Suppose that p is an element of an exponential family with sufficient statistic F_1 and parameterized by elements of 1. Suppose S is a set with smooth boundary @S which has reach $_{@S} > 0$. Suppose that 1_S is not an affine function of F_1 , so there exists $_1 > 0$ such that

$$:= 1 \qquad \qquad ^{p} \frac{ }{1 - 1 + 2} \stackrel{q}{\stackrel{x_{2} \otimes_{S}}{p_{r}(X \times 2S)(1 - p_{r}(X \times 2S))}} \stackrel{2}{\stackrel{}{}} . \text{ Define an additional sufficient statistic } F_{2} = 1_{S} \quad \text{ so that the }$$

enlarged exponential family contains distributions of the form

$$p_{(1;2)}(x) / exp(h_1; F_1(x)i + {}_2F_2(x))$$

and consider the MLE and score matching estimators in this exponential family with ground truth $p_{(;0)}$.

Then there exists some w so that the relative (in)efficiency of the score matching estimator compared to the MLE for estimating hw; i admits the following lower bound

$$\frac{\text{hw; }_{\text{SM}} \text{ wi}}{\text{hw; }_{\text{MLE}} \text{ wi}} \qquad \frac{\text{c}^0}{\text{minfPr}(\text{X 2 S}); \text{Pr}(\text{X $\not\supseteq$ S})g} \\ \\ \times_{\text{x2@S}} \text{p(x)dx}$$

where
$$c^0 := \frac{c^{-d}}{1 + k_{F,1} k_{O,P}}$$
.

Remark 5. If we choose S to be the set achieving the worst isoperimetric constant, then the right hand side of the bound is simply $c^0 C_{1S}$. (See the appendix for details.) Finally, we observe that although c^0 is exponentially small in d, the bound is still useful in high dimensions because in the bad cases of interest CIS is often exponentially large in d. For example, this is the case for a mixture (d) separation between the means (see e.g. Chen et al. [2021a]).

Remark 6. The assumption $_1 > 0$ is a quantitative way of saying that the function 1_S , the cut we are using to define the new sufficient statistic F2, is not already a linear combination of the existing sufficient statistics. The assumptions will always holds with some 1 0 by the Cauchy-Schwarz inequality. The equality case is when 1s is an affine function of hw₁; F₁i — if such a linear dependence exists, the parameterization is degenerate and the coefficient of F₂ is not identifiable as ! 0.

Proof sketch. The proof of the theorem proceed in two parts: we lower bound hw; MLE wi and upper bound hw; MLE WI. The former part, which ends up to be somewhat involved, proceeds by proving a lower bound on the spectral norm of SM (the full proof is in Subsection B.1) — by picking a direction in which the quadratic form is large. The upper bound on 2 (w) (the full proof is in Subsection B.2) will proceed by relating the Fisher matrix for the augmented sufficient statistic (F_1 ; F_2) with the Fisher matrix for the original sufficient statistic F_1 . Since the Fisher matrix is a covariance matrix in exponential families, this is where the numerator minfPr(X 2 S); Pr(X ≥ S)g, which is up to constant factors the variance of 1_S , naturally arises in the theorem statement.

For the lower bound, it is clear that we should select w which changes the distribution a lot, but not the observed gradients. The w we choose that satisfies these desiderata is proportional to $E[(JF)_X(JF)^T](0;1)$. This w also has the property that it results in a simple expression for the quadratic form hw; SM wi, using the fact

$$_{SM} = E[(JF)_X(JF)_X^T]^{-1}_{(JF)_X(JF)_X^{\mp}F}E[(JF)_X(JF)^T]^{-1}_X$$

The result of this calculation (details in Lemma 4, Appendix B) is that

hw;
$$_{SM}$$
 wi $_{SM}$ $_{SM}$

Note_Rthat the key term Pr[d(X; @S)], when divided by and in the limit ! 0, corresponds to the surface area x2@S p(x)dx of the cut. Showing that the other terms do not "cancel" this one out and determining the precise dependence on requires a differential-geometric argument, which is somewhat more intricate. The two key ideas are to use the divergence theorem (or generalized Stokes theorem) to rewrite the numerator as a more interpretable surface integral and then rigorously argue that as ! 0 and we "zoom in" to the manifold, we can compare to the case when the surface looks flat. The quantitative version of this argument involves geometric properties of the manifold (precisely, the curvature and reach). For example, Lemma 6 makes rigorous the statement that well-conditioned (i.e. large reach) manifolds are locally flat. More details, as well as the full proof is included in Appendix B.

Example application. We provide an instantiation of the theorem for a simple example of a bimodal distribution:

Example 2. A concrete example in one dimension with a single sufficient statistic is

$$F_1(x) = \frac{1}{8a^2}(x - a)^2(x + a)^2 = x^4 = 8a^2 + x^2 = 4 - a^2 = 8$$

and = (1; 0) for a parameter a > 1 to be taken large. This looks similar to a mixture of standard Gaussians centered at a and a. Specializing Theorem 3 to this case, we get:

Corollary 1. There exists absolute constants $_0 > 0$ and $_0 > 0$ so that the following is true. Suppose that $_0 > 0$, $_0 < 0$, and expanded exponential family $_0 < 0$, with $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$, $_0 < 0$

$$\frac{\text{hw; sm wi}}{\text{hw; mle wi}}$$
 ce $a^2=8$:

Proof of Corollary 1. First observe that

Z₁
$$Z_1$$
 Z_1 Z_1 Z_1 Z_1 Z_1 Z_1 Z_1 Z_1 Z_1 Z_2 Z_1 Z_2 Z_2 Z_3 Z_4 Z_4 Z_5 Z_5

where C is a positive constant independent of a. Using that $F_1(x) = (1=8)(x-a)^2(x=a+1)^2$ it then follows that

$$Pr(X 2 [a 1; a+1]) = \frac{R_{a+1}}{R_{a+1}} e^{F_1(x)} dx \frac{e^{(1=8)(x=a+1)^2}}{C} C > 0$$

where C^0 is a positive constant independent of a. From this, we see by the law of total variance that $Var(F_1)$ $Var(F_1 \ j \ X \ 2 \ [a \ 1; a + 1]) \ Pr(X \ 2 \ [a \ 1; a + 1]) \ C^{00} > 0$ where $C^{00} > 0$ is another positive constant independent of a. Hence $k_1^F \ k_D^F = O(1)$ independent of a. Also, if we define S = fx : x > 0g then

$$Cov(F_1(x); 1_S) = 0$$

becuase $F_1(x)$ is even, 1_S is odd and the distribution is symmetric about zero. So we can take $_1=1$ in the statement of Theorem 3. Therefore, applying Theorem 3 to S and using that $F_1(0)=a^2=8$, we therefore get for smaller than an absolute constant, that the inefficiency is lower bounded by $(e^{a^2=8}=)$. By taking equal to a fixed constant we get the result.

In Section 7, we perform simulations which show the performance of score matching indeed degrades exponentially as a becomes large.

5 Discrete Analogues: Pseudolikelihood, Glauber Dynamics, and Approximate Tensorization

5.1 Pseudolikelihood

Several authors have proposed variants of score matching for discrete probability distributions, e.g. Lyu [2009], Shao et al. [2019], Hyvärinen [2007b]. Furthermore, Hyvärinen [2005, 2006, 2007b,a] pointed out some connections between pseudolikelihood methods (a classic alternative to maximum likelihood in statistics Besag [1975, 1977]), Glauber dynamics (a.k.a. Gibbs sampler, see Preliminaries), and score matching. Finally, just like the log-Sobolev inequality controls the rapid mixing of Langevin dynamics, there are functional inequalities [Gross, 1975, Bobkov and Tetali, 2006] which bound the mixing time of Glauber dynamics. Thus, we ask: Is there a discrete analogue of the relationship between score matching and the log-Sobolev inequality?

