

# NCE of Unnormalized Statistical Models , with Applications to Natural Image Statistics

**Michael U. Gutmann**

MICHAEL.GUTMANN@HELSINKI.FI

**Aapo Hyvärinen**

AAPO.HYVARINEN@HELSINKI.FI

*Department of Computer Science*

*Department of Mathematics and Statistics*

*Helsinki Institute for Information Technology HIIT*

*University of Helsinki, Finland*

**Editor:** Yoshua Bengio

## Abstract

We consider the task of estimating, from observed data, a probabilistic model that is parameterized by a finite number of parameters. In particular, we are considering the situation where the model pdf is unnormalized . That is, the model is only specified up to the partition function . The partition function normalizes a model so that it integrates to one for any choice of the parameters. However, it is often impossible to obtain it in closed form. Gibbs distributions, Markov and multi-layer networks are examples of models where analytical normalization is often impossible . MLE can then not be used without resorting to numerical approximations which are often computationally expensive . We propose here a new objective function for the estimation of both normalized and unnormalized models . The basic idea is to perform nonlinear logistic regression to discriminate between the observed data and some artificially generated noise. With this approach , the normalizing partition function can be estimated like any other parameter. We prove that the new estimation method leads to a consistent (convergent ) estimator of the parameters . For large noise sample sizes , the new estimator is furthermore shown to behave like the mle. In the estimation of unnormalized models , there is a trade -off between statistical and computational performance . We show that the new method strikes a competitive trade -off in comparison to other estimation methods for unnormalized models. As an application to real data, we estimate novel two-layer models of natural image statistics with spline nonlinearities.

**Keywords:** unnormalized models, partition function, computation, estimation, natural image statistics

## 1. Introduction

This paper is about parametric density estimation, where the general setup is as follows. A sample  $X=(\mathbf{x}_1, \dots, \mathbf{x}_{T_d})$  of a random vector  $\mathbf{x} \in \mathbb{R}^n$  is observed which follows an unknown pdf  $p_d$ . The data-pdf  $p_d$  is modeled by a parameterized family of functions  $\{p_m(\cdot; \boldsymbol{\theta})\}_{\boldsymbol{\theta}}$  where  $\boldsymbol{\theta}$  is a vector of parameters . It is commonly assumed that  $p_d$  belongs to this family . In other words,  $p_d(\cdot) = p_m(\cdot; \boldsymbol{\theta}^*)$  for some parameter  $\boldsymbol{\theta}^*$ . The parametric density estimation problem is then about finding  $\boldsymbol{\theta}^*$  from the observed sample  $X$ .

Any estimate  $\hat{\theta}$  must yield a properly normalized pdf  $p_m(\cdot; \hat{\theta})$  :

$$\int p_m(\mathbf{u}; \hat{\theta}) d\mathbf{u} = 1, \quad p_m(\cdot; \hat{\theta}) \geq 0. \quad (1)$$

If the model  $p_m(\cdot; \theta)$  is such that the constraints hold for all  $\theta$ , and not only  $\hat{\theta}$ , we say that the model is **normalized**. If the model is specified such that the positivity constraint but not the normalization constraint is satisfied for all parameters, we say that the model is **unnormalized**. By assumption there is, however, at least one value of the parameters for which an unnormalized model integrates to one, namely  $\theta^*$ .

In order to highlight that a model, parameterized by  $\alpha$ , is unnormalized, we denote it by  $p_m^0(\cdot; \alpha)$ .

The **partition function**

$$Z(\alpha) = \int p_m^0(\mathbf{u}; \alpha) d\mathbf{u}, \quad (2)$$

can be used to convert an unnormalized model  $p_m^0(\cdot; \alpha)$  into a normalized one:  $p_m^0(\cdot; \alpha)/Z(\alpha)$  integrates to one for every value of  $\alpha$ . Examples of distributions which are often specified by means of an unnormalized model and the partition function are Gibbs distributions, Markov networks or multilayer networks. The function  $\alpha \mapsto Z(\alpha)$  is, however, defined via an integral. Unless  $p_m^0(\cdot; \alpha)$  has some particularly convenient form, the integral cannot be computed analytically so that the function  $Z(\alpha)$  is not available in closed form. For low-dimensional problems, numerical integration can be used to approximate the function  $Z(\alpha)$  to a very high accuracy but for high-dimensional problems this is computationally expensive. Our paper deals with density estimation in this case, that is, with density estimation when the computation of the partition function is analytically intractable and computationally expensive.

Several solutions for the estimation of unnormalized models which cannot be normalized in closed form have been suggested so far. Geyer (1994) proposed to approximate the calculation of the partition function by means of importance sampling and then to maximize the approximate log-likelihood (Monte Carlo maximum likelihood). Approximation of the gradient of the log-likelihood led to another estimation method (contrastive divergence by Hinton, 2002). Estimation of the parameter  $\alpha$  directly from an unnormalized model  $p_m^0(\cdot; \alpha)$  has been proposed by Hyvärinen (2005). This approach, called score matching, avoids the problematic integration to obtain the partition function altogether. All these methods need to balance the accuracy of the estimate and the time to compute the estimate.