The answer is yes. To explain further, we need a key concept recently introduced by Marton [2013, 2015] and Caputo et al. [2015]: if $\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$; ::: ($\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$ are arbitrary measure spaces, we say a distribution q on is satisfies approximation tensorization of entropy with constant C_{AT} (q) if

$$KL(p;q) C_{AT}(q) E_{X_ip_i}[KL(p(X_ijX_i);q(X_ijX_i))]:$$
(12)

This inequality is sandwiched between two discrete versions of the log-Sobolev inequality (Proposition 1.1 of Caputo et al. [2015]): it is weaker than the standard discrete version of the log-Sobolev inequality [Diaconis and Saloff-Coste, 1996] and stronger than the Modified Log-Sobolev Inequality [Bobkov and Tetali, 2006] which characterizes exponential ergodicity of the Glauber dynamics.⁴ We define a restricted version C_{AT} (q; P) analogously to the restricted log-Sobolev constant.

Finally, we recall the pseudolikelihood objective [Besag, 1975] based on entrywise conditional probabilities: $L_p(q) := \int_{i=1}^{d} E_{Xp}[\log q(X_i j X_i)]$. With these definition in place, we have:

Proposition 4. We have KL(p;q) $C_{AT}(q)(L_p(p)$ $L_p(q))$ and more generally for any class P containing p, we have KL(p;q) $C_{AT}(q;P)(L_p(p)$ $L_p(q))$.

Proof. Observe that $L_p(p)$ $L_p(q) = \bigcap_{i=1}^{p} E_{X_i,jp} [KL(p(X_i j X_i); q(X_i j X_i))]$, so the result follows by expanding the definition.

Thus, just as the score matching objective is a relaxation of maximum likelihood through the log-Sobolev inequality, pseudolikelihood is a relaxation through approximate tensorization of entropy.

Remark 7. Pseudolikelihood methods (and variants like node-wise regression) are one of the dominant approaches to fitting fully-observed graphical models, e.g. [Wu et al., 2019, Lokhov et al., 2018, Klivans and Meka, 2017, Kelner et al., 2020]. Like score matching, pseudolikelihood methods do not require computing normalizing constants which can be slow or computationally hard (e.g. Sly and Sun [2012]). Pseudolikelihood is applicable in both discrete and continuous settings, as is our connection with approximate tensorization.

We state explicitly the analogue of Theorem 1 for pseudolikelihood, which follows from the same proof by replacing Proposition 1 with Proposition 4.

Theorem 4. Suppose that P is a class of probability distributions containing p and $C_{AT}(P; P) := \sup_{q \neq P} C_{AT}(q; P)$ is the worst-case (restricted) approximate tensorization constant in the class of distributions (e.g. bounded by a constant if all of the distributions in P satisfy a version of Dobrushin's uniqueness condition [Marton, 2015]). Let

$$R_{n} := E_{X_{1}; \dots; X_{n}; j; \dots; n} \sup_{q \ge P} \frac{1}{n} \prod_{i=1}^{X} \frac{2}{4} X \log_{j=1}^{d} q((X_{i})_{j} j(X_{i})_{j}) 5$$

⁴In most cases where the MLSI is known, approximate tensorization of entropy is also, e.g. Chen et al. [2021b], Anari et al. [2021a], Marton [2015], Caputo et al. [2015].

be the expected Rademacher complexity of the class given n samples $X_1; :::; X_n$ pi.i.d. and independent $_1; :::;_n$ U nif1g i.i.d. Rademacher random variables. Let p be the pseudolikelihood estimator from n samples, i.e. p = arg min_{q2P} L_p(q). Then

$$EKL(p; p^n) 2C_{AT}(P; P)R_n$$
:

In particular, if $C_{AT} < 1$ then $\lim_{n \neq 1} E_{KL}(p; p^n) = 0$ as long as $\lim_{n \neq 1} R_n = 0$.

5.2 Ratio Matching

[Hyvärinen, 2007b] proposed a version of score matching for distributions on the hypercube f1g^d and observed that the resulting method ("ratio matching") bears similarity to pseudolikelihood. A similar calculation as the proof of Proposition 4 allows us to arrive at ratio matching based on a strengthening of approximate tensorization studied in [Marton, 2015]. Our derivation seems more conceptual than the original derivation, explains the similarity to pseudolikelihood, and establishes some useful connections.

Marton [2015] studied a strengthened version of approximate tensorization of the form

$$KL(p;q) C_{AT2}(q) \sum_{i=1}^{X^d} E_{X_ip_i} TV^2(p(X_i j X_i); q(X_i j X_i))$$
(13)

where T V denotes the total variation distance (see Cover [1999]). (This is known to hold for a class of distributions q satisfying a version of Dobrushin's condition and marginal bounds [Marton, 2015].) This inequality is stronger than the standard approximate tensorization because of Pinsker's inequality $TV^2(P; Q)$. KL(P; Q) [Cover, 1999]. In the case of distributions on the hypercube, we have

$$\begin{split} TV^2(p(X_i j X_i); q(X_i j X_i)) &= jp(X_i = +1 j X_i) \quad q(X_i = +1 j X_i)j^2 \\ &= E_{X_i p_{X_i j X_i}} j1(X_i = +1) \quad q(X_i = +1 j X_i)j^2 \\ &= E_{X_i p_{X_i j X_i}} j1(X_i = +1) \quad p(X_i = +1 j X_i)j^2 \end{split}$$

where in the last step we used the Pythagorean theorem applied to the p_{X+iX+}-orthogonal decomposition

$$1(X_i = +1) \quad q(X_i = +1 j X_i) = [1(X_i = +1) \quad p(X_i = +1 j X_i)] + [p(X_i = +1 j X_i) \quad q(X_i = +1 j X_i)]$$

Hence, there exists a constant $K_p^{\ C}$ not depending on q such that

where we define the ratio matching objective function to be

$$M_{p}(q) := \sum_{i=1}^{X^{d}} E_{X p} j 1 (X_{i} = +1) \quad q(X_{i} = +1 j X_{i}) j^{2}$$
(15)

This objective is now straightforward to estimate from data, by replacing the expectation with the average over data. Analogous to before, we have the following proposition:

Proposition 5. We have

$$KL(p;q)$$
 $C_{AT2}(q)(M_p(q)$ $M_p(p))$

and more generally for any class P containing p, we have KL(p;q) $C_{AT2}(q;P)(M_p(q)$ $M_p(p))$.

We now show how to rewrite $M_p(q)$ to match the formula from the original reference. Observe

$$M_{p}(q) = \frac{1}{4} X^{d} E_{Xp} j X_{i} \qquad E_{q}[X_{i} j X_{i}] j^{2} = \frac{1}{4} X^{d} E_{Xp} j 1 \qquad X_{i} E_{q}[X_{i} j X_{i}] j^{2} = 1$$

Observe that for any z 2 f1g we have

$$zE[X_{i} j X_{i}] = \frac{q(X_{i} = z j X_{i}) - q(X_{i} = z j X_{i})^{q}}{q(X_{i} = z j X_{i}) + q(X_{i} = z j X_{i})}$$

and

$$1 z E_{q}[X_{i} j X_{i}] = \frac{2q(X_{i} = z j X)}{q(X_{i} = z j X_{i}) + q(X_{i} = z j X_{i})}$$

$$= \frac{2}{1 + q(X_{i} = z j X_{i}) = q(X_{i} = z j X_{i})}$$

Also for z 2 f1g^d we have $q(X_i = z_i j X_i = z_i) = q(X_i = z_i j X_i = z_i) = q(z) = q(z_i)$ where z is represents z with coordinate if lipped, so

$$M_p(q) = X_{i=1}^d E_{Xp} \frac{1}{1 + q(X) = q(X_i)}$$

which matches the formula in Theorem 1 of Hyvärinen [2007b].

Summarizing, minimizing the ratio matching objective makes the right hand side of the strengthened tensorization estimate (13) small, so when $C_{AT\,2}(q)$ is small it will imply successful distribution learnig in KL. (The obvious variant of Theorem 4 will therefore hold.) In this way ratio matching can also be understood as a relaxation of maximum likelihood.

6 Related work

Score matching was originally introduced by Hyvärinen [2005], who also proved that the estimator is asymptotically consistent. In [Hyvärinen, 2007b], the authors propose estimators that are defined over bounded domains. [Song and Ermon, 2019] scaled the techniques to neurally parameterized energy-based models, leveraging score matching versions like denoising score matching Vincent [2011], which involves an annealing strategy by convolving the data distribution with Gaussians of different variances, and sliced score matching [Song et al., 2020]. The authors conjectured that annealing helps with multimodality and low-dimensional manifold structure in the data distribution — and our paper can be seen as formalizing this conjecture.

The connection between Hyvarinen's score matching objective and the relative Fisher information in (7) is known in the literature — see e.g. [Shao et al., 2019, Nielsen, 2021, Barp et al., 2019, Vempala and Wibisono, 2019, Yamano, 2021]. Relatedly, Hyvarinen [2007a] pointed out some connections between the score matching objective and contrastive divergence using the lens of Langevin dynamics. We also remark that since $I(pjq) = \frac{d}{dt}KL(p_t;q)$ $j_{t=0}$ for p_t the output of Langevin dynamics at time t, score matching can be interpreted as finding a q to minimize the contraction of the Langevin dynamics for q started at p. Previously, [Guo, 2009, Lyu, 2009] observed that the score matching objective can be interpreted as the infinitesimal change in KL divergence as we add Gaussian noise — see Appendix A for an explanation why these two quantities are equal. We note that Hyvarinen [2008] also gave a related interpretation of score matching in terms of adding an infinitesimal amount of Gaussian noise.

In the discrete setting, it was recently observed that approximate tensorization has applications to identity testing of distributions in the "coordinate oracle" query model [Blanca et al., 2022], which is another application of approximate tensorization outside of sampling otherwise unrelated to our result. Finally, [Block et al., 2020, Lee et al., 2022a] show guarantees on running Langevin dynamics, given estimates on r log p that are only -correct in the L₂(p) sense. They show that when the Langevin dynamics are run for some moderate amount of time, the drift between the true Langevin

dynamics (using r log p exactly) and the noisy estimates can be bounded. Recent concurrent works [Lee et al., 2022b, Chen et al., 2022] show results of a similar flavor for denoising diffusion model score matching, specifically when the forward SDE is an Ornstein-Uhlenbeck process.

7 Simulations

7.1 Exponential family experiments.

Fitting a bimodal distribution with a cut statistic. First, we show the result of fitting a bimodal distribution (as in Example 2) from an exponential family. In Figure 1, the difference of the two sufficient statistics we consider corresponds to the cut statistic used in our negative result (Theorem 3). As predicted (Corollary 1) score matching performs poorly compared to the MLE as the distance between modes grows.

In Figure 2, we illustrate the distribution of the errors in the bimodal experiment with the cut statistic. As expected based on the theory, the direction where score matching with large offset performs very poorly corresponds to the difference between the two sufficient statistics, which encodes the sparse cut in the distribution.

Fitting a bimodal distribution without a cut statistic. In Figure 3 we show the result of fitting the same bimodal distribution using score matching, but we remove the second sufficient statistic (which is correlated with the sparse cut in the distribution). In this case, score matching fits the distribution nearly as well as the MLE. This is consistent with our theory (e.g. the failure of score matching in Theorem 3 requires that we have a sufficient statistic approximately representing the cut) and justifies some of the distinctions we made in our results: even though the Poincaré constant is very large, the asymptotic variance of score matching within the exponential family is upper bounded by the restricted Poincaré constant (see Theorem 2) which is much smaller.

Example 3 (Application of Theorem 2 to this example). To briefly expand the last point, we show how to apply Theorem 2 in this example (Example 2, where we have not added a bad cut statistic.) The restricted Poincare constant for applying Theorem 2 will be

$$C := \frac{Var(F_1(X))}{E(F_1^{c}(X))^2} = \frac{Var(X^2 - X^4 = 2a^2)}{E(2X - 2X^3 = a^2)^2}$$
 (16)

which asymptotically goes to a constant, rather than blowing up exponentially, as a goes to infinity. (This can be made formal using arguments as in the proof of Corollary 1; informally, the distribution is similar to a mixture of two standard Gaussians centered at a so the numerator is close to $Var_{ZN(0;1)}((a+Z)^2 (a+Z)^4=2a^2) = Var(2aZ+Z)(4aZ+6Z+4Z=a+Z)=2)=(1)$ 4and the denominator is approximately $E_{ZN(0;1)}(2(a+Z)-2(a+Z))=2$ 0 and $E_{ZN(0;1)}(2(a+Z)-2(a+Z))=2$ 1.) 3 2 2

Given this bound on the restricted Poincaré constant, we can apply Theorem 2. Based on similar reasoning to above, one can show that $EF_1^0(X)^4 = (1=4a^2)^4E((X=a)(X+a)^2+(X=a)^2(X+a))^4 = (1)$ and $EF_1^{00}(X)^2 = E(3x^2=2a^2+1=2)^2 = (1)$, so we conclude that $k_{SM}k_{OP} = O(k_{MLE}k^2_{O})_P$ This proves that score matching will perform not much worse than the MLE, as we saw in the experimental result of Figure 3.

Remark 8. Example 3 shows a case where there is a large gap between the restricted and unrestricted Poincaré constants. This also implies a completely analogous gap between appropriate restricted and unrestricted log-Sobolev constants, as used e.g. in the context of Theorem 1. To elaborate, we know that the unrestricted log-Sobolev constant blows up exponentially in a, just like the unrestricted Poincare constant, because C_{LS} C_P =2 [Van Handel, 2014]. On the other hand, if we fix the ground truth distribution p_a consider the class of distributions

$$P_r = fp_{a^0} : ja \quad a^0 j rg;$$

we have that

$$\lim_{r \downarrow 0} C_{LS}(q; P_r) = C=2$$

where C is the constant defined in (16) in terms of a (and which is O(1) as a ! 1). This is because from the definition as an exponential family, we have

$$p_a(x)=p_{a^0}(x) = \frac{\exp((a - a^0)F_1(x))}{E_{a^0}\exp(((a - a^0)F_1(x)))}$$

SO

$$\lim_{a^0!a} \frac{\mathsf{KL}(p_a;p_{a^0})}{\mathsf{I}(p_a\;j\;p_{a^0})} = \lim_{a^0!a} \frac{(a-a^0)^2 \mathsf{Var}_{p_{a^0}}(\mathsf{F}_1(x))}{2(a-a^0)^2 \mathsf{E}_{p_{a^0}} \mathsf{krF}_1(x) \mathsf{k}^2} = C = 2$$

where the first equality is by a standard Taylor expansion argument (see proof of Lemma 3.28 of [Van Handel, 2014]).

Fitting a unimodal distribution with rapid oscilation. In Figure 5, we demonstrate what happens when the distribution is unimodal (and has small isoperimetric constant), but the sufficient statistic is not quantitatively smooth. More precisely, we consider the case $p(x) / e^{-0x^2-2} e^{-1\sin(!x)}$ as ! increases. In the figure, we used the formulas from asymptotic normality to calculate the distribution over parameter estimates from 100,000 samples. We also verified via simulations that the asymptotic formula almost exactly matches the actual error distribution.

The result is that while the MLE can always estimate the coefficient $_1$ accurately, score matching performs much worse for large values of !. This demonstrates that the dependence on smoothness in our results (in particular, Theorem 2) is actually required, rather than being an artifact of the proof. Conceptually, the reason score matching fails even when though the distribution has no sparse cuts is this: the gradient of the log density becomes harder to fit as the distribution becomes less smooth (for example, the Rademacher complexity from Theorem 1 will become larger as it scales with r_{x} log p and r^{2} log p).