In this paper,<sup>1</sup> we propose a new estimation method for unnormalized models. The idea is to consider  $Z$ , or  $c = \ln 1/Z$ , not any more as a function of  $\alpha$  but as an additional parameter of the model. That is, we extend the unnormalized model  $p_m^0(\cdot; \alpha)$  to include a normalizing parameter  $c$  and estimate

$$\ln p_m(\cdot; \theta) = \ln p_m^0(\cdot; \alpha) + c,$$

with parameter vector  $\theta = (\alpha, c)$ . The estimate  $\hat{\theta} = (\hat{\alpha}, \hat{c})$  is then such that the unnormalized model  $p_m^0(\cdot; \hat{\alpha})$  matches the shape of  $p_d$ , while  $\hat{c}$  provides the proper scaling so that Eq (1) holds.

---

<sup>1</sup> Preliminary versions were presented at AISTATS (Gutmann and Hyvärinen, 2010) and ICANN (Gutmann and Hyvärinen, 2009).

Unlike in the approach based on the partition function, we aim not at normalizing  $p_m^0(\cdot; \alpha)$  for all  $\alpha$  but only for  $\hat{\alpha}$ . This avoids the problematic integration in the definition of the partition function  $\alpha \mapsto Z(\alpha)$ . Such a separate estimation of shape and scale is, however, not possible for MLE. The reason is that the likelihood can be made arbitrarily large by setting the normalizing parameter  $c$  to larger and larger numbers. The new estimation method which we propose here is based on the maximization of a well defined objective function. There are no constraints in the optimization so that powerful optimization techniques can be employed. The intuition behind the new objective function is to learn to classify between the observed data and some artificially generated noise. We approach thus the density estimation problem, which is an unsupervised learning problem, via supervised learning. The new method relies on noise which the data is contrasted to, so that we will refer to it as “noise-contrastive estimation”.

The paper is organized in four main sections. In Section 2, we present NCE and prove fundamental statistical properties such as consistency. In Section 3, we validate and illustrate the derived properties on artificial data. We use artificial data also in Section 4 in order to compare the new method to the aforementioned estimation methods with respect to their statistical and computational efficiency. In Section 5, we apply NCE to real data. We estimate two-layer models of natural images and also learn the nonlinearities from the data. This section is fairly independent from the other ones. The reader who wants to focus on natural image statistics may not need to go first through the previous sections. On the other hand, the reader whose interest is in estimation theory only can skip this section without missing pieces of the theory although the section provides, using real data, a further illustration of the workings of unnormalized models and the new estimation method. Section 6 concludes the paper.

## 2. NCE

This section presents the theory of NCE. In Section 2.1, we motivate NCE and relate it to supervised learning. The definition of NCE is given in Section 2.2. In Section 2.3, we prove that the estimator is consistent for both normalized and unnormalized models, and derive its asymptotic distribution. In Section 2.4, we discuss practical aspects of the estimator and show that, in some limiting case, the estimator performs as well as MLE.

### 2.1 Density Estimation by Comparison

Density estimation is much about characterizing properties of the observed data  $X$ . A convenient way to describe properties is to describe them relative to the properties of some reference data  $Y$ . Assume that the reference (noise) data  $Y$  is an i.i.d. sample  $(y_1, \dots, y_{T_n})$  of a rv  $y \in \mathbb{R}^n$  with pdf  $p_n$ . A relative description of the data  $X$  is then given by the ratio  $p_d/p_n$  of the two density functions. If the reference distribution  $p_n$  is known, one can, of course, obtain  $p_d$  from the ratio  $p_d/p_n$ . In other words, if one knows the differences between  $X$  and  $Y$ , and also the properties of  $Y$ , one can deduce from the differences the properties of  $X$ .

Comparison between two data sets can be performed via classification: In order to discriminate between two data sets, the classifier needs to compare their properties. In the following, we show that training a classifier based on logistic regression provides a relative description of  $X$  in the form of an estimate of the ratio  $p_d/p_n$ .

augmented data

Denote by  $U = (\mathbf{u}_1, \dots, \mathbf{u}_{T_d+T_n})$  the union of  $X$  and  $Y$ , and assign to each data point  $\mathbf{u}_t$  a binary class label  $C_t$ :  $C_t = 1$  if  $\mathbf{u}_t \in X$  and  $C_t = 0$  if  $\mathbf{u}_t \in Y$ . In logistic regression, the posterior probabilities of the classes given the data are estimated. As the pdf  $p_d$  of the data  $\mathbf{x}$  is unknown, we model the class-conditional probability  $p(\cdot | C = 1)$  with  $p_m(\cdot; \boldsymbol{\theta})$ .<sup>2</sup> The class-conditional probability densities are

augmented distr.

$$p(\mathbf{u} | C = 1; \boldsymbol{\theta}) = p_m(\mathbf{u}; \boldsymbol{\theta}), \quad p(\mathbf{u} | C = 0) = p_n(\mathbf{u}).$$

The prior probabilities are  $P(C = 1) = T_d / (T_d + T_n)$  and  $P(C = 0) = T_n / (T_d + T_n)$ . The posterior probabilities for the classes are therefore

$$P(C = 1 | \mathbf{u}; \boldsymbol{\theta}) = \frac{p_m(\mathbf{u}; \boldsymbol{\theta})}{p_m(\mathbf{u}; \boldsymbol{\theta}) + \nu p_n(\mathbf{u})}, \quad P(C = 0 | \mathbf{u}; \boldsymbol{\theta}) = \frac{\nu p_n(\mathbf{u})}{p_m(\mathbf{u}; \boldsymbol{\theta}) + \nu p_n(\mathbf{u})}, \quad (3)$$

where  $\nu := P(C = 0) / P(C = 1) \sim T_n / T_d$ .  $h(\mathbf{u}; \boldsymbol{\theta}) := P(C = 1 | \mathbf{u}; \boldsymbol{\theta})$  0/1-classifier of data  $U$

$$G(\mathbf{u}; \boldsymbol{\theta}) := \ln p_m(\mathbf{u}; \boldsymbol{\theta}) - \ln p_n(\mathbf{u}), \quad \text{model-noise LR} \quad (4)$$

==>

$$h(\mathbf{u}; \boldsymbol{\theta}) = r_\nu(G(\mathbf{u}; \boldsymbol{\theta})), \quad (5)$$

where

$$r_\nu(u) = \frac{1}{1 + \nu \exp(-u)} = \text{expit}(u - \ln \nu) \quad (6)$$

The conditional log-likelihood of aug. dstr.  $p(C | \mathbf{u}, \boldsymbol{\theta})$   
not joint log-likelihood!!!

$$\ell(\boldsymbol{\theta}) = \sum_{t=1}^{T_d} \ln [h(\mathbf{x}_t; \boldsymbol{\theta})] + \sum_{t=1}^{T_n} \ln [1 - h(\mathbf{y}_t; \boldsymbol{\theta})]. \quad = - \sum_t \mathbf{Hb}(\mathbf{u}_t, \mathbf{h}(\mathbf{x}_t, \boldsymbol{\Theta})) \quad (7)$$

Optimizing  $\ell(\boldsymbol{\theta})$  wrt  $\boldsymbol{\theta}$  leads to an estimate  $G(\cdot; \hat{\boldsymbol{\theta}})$  of the log-ratio  $\ln(p_d/p_n)$ . That is, an approximate description of  $X$  relative to  $Y$  can be obtained by optimization of Eq (7). The sign-flipped objective function,  $-\ell(\boldsymbol{\theta})$ , is the cross-entropy error function (Bishop, 1995).

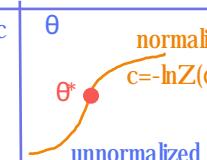
Density estimation, an unsupervised learning problem, can be performed by logistic regression, that is, supervised learning. [Hastie et al. \(2009, Chapter 14.2.4, pp. 495–497\)](#).

---

2. Classically,  $p_m(\cdot; \boldsymbol{\theta})$  would, in the context of this section, be a normalized pdf. In our paper, however,  $\boldsymbol{\theta}$  may include a parameter for the normalization of the model.

## 2.2 Definition of the Estimator

Given an unnormalized statistical model  $p_m^0(\cdot; \boldsymbol{\alpha})$ , we include for normalization an additional parameter  $c$  into the model:



$$\ln p_m(\cdot; \boldsymbol{\theta}) = \ln p_m^0(\cdot; \boldsymbol{\alpha}) + c,$$

where  $\boldsymbol{\theta} = (\boldsymbol{\alpha}, c)$ . The parameter  $c$  scales the unnormalized model  $p_m^0(\cdot; \boldsymbol{\alpha})$  so that Eq (1) can be fulfilled. After learning,  $\hat{c}$  provides an estimate for  $\ln 1/Z(\boldsymbol{\alpha})$ .

$X = (\mathbf{x}_1, \dots, \mathbf{x}_{T_d})$ : the observed data set that consists of  $T_d$  independent

observations of  $\mathbf{x} \in \mathbb{R}^n$ .  $Y = (\mathbf{y}_1, \dots, \mathbf{y}_{T_n})$ : an artificially generated data set that consists of  $T_n = \nu T_d$  independent observations of noise  $\mathbf{y} \in \mathbb{R}^n$  with known distribution  $p_n$ . The estimator  $\hat{\boldsymbol{\theta}}_T$  maximizes

$$J_T(\boldsymbol{\theta}) = \frac{1}{T_d} \left\{ \sum_{t=1}^{T_d} \ln [h(\mathbf{x}_t; \boldsymbol{\theta})] + \sum_{t=1}^{T_n} \ln [1 - h(\mathbf{y}_t; \boldsymbol{\theta})] \right\}, \quad (8)$$

where  $h(\cdot; \boldsymbol{\theta})$  was defined in (5).  $J_T$  is, up to the division by  $T_d$ , the log-likelihood in (7)

$$J_T(\boldsymbol{\theta}) = \frac{1}{T_d} \sum_{t=1}^{T_d} \ln [h(\mathbf{x}_t; \boldsymbol{\theta})] + \nu \frac{1}{T_n} \sum_{t=1}^{T_n} \ln [1 - h(\mathbf{y}_t; \boldsymbol{\theta})]. \quad (9)$$

$h(\cdot; \boldsymbol{\theta}) \rightarrow 0$ ,  $G(\cdot; \boldsymbol{\theta}) \rightarrow -\infty$ ;  $\rightarrow 1$ ,  $G(\cdot; \boldsymbol{\theta}) \rightarrow \infty$ .  $J_T \leq 0$ ,  $= 0$ , for all  $t$ ,  $h(\mathbf{x}_t; \boldsymbol{\theta})/h(\mathbf{y}_t; \boldsymbol{\theta}) \rightarrow 1/0$ , respectively. Intuitively, this means that logistic regression has learned to discriminate between the two sets as well as possible.