7.2 Score matching with neural networks

Fitting a mixture of Gaussians with a one-layer network. We also show that empirically, our results are robust even beyond exponential families. In Figure 4 we show the results of fitting a mixture of two Gaussians via score matching⁵, where the score function is parameterized as a one hidden-layer network with tanh activations. We see that the predictions of our theory persist: the distribution is learned successfully when the two modes are close and is not when the modes are far. This matches our expectations, since the Poincare, log-Sobolev, and isoperimetric constants blow up exponentially in the distance between the two modes (see e.g. Chen et al. [2021a]) and the neural network is capable of detecting the cut between the two modes.

In the right hand side example (the one with large separation between modes), the shape of the two Gaussian components is learned essentially perfectly — it is only the relative weights of the two components which are wrong. This closely matches the idea behind the proof of the lower bound in Theorem 3; informally, the feedforward network can naturally represent a function which detects the cut between the two modes of the distribution, i.e. the additional bad sufficient statistic F_2 from Theorem 3. The fact that the shapes are almost perfectly fit where the distribution is concentrated indicates that the test loss J_p is near its minimum. Recall from (1) that the suboptimality of a distribution q in score matching loss is given by $J_p(q) = E_p kr \log p = r \log q k^2$. If we let q be the distribution recovered by score matching, we see from the figure that the slopes of the distribution were correctly fit wherever p is concentrated, so $E_p kr \log p = r \log q k^2$ is small. However near-optimality of the test loss $J_p(q)$ does not imply that q is actually close to p: the test loss does not heavily depend on the behavior of $\log q$ in between the two modes, but the value of $r \log q$ in between the modes affects the relative weight of the two modes of the distribution, leading to failure.

Both models illustrated in the figure have 2048 tanh units and are trained via SGD on fresh samples for 300000 steps. After training the model, the estimated distribution is computed from the learned score function using numerical integration.

⁵We note that this experiment is similar in flavor to plots in (Figure 2) in Song and Ermon [2019], where they show that the score is estimated poorly near the low-probability regions of a mixture of Gaussians. In our plots, we numerically integrate the estimates of the score to produce the pdf of the estimated distribution.

8 Conclusion

In this paper, we studied the statistical efficiency of score matching and identified a close connection to functional inequalities which characterize the ergodicity of Langevin dynamics. For future work, it would be interesting to characterize formally the improvements conferred by annealing strategies like [Song and Ermon, 2019], like it has been done in the setting of sampling using Langevin dynamics [Lee et al., 2018].

Acknowledgements. We are grateful to Lester Mackey and Aapo Hyvärinen, as well as to the anonymous reviewers, for feedback on an earlier draft.

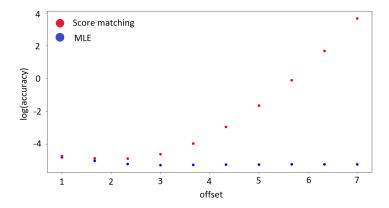


Figure 1: Statistical efficiency of score matching vs MLE for fitting the distribution with ground truth parameters (0; 1) = (1; 0) of the form $p(x) / e^{0(x^2 - x^4 = (2a^2)) + 1(x^2 - x^4 = (2a^2) + erf(x))}$ as we vary the offset a between 1 and 7 and train with fixed number of samples (10^5) . We see score matching (red) performs very poorly compared to the MLE (blue) as the offset (distance between modes) grows, by plotting the log of the Euclidean distance to the true parameter for both estimators.

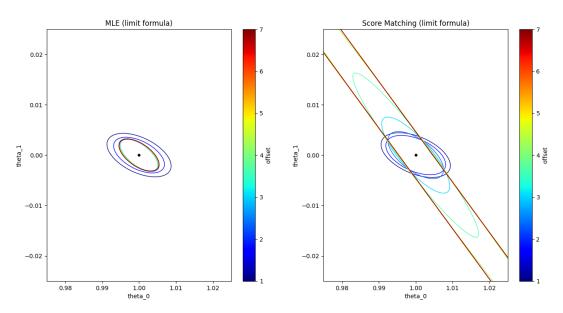


Figure 2: Level sets for the distribution over estimates in the same example as Figure 1. We see that as the distance a between modes increases, the direction of large variance for the score matching estimator (right figure) corresponds to the difference of the sufficient statistics which encodes the sparse cut in the distribution. On the other hand, the MLE (left figure) does not exhibit this behavior and has low variance in all directions.

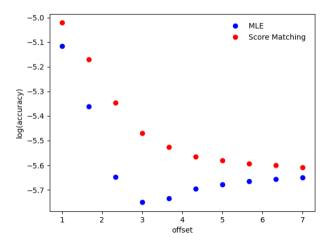
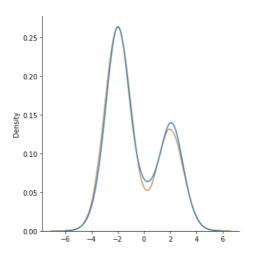


Figure 3: Here we see the result of running an identical experiment to Figure 1, only we remove the second sufficient statistic, so our distribution is now $p(x) / e^{o(x^2 - x^4 = (2a^2))}$ where o = 1 and we again vary the offset a between 1 and 7. With only the single sufficient statistic, score matching performs comparably to MLE.



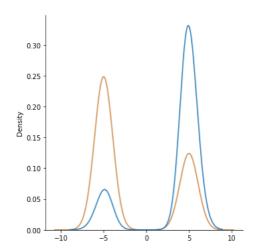


Figure 4: Training a single hidden-layer network to score match a mixture of Gaussians (ground truth orange, score matching output blue) succeeds at learning the distribution when the modes are close (left, small isoperimetric constant), but not when they are distant (right, large isoperimetric constant) in which case it weighs the modes incorrectly.

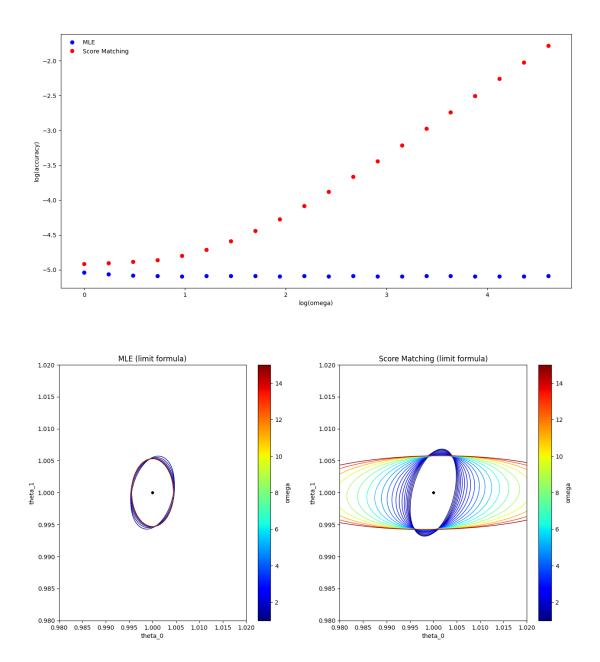


Figure 5: Score matching vs MLE for a distribution with a rapidly oscillating sufficient statistic, $p(x) / e^{0x^2-2} e^{1\sin(!x)}$ where $o^{1} = 0$, and increasing!. On the top, for increasing! we show a log-log plot of the average Euclidean distance in parameter space between and the output of each estimator. On the bottom, for each value of!, we draw a level set of the distribution within which a fixed fraction of returned estimates lie (MLE left, score matching right). Score matching becomes increasingly inaccurate as! increases while the MLE stays extremely accurate.