## 2.3 Properties of the Estimator

We characterize here the behavior of the estimator  $\hat{\boldsymbol{\theta}}_T$  for large sample sizes  $T_d$  and fixed ratio  $\nu$ . Since  $\nu$  is kept fixed,  $T_n = \nu T_d$  will also increase as  $T_d$  increases. The weak law of large numbers shows that as  $T_d$  increases  $J_T(\boldsymbol{\theta})$  converges in probability to  $J$ ,

$$J(\boldsymbol{\theta}) := E \{ \ln [h(\mathbf{x}; \boldsymbol{\theta})] \} + \nu E \{ \ln [1 - h(\mathbf{y}; \boldsymbol{\theta})] \}. \quad (10)$$

the objective  $J$  seen as a function of  $f_m(\cdot) = \ln p_m(\cdot; \boldsymbol{\theta})$ ,

$$\tilde{J}(f_m) := E \{ \ln [r_\nu(f_m(\mathbf{x}) - \ln p_n(\mathbf{x}))] \} + \nu E \{ \ln [1 - r_\nu(f_m(\mathbf{y}) - \ln p_n(\mathbf{y}))] \} \quad (11)$$

We start the characterization of the estimator  $\hat{\boldsymbol{\theta}}_T$  by describing the optimization landscape for  $f_m$ . The following theorem shows that the data-pdf  $p_d$  can be found by maximization of  $\tilde{J}$ , that is by learning a nonparametric classifier under the ideal situation of an infinite amount of data.

**Theorem 1 (Nonparametric estimation)**  $\tilde{J}$  attains a maximum at  $f_m = \ln p_d$ . There are no other extrema if  $p_n \gg p_d$ .

A fundamental point in the theorem is that the maximization is performed without any normalization constraint for  $f_m$ . This is in stark contrast to MLE, where  $\exp(f_m)$  must integrate to one. With our objective function, no such constraints are necessary. The maximizing pdf is found to have unit integral automatically.

The positivity condition for  $p_n$  in the theorem tells us that the data-pdf  $p_d$  cannot be inferred at regions in the data space where there are no contrastive noise samples. For example, the estimation of a pdf  $p_d$  which is nonzero only on the positive real line by means of a noise distribution  $p_n$  that has its support on the negative real line is impossible. The positivity condition can be easily fulfilled by taking, for example, a Gaussian as cnd.

In practice, the amount of data is limited and a finite number of parameters  $\theta \in \mathbb{R}^m$  specify  $p_m(\cdot; \theta)$ . This has two consequences for any estimation method that is based on optimization: First, it restricts the space where the data-pdf  $p_d$  is searched for. Second, it may introduce local maxima into the optimization landscape. For the characterization of the estimator in this situation, it is normally assumed that  $p_d$  follows the model, so that there is a  $\theta^*$  with  $p_d(\cdot) = p_m(\cdot; \theta^*)$ .

Our second theorem shows that  $\hat{\theta}_T$ , the value of  $\theta$  which (globally) maximizes  $J_T$ , converges to  $\theta^*$ . The correct estimate of  $p_d$  is thus obtained as the sample size  $T_d$  increases. For unnormalized models, the conclusion of the theorem is that maximization of  $J_T$  leads to the correct estimates for both the parameter  $\alpha$  in the unnormalized pdf  $p_m^0(\cdot; \alpha)$  and the normalizing parameter  $c$ .

**Theorem 2 (Consistency)** *If conditions (a) to (c) are fulfilled then  $\hat{\theta}_T$  converges in probability to  $\theta^*$*

$$(a) p_n >> p_d$$

$$(b) \sup_{\theta} |J_T(\theta) - J(\theta)| \xrightarrow{P} 0$$

(c) The matrix  $\mathcal{I}_{\nu} = \int g(u)g(u)^T P_{\nu}(u)p_d(u)du$  has full rank, where

$$g(u) = \nabla_{\theta} \ln p_m(u; \theta)|_{\theta^*}, \quad P_{\nu}(u) = \frac{\nu p_n(u)}{p_d(u) + \nu p_n(u)}.$$


---

The proof is given in Appendix A.3. Condition (a) is inherited from Theorem 1. Conditions (b) and (c) have their counterparts in MLE (see for example Wasserman, 2004, Theorem 9.13): We need in (b) uniform convergence in probability of  $J_T$  to  $J$ ; in MLE, uniform convergence of the log-likelihood to the Kullback-Leibler divergence is required likewise. Condition (c) assures that for large sample sizes, the objective function  $J_T$  becomes peaked enough around the true value  $\theta^*$ . This imposes a constraint on the model  $p_m(\cdot; \theta)$  via the vector  $g$ . A similar constraint is required in MLE.

The next theorem describes the distribution of the estimation error  $(\hat{\theta}_T - \theta^*)$  for large sample sizes. The proof is given in Appendix A.4.

**Theorem 3 (Asymptotic normality)**  $\sqrt{T_d}(\hat{\theta}_T - \theta^*)$  is asymptotically normal with mean zero and covariance matrix  $\Sigma$ ,

$$\Sigma = \mathcal{I}_{\nu}^{-1} - \left(1 + \frac{1}{\nu}\right) \mathcal{I}_{\nu}^{-1} E(P_{\nu}g) E(P_{\nu}g)^T \mathcal{I}_{\nu}^{-1},$$

where  $E(P_{\nu}g) = \int P_{\nu}(u)g(u)p_d(u)du$ .

**Corollary 4** For large sample sizes  $T_d$ , the mse  $E\left(||\hat{\theta}_T - \theta^*||^2\right)$  equals  $\text{tr}(\Sigma)/T_d$ .

■

## 2.4 Choosing the Noise

Theorem 3 shows that the noise distribution  $p_n$  and the ratio  $\nu = T_n/T_d$  have an influence on the accuracy of the estimate  $\hat{\theta}_T$ . A natural question to ask is what, from a statistical standpoint, the best choice of  $p_n$  and  $\nu$  is. Our result on consistency (Theorem 2) also includes a technical constraint for  $p_n$  but this one is so mild that many distributions will satisfy it.

Theorem 2 shows that, for a given samples size  $T_d$ ,  $P_\nu \rightarrow 1$  as  $T_n$  is made larger and larger. This implies that for large  $\nu$ , the covariance matrix  $\Sigma$  does not depend on the choice of  $p_n$ .

**Corollary 5** For  $\nu \rightarrow \infty$ ,  $\Sigma$  is independent of the choice of  $p_n$  and

$$\Sigma = \mathcal{I}^{-1} - \mathcal{I}^{-1} E(\mathbf{g}) E(\mathbf{g})^T \mathcal{I}^{-1},$$

where  $E(\mathbf{g}) = \int \mathbf{g}(\mathbf{u}) p_d(\mathbf{u}) d\mathbf{u}$  and  $\mathcal{I} = \int \mathbf{g}(\mathbf{u}) \mathbf{g}(\mathbf{u})^T p_d(\mathbf{u}) d\mathbf{u}$ .

The asymptotic distribution of the estimation error becomes thus independent from  $p_n$ . Hence, as the size of the contrastive -noise sample  $Y$  increases , the choice of the cnd becomes less and less important. Moreover, for normalized models, we have the result that the estimation error has the same distribution as the estimation error in MLE.

**Corollary 6** For normalized models, NCE is,  $\nu \rightarrow \infty$ , asymptotically Fisher-efficient for all  $p_n$ .

**Proof** For normalized models, no normalizing parameter  $c$  is needed. In Corollary 5, the function  $\mathbf{g}$  is then the score function as in MLE, and the matrix  $\mathcal{I}$  is the FIM. Since  $E(\mathbf{g})=0$ ,  $\Sigma$  is the inverse of the FIM.

■

The corollaries above give one answer to the question on how to choose the noise distribution  $p_n$  and the ratio  $\nu$ : If  $\nu$  is made large enough, the actual choice of  $p_n$  is not of great importance. Note that this answer considers only estimation accuracy and ignores the computational load associated with the processing of noise. In Section 4, we will analyze the trade-off between estimation accuracy and computation time.

given  $\nu$ , one could try to find the noise distribution which minimizes MSE  $E ||\hat{\theta}_T - \theta^*||^2$ . However, this minimization turns out to be quite difficult. Intuitively, one could think that a good candidate for the noise distribution  $p_n$  is a distribution which is close to the data distribution  $p_d$ . If  $p_n$  is too different from  $p_d$ , the classification problem might be too easy and would not require the system to learn much about the structure of the data. This intuition is partly justified by the following theoretical result:

**Corollary 7** If  $p_n = p_d$  then  $\Sigma = \left(1 + \frac{1}{\nu}\right) \left(\mathcal{I}^{-1} - \mathcal{I}^{-1} E(\mathbf{g}) E(\mathbf{g})^T \mathcal{I}^{-1}\right)$ . ■

For normalized models, for  $\nu=1$ ,  $\Sigma$  is two times the inverse of the FIM, and that for  $\nu=10$ , the ratio is already down to 1.1. For a noise distribution that is close to the data distribution, we have thus even for moderate values of  $\nu$  some guarantee that the MSE is reasonably close to the theoretical optimum.

To get estimates with a small estimation error, the foregoing discussion suggests the following

1. Choose noise for which an analytical expression for  $\ln p_n$  is available.
2. Choose noise that can be sampled easily.
3. Choose noise that is in some aspect, for example wrt its covariance structure, similar to the data.
4. Make the noise sample size as large as computationally possible.

Suitable noise distributions : Gaussian , Gaussian mixture , or ICA distributions . Uniform distributions are also suitable as long as their support includes the support of the data distribution so that condition (a) in Theorem 2 holds.

### 3. Simulations to Validate and Illustrate the Theory

In this section,3 we validate and illustrate the theoretical properties of NCE. In Section 3.1, we focus on the consistency of the estimator . In Section 3.2, we validate our theoretical results on the distribution of the estimation error, and investigate its dependency on the ratio  $\nu$  between noise and data sample size. In Section 3.3, we study how the performance of the estimator scales with the dimension of the data.

#### 3.1 Consistency

For the illustration of consistency, we estimate the parameters of a zero mean mv Gaussian

$$\ln p_d(\mathbf{x}) = -\frac{1}{2} \mathbf{x}^T \boldsymbol{\Lambda}^* \mathbf{x} + c^*, \quad c^* = \left( +\frac{1}{2} \ln |\det \boldsymbol{\Lambda}^*| - \frac{n}{2} \ln(2\pi) \right), \quad (12)$$

where  $c^*$  normalizes  $p_d$  to integrate to one. The precision matrix  $\boldsymbol{\Lambda}^*$ . The dim of  $\mathbf{x}$  is  $n=5$ .