References

- Nima Anari, Vishesh Jain, Frederic Koehler, Huy Tuan Pham, and Thuy-Duong Vuong. Entropic independence in high-dimensional expanders: Modified log-sobolev inequalities for fractionally log-concave polynomials and the ising model. arXiv preprint arXiv:2106.04105, 2021a.
- Nima Anari, Vishesh Jain, Frederic Koehler, Huy Tuan Pham, and Thuy-Duong Vuong. Entropic independence ii: optimal sampling and concentration via restricted modified log-sobolev inequalities. arXiv preprint arXiv:2111.03247, 2021b.
- Dominique Bakry, Ivan Gentil, Michel Ledoux, et al. Analysis and geometry of Markov diffusion operators, volume 103. Springer, 2014.
- Alessandro Barp, François-Xavier Briol, Andrew Duncan, Mark Girolami, and Lester Mackey. Minimum stein discrepancy estimators. Advances in Neural Information Processing Systems, 32, 2019.
- Peter L Bartlett, Olivier Bousquet, and Shahar Mendelson. Local rademacher complexities. The Annals of Statistics, 33(4):1497–1537, 2005.
- Julian Besag. Statistical analysis of non-lattice data. Journal of the Royal Statistical Society: Series D (The Statistician), 24(3):179–195, 1975.
- Julian Besag. Efficiency of pseudolikelihood estimation for simple gaussian fields. Biometrika, pages 616-618, 1977.
- Antonio Blanca, Zongchen Chen, Daniel Štefankovič, and Eric Vigoda. Identity testing for high-dimensional distributions via entropy tensorization. arXiv preprint arXiv:2207.09102, 2022.
- Adam Block, Youssef Mroueh, and Alexander Rakhlin. Generative modeling with denoising auto-encoders and langevin sampling. arXiv preprint arXiv:2002.00107, 2020.
- Sergey G Bobkov. An isoperimetric inequality on the discrete cube, and an elementary proof of the isoperimetric inequality in gauss space. The Annals of Probability, 25(1):206–214, 1997.
- Sergey G Bobkov and Prasad Tetali. Modified logarithmic sobolev inequalities in discrete settings. Journal of Theoretical Probability, 19(2):289–336, 2006.
- Pietro Caputo, Georg Menz, and Prasad Tetali. Approximate tensorization of entropy at high temperature. In Annales de la Faculté des sciences de Toulouse: Mathématiques, volume 24, pages 691–716, 2015.
- Hong-Bin Chen, Sinho Chewi, and Jonathan Niles-Weed. Dimension-free log-sobolev inequalities for mixture distributions. Journal of Functional Analysis, 281(11):109236, 2021a.
- Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru R Zhang. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. arXiv preprint arXiv:2209.11215, 2022.
- Zongchen Chen, Kuikui Liu, and Eric Vigoda. Optimal mixing of glauber dynamics: Entropy factorization via high-dimensional expansion. In Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing, pages 1537–1550, 2021b.
- Thomas M Cover. Elements of information theory. John Wiley & Sons, 1999.
- Persi Diaconis and Laurent Saloff-Coste. Logarithmic sobolev inequalities for finite markov chains. The Annals of Applied Probability, 6(3):695–750, 1996.
- Rick Durrett. Probability: theory and examples, volume 49. Cambridge university press, 2019.
- Herbert Federer. Curvature measures. Transactions of the American Mathematical Society, 93(3):418-491, 1959.

- Peter GM Forbes and Steffen Lauritzen. Linear estimating equations for exponential families with application to gaussian linear concentration models. Linear Algebra and its Applications, 473:261–283, 2015.
- Véronique Gayrard, Anton Bovier, Michael Eckhoff, and Markus Klein. Metastability in reversible diffusion processes i: Sharp asymptotics for capacities and exit times. Journal of the European Mathematical Society, 6(4):399–424, 2004.
- Véronique Gayrard, Anton Bovier, and Markus Klein. Metastability in reversible diffusion processes ii: Precise asymptotics for small eigenvalues. Journal of the European Mathematical Society, 7(1):69–99, 2005.
- Alfred Gray. Tubes, volume 221. Springer Science & Business Media, 2003.
- Leonard Gross. Logarithmic sobolev inequalities. American Journal of Mathematics, 97(4):1061-1083, 1975.
- Dongning Guo. Relative entropy and score function: New information-estimation relationships through arbitrary additive perturbation. In 2009 IEEE International Symposium on Information Theory, pages 814–818. IEEE, 2009.
- Lars Hörmander. The analysis of linear partial differential operators I: Distribution theory and Fourier analysis. Springer, 2015.
- Elton P Hsu. Stochastic analysis on manifolds. Number 38. American Mathematical Soc., 2002.
- Aapo Hyvärinen. Estimation of non-normalized statistical models by score matching. Journal of Machine Learning Research, 6(4), 2005.
- Aapo Hyvärinen. Consistency of pseudolikelihood estimation of fully visible boltzmann machines. Neural Computation, 18(10):2283–2292, 2006.
- Aapo Hyvärinen. Connections between score matching, contrastive divergence, and pseudolikelihood for continuous-valued variables. IEEE Transactions on neural networks, 18(5):1529–1531, 2007a.
- Aapo Hyvärinen. Some extensions of score matching. Computational statistics & data analysis, 51(5):2499–2512, 2007b.
- Aapo Hyvärinen. Optimal approximation of signal priors. Neural Computation, 20(12):3087-3110, 2008.
- Jonathan Kelner, Frederic Koehler, Raghu Meka, and Ankur Moitra. Learning some popular gaussian graphical models without condition number bounds. Advances in Neural Information Processing Systems, 33:10986–10998, 2020.
- Adam Klivans and Raghu Meka. Learning graphical models using multiplicative weights. In 2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS), pages 343–354. IEEE, 2017.
- Michel Ledoux. The geometry of markov diffusion generators. In Annales de la Faculté des sciences de Toulouse: Mathématiques, volume 9, pages 305–366, 2000.
- Holden Lee, Andrej Risteski, and Rong Ge. Beyond log-concavity: Provable guarantees for sampling multimodal distributions using simulated tempering langevin monte carlo. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, Advances in Neural Information Processing Systems, vol-ume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper/2018/file/c6ede20e6f597abf4b3f6bb30cee16c7-Paper.pdf.
- Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence for score-based generative modeling with polynomial complexity. arXiv preprint arXiv:2206.06227, 2022a.
- Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence of score-based generative modeling for general data distributions. arXiv preprint arXiv:2209.12381, 2022b.
- David A Levin and Yuval Peres. Markov chains and mixing times, volume 107. American Mathematical Soc., 2017.

- Andrey Y Lokhov, Marc Vuffray, Sidhant Misra, and Michael Chertkov. Optimal structure and parameter learning of ising models. Science advances, 4(3):e1700791, 2018.
- Siwei Lyu. Interpretation and generalization of score matching. In Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence, pages 359–366, 2009.
- Katalin Marton. An inequality for relative entropy and logarithmic sobolev inequalities in euclidean spaces. Journal of Functional Analysis, 264(1):34–61, 2013.
- Katalin Marton. Logarithmic sobolev inequalities in discrete product spaces: a proof by a transportation cost distance. arXiv preprint arXiv:1507.02803, 2015.
- Frank Nielsen. Fast approximations of the jeffreys divergence between univariate gaussian mixtures via mixture conversions to exponential-polynomial distributions. Entropy, 23(11):1417, 2021.
- Partha Niyogi, Stephen Smale, and Shmuel Weinberger. Finding the homology of submanifolds with high confidence from random samples. Discrete & Computational Geometry, 39(1):419–441, 2008.
- Shai Shalev-Shwartz and Shai Ben-David. Understanding machine learning: From theory to algorithms. Cambridge university press, 2014.
- Stephane Shao, Pierre E Jacob, Jie Ding, and Vahid Tarokh. Bayesian model comparison with the hyvärinen score: Computation and consistency. Journal of the American Statistical Association, 2019.
- Allan Sly and Nike Sun. The computational hardness of counting in two-spin models on d-regular graphs. In 2012 IEEE 53rd Annual Symposium on Foundations of Computer Science, pages 361–369. IEEE, 2012.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. Advances in Neural Information Processing Systems, 32, 2019.
- Yang Song, Sahaj Garg, Jiaxin Shi, and Stefano Ermon. Sliced score matching: A scalable approach to density and score estimation. In Uncertainty in Artificial Intelligence, pages 574–584. PMLR, 2020.
- Alexis Akira Toda. Operator reverse monotonicity of the inverse. The American Mathematical Monthly, 118(1):82–83, 2011.
- Aad W Van der Vaart. Asymptotic statistics, volume 3. Cambridge university press, 2000.
- Ramon Van Handel. Probability in high dimension. Technical report, PRINCETON UNIV NJ, 2014.
- Santosh Vempala and Andre Wibisono. Rapid convergence of the unadjusted langevin algorithm: Isoperimetry suffices. Advances in neural information processing systems, 32, 2019.
- Pascal Vincent. A connection between score matching and denoising autoencoders. Neural computation, 23(7): 1661–1674, 2011.
- Hermann Weyl. On the volume of tubes. American Journal of Mathematics, 61(2):461–472, 1939.
- Shanshan Wu, Sujay Sanghavi, and Alexandros G Dimakis. Sparse logistic regression learns all discrete pairwise graphical models. Advances in Neural Information Processing Systems, 32, 2019.
- Takuya Yamano. Skewed jensen—fisher divergence and its bounds. Foundations, 1(2):256–264, 2021.