As we are mostly interested estimation of unnormalized models

$$\ln p_m^0(\mathbf{x}; \boldsymbol{\alpha}) = -\frac{1}{2} \mathbf{x}^T \boldsymbol{\Lambda} \mathbf{x}$$

without knowing how to normalize it in closed form. This unnormalized model is a pairwise Markov network with quadratic node and edge potentials (Koller and Friedman , 2009 , Chapter 7). The parameter vector  $\boldsymbol{\alpha} \in \mathbb{R}^{15}$  contains the coefficients of the lower-triangular part of  $\boldsymbol{\Lambda}$  as the matrix is symmetric. For NCE, we add an additional normalizing parameter  $c$  to the model.

---

3. Matlab code for this and the other sections can be downloaded from the homepage of the first author.

The model that we estimate is thus

$$\ln p_m(\mathbf{x}; \boldsymbol{\theta}) = \ln p_m^0(\mathbf{x}; \boldsymbol{\alpha}) + c.$$

The model has 16 parameters given by  $\boldsymbol{\theta} = (\boldsymbol{\alpha}, c)$ . They are estimated by maximization of the objective function  $J_T(\boldsymbol{\theta})$  in Eq (8). We used  $N(0, 1)$  for  $p_n$ . The optimization was performed with the nonlinear conjugate gradient algorithm of Rasmussen (2006).

### 3.1.1 RESULTS

The presented results are an average over 500 estimation problems where the true precision matrix  $\Lambda^*$  was drawn at random with the condition number being controlled to be smaller than ten. The sampling of  $\Lambda^*$  was performed by randomly sampling its eigenvalues and eigenvectors: We drew the eigenvalues from an uniform distribution on the interval [0.1 0.9]. The orthonormal matrix  $\mathbf{E}$  with the eigenvectors was created by orthogonally projecting a matrix  $\mathbf{M}$  with elements drawn independently from a standard Gaussian onto the set of orthonormal matrices:  $\mathbf{E} = (\mathbf{MM}^T)^{-1/2} \mathbf{M}$ .

Figure 1(a) and (b) show the MSE for  $\boldsymbol{\alpha}$ , which contains the elements of the precision matrix  $\Lambda$ , and the normalizing parameter  $c$ , respectively. The MSE as a function of the data sample size  $T_d$  decays linearly on a log-log scale. This illustrates our result of consistency of the estimator, stated as Theorem 2, as convergence in quadratic mean implies convergence in probability. The plots also show that taking more noise samples  $T_n$  than data samples  $T_d$  leads to more and more accurate estimates. The performance for nce with  $\nu = T_n / T_d$  equal to one is shown in blue with circles as markers. For that value of  $\nu$ , there is a clear difference compared to MLE (black triangles in Figure 1(a)). However, the accuracy of the estimate improves strongly for  $\nu = 5$  (green squares) or  $\nu = 10$  (red diamonds) where the performance is rather close to the performance of MLE.

Another way to visualize the results is by showing the K-L divergences between the 500 true and estimated distributions. Figure 2 shows boxplots of the divergences for  $\nu = 1$  (blue) and  $\nu = 10$  (red). The results for MLE are shown in black. In line with the visualization in Figure 1, the estimated distribution becomes closer to the true distribution as the sample size increases. Moreover, the divergences become clearly smaller as  $\nu$  is increased from 1 to 10.

For unnormalized models, there is a subtlety in the computation of the divergence. With a validation set of size  $T_v$ , a sample version  $D_{KL}$  of the K-L div:

$$D_{KL} = \frac{1}{T_v} \sum_{t=1}^{T_v} \ln p_d(\mathbf{x}_t) - \left( \frac{1}{T_v} \sum_{t=1}^{T_v} \ln p_m^0(\mathbf{x}_t; \hat{\boldsymbol{\alpha}}) + \ln 1/Z(\hat{\boldsymbol{\alpha}}) \right).$$

The first term is the rescaled log-likelihood (average, sign-inverted log-loss) for the true distribution. The term in parentheses is the rescaled log-likelihood  $L$  of the estimated model. In the estimation of unnormalized models, we do not assume to know the mapping  $\boldsymbol{\alpha} \rightarrow Z(\boldsymbol{\alpha})$  so that  $L$  cannot be computed. With NCE, we can obtain an estimate,

$$\hat{L} = \frac{1}{T_v} \sum_{t=1}^{T_v} \ln p_m^0(\mathbf{x}_t; \hat{\boldsymbol{\alpha}}) + \hat{c}, \quad (13)$$

Figure 2(a) shows that the estimated  $D_{KL}$  is sometimes negative which means that  $\hat{L}$  is sometimes larger than the rescaled log-likelihood of the true distribution.

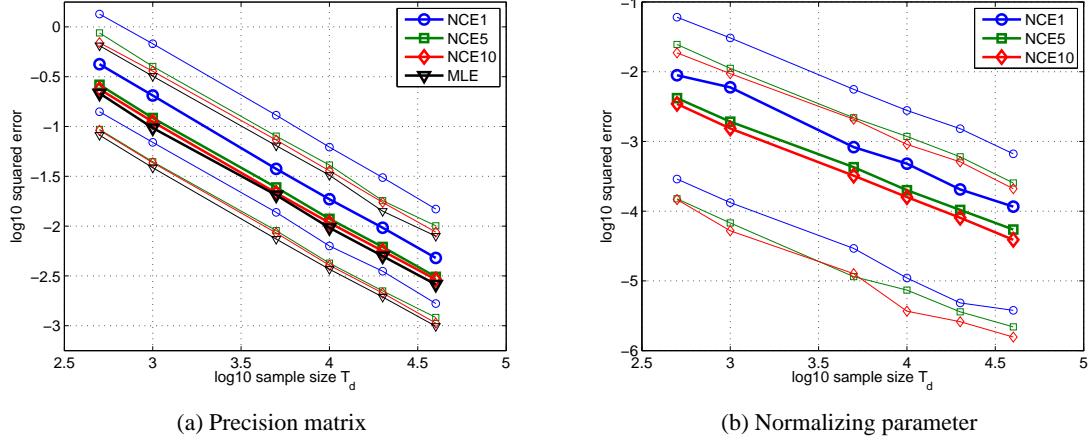


Figure 1: Validation of the theory of nce: Estimation errors for a 5D Gaussian distribution. Figures (a) and (b) show the MSE for the precision matrix  $\Lambda$  and the normalizing parameter  $c$ , respectively . The performance of NCE approaches the performance of MLE (black triangles ) as the ratio  $\nu = T_n / T_d$  increases : the case of  $\nu = 1$  is shown with blue circles ,  $\nu = 5$  with green squares , and  $\nu = 10$  with red diamonds . The thicker curves are the median of the performance for 500 random precision matrices with  $\text{cond} < 10$ . The finer curves show the 0.9 and 0.1 quantiles of the logarithm of the squared estimation error.

This happens because  $\hat{c}$  can be an over or underestimate of  $\ln 1/Z(\alpha^*)$ . This result follows from Figure 2(b) where we have computed  $D_{KL}$  with the analytical expression for  $\ln 1/Z(\alpha^*)$ , which is available for the Gaussian model considered here, see Eq (12).

### 3.2 Distribution of the Estimation Error

We validate and illustrate further properties of our estimator using the ICA model (see for example Hyvärinen et al., 2001b)

$$\mathbf{x} = \mathbf{A}\mathbf{s}. \quad (14)$$

In this subsection,  $n = 4$ , that is  $\mathbf{x} \in \mathbb{R}^4$ , and  $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_4)$  is a  $4 \times 4$  mixing matrix. The sources in the vector  $\mathbf{s} \in \mathbb{R}^4$  are identically distributed and independent from each other so that the data log-pdf  $\ln p_d$  is

$$\ln p_d(\mathbf{x}) = \sum_{i=1}^n f(\mathbf{b}_i^\star \mathbf{x}) + c^\star. \quad (15)$$

The  $i$ -th row of the matrix  $\mathbf{B}^* = \mathbf{A}^{-1}$  is denoted by  $\mathbf{b}_i^*$ . We consider here Laplacian sources of unit variance and zero mean. The nonlinearity  $f$  and the constant  $c^*$ , which normalizes  $p_d$  to integrate to one, are in this case given by

$$f(u) = -\sqrt{2}|u|, \quad c^* = \ln |\det \mathbf{B}^*| - \frac{n}{2} \ln 2. \quad (16)$$

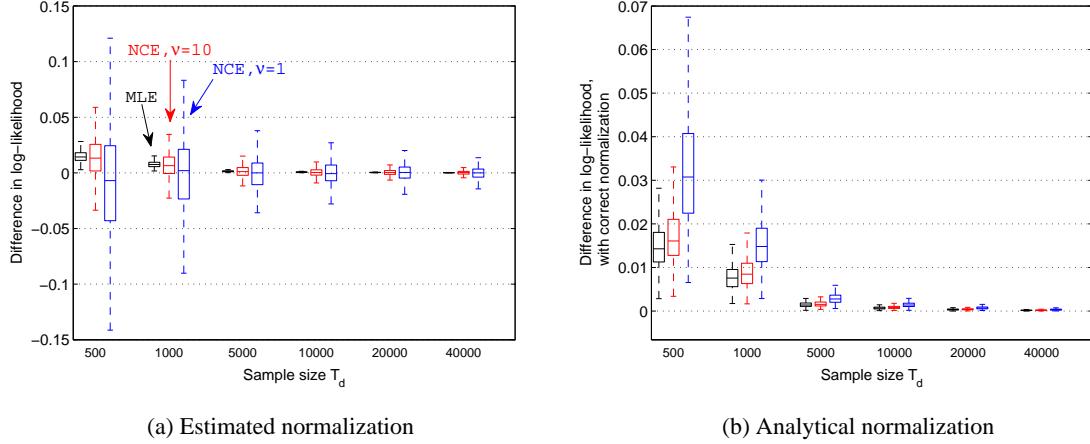


Figure 2: Validation of the theory of nce: Distributions of the KL divergences between the true and estimated 5D Gaussians . For each sample size, from left to right, the results for MLE are shown in black , the results for NCE with  $\nu = 10$  in red , and the results for  $\nu = 1$  in blue . The size  $T_v$  of the validation set was 100000 . For MLE, the results shown in Figures (a) and (b) are the same. For NCE, the divergences in Figure (a) were computed using the estimate  $\hat{c}$  of  $\ln 1/Z(\boldsymbol{\alpha}^{\hat{c}})$ . In Figure (b), the analytical expression for  $\ln 1/Z(\boldsymbol{\alpha}^{\hat{c}})$  was used.