A Recovering an interpretation of score matching

We remarked that if we use the fact $I(pjq) = \frac{d}{dt} KL(p_t;q) j_{t=0}$, the score matching objective has a natural interpretation in terms of select q to minimize the contraction of the Langevin dynamics for q started at p. On the other hand, Guo [2009] and Lyu [2009] previously observed that the score matching objective can be interpreted as the infinitesimal change in KL divergence between p and q as we add noise to both of them, which is closely related to the de Bruijn identity. We now explain why these two quantities are equal by giving a proof of their equality (which is shorter than the one you get by going through the proof in Lyu [2009]).

Before giving the formal proof, we give some intuition for why the statement should be true. The Langevin dynamics approximately adds a noise of size N(0; 2t) and subtracts a gradient step along $r \log q$, and this dynamics preserves q. For small t, the gradient step is essentially reversible and preserves the KL. So heuristically, reversing the gradient step gives $KL(p_t; q)$ KL(N(0; 2t) p; N(0; 2t) q). We now give the formal proof.

Lemma 3. Assuming smooth probability densities p(x) and q(x) decay sufficiently fast at infinity,

$$\frac{d}{dt}KL(p_t;q) = \frac{d}{dt}KL(p N(0;2t);q N(0;2t))$$

where denotes convolution.

Proof. Recalling from Section 2 that $H_t = e^{tL}$ we have that $\frac{d}{dt} \frac{p_t}{q} = \frac{d}{dt} H_t \frac{p}{q} = L_q \frac{p}{q}$. Since $KL(p_t; q) = E_q[\frac{p_t}{q} \log \frac{p_t}{q}]$ and $\frac{d}{dx}[x \log x] = \log x + 1$, it follows by the chain rule that

$$\begin{split} \frac{d}{dt} K L(p_t;q) &= E_q \quad \log \frac{p}{q} + 1 \quad L \frac{p}{q} \quad = E_q \quad \log \frac{p}{q} + 1 \quad \text{hrlogq;} \\ &= E_q \quad \log \frac{p}{q} + 1 \quad \text{hrlogq;} \\ r \frac{p}{q} i + \frac{p}{q} \quad \frac{pq}{q^2} \end{split}$$

D E where in the last step we used the quotient rule $p = \frac{p}{q} - \frac{p}{q} - 2$ r log q; r $\frac{p}{q} - \frac{p}{q^2}$ the other hand, by using the Fokker-Planck equation $\frac{Q}{Q}(p N(0; 2t)) = p$ (Lemma 2 of Lyu [2009]) and the chain rule we have Z

$$\frac{\text{dit}}{\mathbb{E}} \begin{array}{c} \text{KL(p N(0;2t);q N(0;2t))} = \\ \text{dt} \frac{\text{d}}{\mathbb{E}} \\ \text{(q)} \\ \text{p} \\ \text{log} \\ \text{p} \\ \text{q} \\ \text{q} \\ \text{q} \\ \text{p} \\ \text{q} \\ \text{p} \\ \text{q} \\ \text{p} \\ \text{q} \\ \text{p} \\ \text{q} \\ \text{p} \\ \text{p} \\ \text{q} \\ \text{p} \\ \text{p} \\ \text{q} \\ \text{p} \\ \text{q} \\ \text{p} \\ \text{p} \\ \text{p} \\ \text{p} \\ \text{q} \\ \text{p} \\ \text{p} \\ \text{p} \\ \text{q} \\ \text{p} \\$$

Since by the chain rule and integration by parts we have

$$E_{q} \log \frac{p}{q} + 1 \operatorname{rlog} q; r \frac{p}{q} = \operatorname{rq}; r \frac{p}{q} \log \frac{p}{q} dx = \operatorname{Z} (q)_{q} \log \frac{p}{q} dx;$$

we see that the two derivatives are indeed equal.

B Proof of Theorem 3 and Applications

We restate Theorem 3 for the reader's convenience and in a slightly more detailed form (we include an upper bound on the covariance of the MLE error which follows from the proof).

Theorem 5 (Inefficiency of score matching in the presence of sparse cuts, Restatement of Theorem 3). There exists an absolute constant c > 0 such that the following is true. Suppose that p is an element of an exponential family with sufficient statistic F_1 and parameterized by elements of p. Suppose p is a set with smooth and compact boundary

@S. Let $_{@S} > 0$ denote the reach of @S (see Section 2) Suppose that 1_S is not an affine function of F_1 , so there exists ₁> 0 such that

$$\sup_{\substack{\text{w}_1: \text{Var}(\text{hw}_1; F_1 i) = 1\\ \text{Nupose that}}} \text{Cov} \quad \text{hw}_1; F_1 i; \\ \frac{1_S}{\text{Var}(1_S)} \qquad 1_1 : \qquad (17)$$

$$\sup_{\substack{\text{w}_1: \text{Var}(\text{hw}_1; F_1 i) = 1\\ \text{Nupose that}}} \text{Suppose that} \quad > 0 \text{ satisfies} \quad < \min_{\substack{(1+k_1k) \text{ sup}_{x:d(x; @ S)} \\ (1+k_1k) \text{ sup}_{x:d(x; @ S)} \\ k(JF_1)_x k_{OP}}; C_d^{\frac{S}{2}}} \text{ and is small enough so that } 0 < := 1$$

$$\frac{p}{1} \quad \text{1} + 2 \frac{q}{-pr(\frac{x^2 @ S}{2})(1-pr(x^2S))} \quad \text{Define an additional sufficient statistic } F_2 = 1_S \quad \text{so that the}$$

$$\frac{p}{1} - \frac{q}{1} + 2 \frac{R}{\Pr(X^2 S_S)(1 - \Pr(X 2 S_S))}$$
. Define an additional sufficient statistic $F_2 = 1_S$ so that the

enlarged exponential family contains distributions of the form

$$p_{(1;2)}(x) / exp(h_1; F_1(x)i + {}_2F_2(x))$$

and consider the MLE and score matching estimators in this exponential family with ground truth $p_{(;0)}$.

Then the asymptotic renormalized covariance matrix MLE of the MLE is bounded above as

and there there exists some w so that the relative (in)efficiency of the score matching estimator compared to the MLE for estimating hw; i admits the following lower bound

$$\frac{\text{hw; smwi}}{\text{hw; mlewi}} \quad \frac{\text{c}^0}{\text{minfPr(X 2 S); Pr(X 2 S)g}} \\ \frac{\text{c}^0}{\text{minfPr(X 2 S); Pr(X 2 S)g}}$$

where
$$c^0 := \frac{c^{-d}}{1 + k_{F_1} k_{OP}}$$
.

B.1 Lower bounding the spectral norm of SM

We recall the new statistic F₂, defined in terms of the mollifier introduced in Section 2:

and the new sufficient statistic is $F(x) = (F_1(x); F_2(x))$. We first show the following lower bound on the largest eigenvalue of $_{S\,M}$, the renormalized limiting covariance of score matching:

Lemma 4 (Largest eigenvalue of $_{S\,M}$). The largest eigenvalue of $_{S\,M}$ satisfies

Proof. We have

$$r_x F_2(x) = r_x (x y) dy; s r_x^2 F_2(x) = r_x^2 (x y) dy:$$

Defining

$$u := E[(JF)_X(JF)_X^{T}](0; 1) = E[(JF)_X r_x F_2(x)]$$

we have, by the variational characterization of eigenvalues of symmetric matrices, that

$$(s_{M}) \frac{hu; E[(JF)_{X}(JF)_{X}^{T}]^{1}_{(JF)_{X}(JF)_{X}^{T}} E[(JF)_{X}(JF)^{T}]^{1}_{X}ui^{max}}{kuk_{2}^{2}} :$$
 (19)

To upper bound the denominator we observe that if B_d is the volume of the unit ball,

$$1(d(x; @S))^{d-1}vol(B(X;))=I_d$$
 (21)

$$8^{d}1(d(x;@S))^{-1}$$
 (22)

and so

where we used the computation of the derivative of . To lower bound the numerator we have

hu;
$$E[(JF)_X(JF)_X^T]^{-1}_{(JF)_X(JF)^T}_{X}^T E[(JF)_X(JF)^T]^{-1}_{X}ui =$$

$$(0; 1)^T_{(JF)_X(JF)_X+F}(0; 1)$$

$$= Eh(0; 1); (JF)_X(JF)_X^T + Fi^2 = E(r_XF_2)^T(J_F)^T +_XF_2^{-2};$$

The integrand is zero except when d(X; @S) so it equals

$$Pr[d(X;@S)] E_{Xjd(X;@S)2[;]} (rF_2)_X (JF_T)_X + F_{T2}^2$$

and combining gives the result.

We now estimate the right hand side of (18) for small , using differential geometric techniques. The main idea is that as we take smaller, we end up zooming into the manifold @S which locally looks closer and closer to being flat. Differential-geometric quantities describing the manifold appear when we make this approximation rigorous. The most involved term to handle ends up to be calculating the expectation $E_{XJd}(X;@S)$ $(rF_2)^T(JF)^T_X+F_2^2$. To do this, we first argue that the term with the Laplacian dominates as ! 0, then by Stokes theorem, we end up integrating hr; dNi over intersections of S with small spheres of radius , where N is a normal to S. Such quantities can be calculated by comparing to the "flat" manifold case — i.e. when N does not change. How far away these quantities are (thus how small needs to be) depends on the curvature of S (or more precisely, the condition number of the manifold). Lemma 6 makes rigorous the statement that well-conditioned manifolds are locally flat and then Lemma 7, which is part of the proof of Weyl's tube formula [Gray, 2003, Weyl, 1939], lets us rigorously say that the tubular neighborhood (that is, a thickening of the manifold) behaves similarly to the flat case.

Lemma 5. There exists an absolute constant c > 0 such that the following is true. For any > 0 satisfying

< min
$$\frac{c^{d}}{(1 + k_{1}k) \sup_{x:d(x;@S)} k(JF_{1})_{x}k_{OP}}; c^{S}_{d}$$

for score matching on the extended family with m + 1 sufficient statistics and distribution p with = (1;0) we have

$$_{\text{max}}(s_{M}) = \frac{c^{d}}{_{@S}p(x)dA}$$

Proof. In the denominator, we can observe by (22) that

$$k(JF)_x k_{OP}^2 kJF_1 k_{OP}^2 + krF_2 k_2 kJF_1 k_{OP}^2 kJ^2 B_d 2^2 B_d$$

where the last inequality holds assuming is sufficiently small that kJF_1k^2 b $^2B^2$.

In the numerator we can observe

where the second-to-last expression is a surface integral which we arrived at by applying the divergence theorem, using that the Laplacian is the divergence of the gradient, and in the last step we used that — and all of its derivatives vanish on the boundary of the unit sphere.

Using that = (1, 0) we have

$$Z = \frac{1}{B(0;1)} (x + y)^{T} i = \frac{1}{x} = \frac{1}{B(0;1)} kr (u)kk(JF_{1})xk_{OP}kk$$
(23)

$$8^{d-1}k(JF_1)_Xk_{OP}kk: (24)$$

Let p be the point in @(X S) = which is closest in Euclidean distance to the origin. Let n(q) denote the unit normal vector at point q oriented outwards (Gauss map). Note that by first-order optimality conditions for p, we must have n(p) = p = kpk. Since dN = n(q)dA where dA is the surface area form, we have

We now show how to lower bounding the integral by showing hq; $\frac{p}{kDk}$ + (n(q) n(p))i is lower bounded.

Let c(t) be a minimal unit-speed geodesic on M := (X @S) = from p to q. Note that M = @S = from p to q. Note that M = @S = from p to q. Note that M = @S = from p to q. Note that M = @S = from p to q. Note that M = @S = from p to q.

hp; qi = hp; pi +
$$\sum_{0}^{Z_{1}} hp; c^{0}(t)idt = hp; pi + \sum_{0}^{Z_{1}} hProj_{T_{c(t)}} p; c^{0}(t)idt$$

where $T_{c(t)}$ is the tangent space to M at the point c(t). Hence by the Cauchy-Schwarz inequality we have

jhp; qi hp; pi
$$\sum_{0}^{Z_{1}} k \operatorname{Proj}_{T_{c(t)}} pkkc^{0}(t)kdt$$
:

By Proposition 6.3 of Niyogi et al. [2008], we have that for t the angle between the tangent spaces T_p and $T_{c(t)}$ that

$$\cos_t 1 = \frac{1}{M}(p; c(t)) = 1 = \frac{t}{M}(p; q)$$
: (25)

Since $\sin^2 t + \cos_2 t = 1$ and p is orthogonal to the tangent space at T_p , it follows that

hence

$$Z_{1}$$
 $k \operatorname{Proj}_{T_{c(t)}} \operatorname{pkkc}^{0}(t) \operatorname{kdt} (2=3) \operatorname{kpk}^{p} (2=M) d_{M}(p;q)^{3=2} + \operatorname{kpk}(1=2_{M}) d_{M}(p;q)^{2}$:

Since kp qk 2, provided that $_{\rm M}$ > 16 we have by Proposition 6.3 of Niyogi et al. [2008] that

$$d_{M}(p;q)_{M}(1)^{p}_{1}-\frac{2kp^{-}qk=_{M})}{2kp^{-}qk=_{M})}$$

Combining, we have for some absolute constant C > 0 that

hp; qi hp; pi(1
$$C^{p}1=M^{-}C=M$$
):

Also, we can compute

$$kn(q)$$
 $n(p)k = p \frac{r}{2 - 2\cos 1} \frac{r}{2 - - \frac{r}{2\cos (p;q)}} - 8$

SO

Hence provided $_{\rm M}$ > ${\rm C}^{\rm 0}$ for some absolute constant ${\rm C}^{\rm 0}$ > 0 and kpk > 0:1, we have

$$Z = \frac{2 (q)}{(1 kqk^{2})^{2}} hq \frac{p}{kpk} + (n(q) n(p))idA$$

$$Z = \frac{(q)}{(1 kqk^{2})^{2}} kpkdA$$

$$Z = \frac{(q)}{kqk^{2})^{2}} kpkdA$$

using that the integrand on the left is always negative. We can further lower bound the integral by considering the intersection of M with a ball of radius $r:=\frac{1-kpk}{2}$ centered at p. We have

$$Z = \frac{(q)}{q^{2B(0;1)\setminus(X \otimes S)}} (q) + kpkdA = \frac{Z}{q^{2B(p;r)\setminus M}} (1 - \frac{(q)}{kqk^{2})^{2}} kpkdA + kpk(cos)^{k}vol(B^{k}(p;r)) = kpk(cos)^{k}r^{k} = \frac{(q)}{q^{2B(p;r)\setminus M}} \frac{(q)}{(1 - kqk^{2})^{2}}$$

$$= kpk(cos)^{k}r^{k} = \frac{B_{k}(q)}{(1 - kqk^{2})^{2}}$$

where k = d 1 is the dimension of M and = $\arcsin(r=2)$ and we applied Lemma 5.3 of Niyogi et al. [2008]. If kpk 2 (0:1; 0:9) this is lower bounded by a constant $C_k > 0$ which is at worst exponentially small in k.

Hence recalling (24) we have for any X with d(X; @S) 2 (0:1; 0:9) and for sufficiently small so that $8^{k+1} k(JF_1)_X k_{OP} kk < C_k = 4$ for any such X, we have that

$$(rF_2)_{\bar{K}}(JF)_{\bar{K}} + F_2^2 + C^0$$

where $C_k^{\,\,c}>0$ is a constant that is at worst exponentially small in k. Therefore

$$E_{X,d}(X;@S)^{2}[;] = (r^{2}F_{X})^{T}(J_{K})^{T} + F_{2}^{2} + C^{0}\frac{Pr(d(X;@S)^{2}(0:1;0:9))}{Pr(d(X;@S)^{2}(0:1;0:9))}$$
:

Combining these estimates, we have for some constant $C_k^{00} > 0$ which is at worst exponentially small in k and sufficiently small (to satisfy the conditions above, including the requirement $M > C_k^{00}$) that

$$_{\text{max}}(s_{M}) = \frac{C_{k}^{00} \Pr(d(X;@S) \ 2 \ (0:1;0:9))}{\Pr(d(X;@S))^{2}}$$
: (26)

ļ

Observe that for any points x; y and = (1; 0) we have by the mean value theorem that

$$p(x)=p(y) = \exp(h_1; F_1(x) \quad F_1(y)) \exp kk \sup_{2[0;1]} k(JF_1)_{x+(1)y} k_{OP} kx \quad yk :$$
 (27)

so the log of the density is Lipschitz. This basically reduces estimating Pr(d(X; @S)) for small to understanding the volume of tubes around @S, which can be done using the same ideas as the proof of Weyl's tube formula [Weyl, 1939, Gray, 2003].

Lemma 6 (Proposition 6.1 of Niyogi et al. [2008]). Let M be a smooth and compact submanifold of dimension q in R^d . At a point p 2 M let B : T_p T_p ! T? dengte the second fundamental form, and for a unit normal vector u, let L_u be the linear operator defined so that hu; B(v; w)i = hv; L_uwi (this matches the notation from Niyogi et al. [2008]). Then

Lemma 7 (Lemma 3.14 of Gray [2003]). Let M be a smooth and compact submanifold of dimension q in R^d . Let exp_p denote the exponential map from the normal bundle at p. The Jacobian determinant of the map

$$M \ (1=_M; 1=_M) \ S_{d \ q \ 1} \ ! \ R^d; \qquad (p;t;u) \ ! \ exp_p(tu)$$

is $det(I tL_u)$.

We can compute

where in the second equality we performed a change of variables and obtained the result by applying Lemma 7. We have

$$det(I tL_{II}) 2 [(1 t=)^k; (1+t=)^k]$$

and so applying (27) we find that if we define $c := kk \sup_{x:d(x;@S)} k(JF_1)_x k_{OP}$ which can be made arbitrarily small by taking sufficiently small, then

$$Pr(d(X; @S) r) 2 [2e^{c}(1 =)^{k}V; 2e^{c}(1 + =)^{k}V]$$
 (28)

where

$$V := p(x)dA: @s$$

Note that $(1+=)^k$ $e^{k=}$ and $(1=)^k$ exp(O(k=)) provided that ==O(1=k). Since $Pr(d(X;@S) \ 2 \ (0:1;0:9)) = Pr(d(X;@S) < 0:9)$ $Pr(d(X;@S) \ 0:1)$ and the distribution we consider has a density, by combining (28) and (26) we find that for sufficiently small we have

$$_{\text{max}}(_{\text{SM}}) C^{000} + \frac{1}{k} \frac{1}{_{@S}p(x)dA}$$

where C_k^{000} is at worst exponentially small in k.

B.2 Relating Fisher matrices of augmented and original sufficient statistics

Next, we show that adding the extra sufficient statistic F_2 has a comparatively minor effect on the efficiency of MLE. Intuitively, to be able to estimate the coefficient of F_2 correctly we just need: (1) the variance of F_2 is large, so that a nonzero coefficient of F_2 can be observed from samples (e.g. when F_2 encodes the cut S, the coefficient can be estimated by looking at the relative weight between S and S^C), and (2) there is no redundancy in the sufficient statistics, e.g. $F_2 = F_1$ since otherwise different coefficients can encode the same distribution. The proof of this uses that the inverse covariance of the MLE has a simple explicit form (the Fisher information, which is the covariance matrix of $(F_1; F_2)$), and conditions (1) and (2) naturally appear when we use this fact.

Quantitatively, we show:

Lemma 8. Suppose that $F = (F_1; F_2)$ is a random vector valued in R^{m+1} with F_1 valued in R^m and F_2 valued in R. Suppose that F_2 is not in the affine of linear combinations of the coordinates of F_1 , i.e. for all $w_1 \ 2 \ R_m$ there exists > 0 such that

$$Cov(hw_1; F_1i; F_2)^2 Var(hw_1; F_1i)Var(F_2)$$
:

Then we have the lower bound

$$_{\mathsf{F}}$$
 (1) $\overset{\mathsf{F}_{\scriptscriptstyle{1}}}{\overset{\mathsf{F}}{\overset{\mathsf{O}}{\mathsf{Var}}}}\mathsf{Var}(\mathsf{F}_{\mathsf{2}})$

in the standard PSD (positive semidefinite) order.

Proof. To show a lower bound on

$$F$$
_ F_1 F_2 F_1 F_1 F_5 2

observe that

$$hw;_Fwi = hw_1;_{F_1}w_1i + 2w_2hw_1;_{F_1F_2}i + w_2F_2$$

so under the assumption we have by the AM-GM inequality that

$$hw_{;F}w_{i}$$
 (1) $[hw_{1}; F_{1}w_{1}i + w_{2}F_{2}]$ 2

and hence $_{\rm F}$ is lower bounded in the PSD order as long as $_{\rm F_1}$ is and $_{\rm F_2}$ is.

The lower bound on $Var(F_2)$ is guaranteed when F_2 corresponds to a cut with large mass on both sides since the variance of F_2 is lower bounded by its variance conditioned on being away from the boundary of S.

B.3 Putting together

Finally, given Lemma 5 and 8, we can complete the proof of Theorem 3.

Proof of Theorem 3. Define = Pr(X 2 S) for the purpose of this proof. Observe that by (28)

$$Var(1_S F_2) E(1_S F_2)^2 Pr(d(X;@S)) 4V where V$$

= $R_{@}$ p(x)dA. We have that

$$Cov(hw_1; F_1i; F_2) = Cov(hw_1; F_1i; 1_S) + Cov(hw_1; F_1i; F_2 1_S)$$

so if w_1 is arbitrary and normalized so that $Var(hw_1; F_1i) = 1$ then we have

Therefore provided > 0 we have

$$f^{1} = 0^{1} - F^{1} = 0$$
 $Var(F_{2})^{-1}$:

On the other hand, by Lemma 5 we have

$$_{\text{max}}(s_{M})_{V}:\frac{c^{d}}{}$$

Hence there exists some w such that

$$\frac{\frac{2}{S\,M}\,(W)}{M^2\,L\,E\,(W)} - \frac{c^d}{maxfk_{F_1}\,k_{OP}^{-1}\,;\,1=(1-)g\,V} - \frac{1}{1+} \left(1-\frac{c^d}{)k_{F_1}\,k_{OP}} - V^{-1} - \frac{(1-)}{(1-)g\,V} + \frac{1}{2} \left(1-\frac{1}{2}\right) \left$$

Using that minf; 1 g=2 (1) 1=4 and dividing c by two gives the result.