As in Section 3.1, we apply NCE to the hypothetical situation where we want to estimate the unnormalized model

$$\ln p_m^0(\mathbf{x}; \boldsymbol{\alpha}) = \sum_{i=1}^n f(\mathbf{b}_i \mathbf{x}) \quad (17)$$

without knowing how to normalize it in closed form. The parameter vector  $\boldsymbol{\alpha} \in \mathbb{R}^{16}$  contains the elements of the row vectors  $\mathbf{b}_i$ . For NCE, we add an additional normalizing parameter  $c$  and estimate the model

$$\ln p_m(\mathbf{x}; \boldsymbol{\theta}) = \ln p_m^0(\mathbf{x}; \boldsymbol{\alpha}) + c,$$

with  $\boldsymbol{\theta} = (\boldsymbol{\alpha}, c)$ . As for the Gaussian case, we estimate  $\boldsymbol{\theta}$  by maximizing  $J_T(\boldsymbol{\theta})$  in Eq (8) with the nonlinear conjugate gradient algorithm of Rasmussen (2006). For the noise distribution  $p_n$ , we used a Gaussian distribution with covariance matrix given by the sample covariance of the data.

### 3.2.1 RESULTS

In Figures 3 and 4, we illustrate Theorem 2 on consistency and Theorem 3 on the asymptotic distribution of the estimator, as well as its corollaries. The results are averages over 500 random estimation problems. The mixing matrices  $\mathbf{A}$  were drawn at random by drawing their elements independently from a standard Gaussian and only accepting matrices which had a condition number smaller than ten.

Figure 3(a) and (b) show the MSE for  $\alpha$ , corresponding to the mixing matrix, and the normalizing parameter  $c$ , respectively. As illustrated for the Gaussian case in Figure 1, this figure visualizes the consistency of NCE. Furthermore, we see again that making  $\nu = T_n/T_d$  larger leads to a reduction of the error. The reduction gets, however, smaller as  $\nu$  increases. On average, changing  $\nu$  from 1 (red curve with asterisks as markers) to 10 (light blue squares) reduces the MSE for the mixing matrix by 53%; relative to  $\nu = 10$ ,  $\nu = 100$  (magenta diamonds) leads to a reduction of 18%. For  $c$ , the relative decrease in the MSE is 60% and 17%, respectively.

In Figure 4(a), we test the theoretical prediction of Corollary 4 that, for large samples sizes  $T_d$ , the MSE decays like  $\text{tr } \Sigma / T_d$ . The covariance matrix  $\Sigma$  can be numerically evaluated according to its definition in Theorem 3.<sup>4</sup> This allows for a prediction of the MSE that can be compared to the MSE obtained in the simulations. The figure shows that the MSE from the simulations (labelled “sim” in the figure) matches the prediction (“pred”) for large  $T_d$ . Furthermore, we see again that for large  $\nu$ , the performance of NCE is close to the performance of MLE. In other words, the trace of  $\Sigma$  is close to the trace of the FIM. Note that for clarity, we only show the curves for  $\nu \in \{0.1, 1, 100\}$ . The curve for  $\nu = 10$  was, as in Figure 3(a) and (b), very close to the curve for  $\nu = 100$ .

In Figure 4(b), we investigate how  $\text{tr } \Sigma$  (the asymptotic variance) depends on  $\nu$ . Note that the covariance matrix  $\Sigma$  includes terms related to the parameter  $c$ . The FIM includes, in contrast to  $\Sigma$ , only terms related to the mixing matrix. For better comparison with MLE, we show thus in the figure  $\text{tr } \Sigma$  both with the contribution of the normalizing parameter  $c$  (blue squares) and without (red circles). For the latter case, the reduced trace of  $\Sigma$ , which we will denote by  $\text{tr } \Sigma_B$ , approaches  $\text{tr } \text{FIM}$ . Corollary 6 stated that NCE is asymptotically Fisher-efficient for large  $\nu$  if the normalizing constant is not estimated. Here, we see that this result also approximately holds for our unnormalized model where the normalizing constant needs to be estimated.

Figure 4(c) gives further details to which extent the estimation becomes more difficult if the model is unnormalized. We computed numerically the asymptotic variance  $\text{tr } \tilde{\Sigma}$  if the model is correctly normalized, and compared it to the asymptotic variance  $\text{tr } \Sigma_B$  for the unnormalized model. The figure shows the distribution of the ratio  $\text{tr } \Sigma_B / \text{tr } \tilde{\Sigma}$  for different values of  $\nu$ . Interestingly, the ratio is almost equal to one for all tested values of  $\nu$ . Hence, additional estimation of the normalizing constant does not really seem to have had a negative effect on the accuracy of the estimates for the mixing matrix.

In Corollary 7, we have considered the hypothetical case where the noise distribution  $p_n$  is the same as the data distribution  $p_d$ . In Figure 4(d), we plot for that situation the asymptotic variance as a function of  $\nu$  (green curve). For reference, we plot again the curve for Gaussian contrastive noise (red circles, same as in Figure 4(b)). In both cases, we only show the asymptotic variance  $\text{tr } \Sigma_B$  for the parameters that correspond to the mixing matrix. The asymptotic variance for  $p_n = p_d$  is, for a given value of  $\nu$ , always smaller than the asymptotic variance for the case where the noise is Gaussian. However, by choosing  $\nu$  large enough for the case of Gaussian noise, it is possible to get estimates which are as accurate as those obtained in the hypothetical situation where  $p_n = p_d$ . Moreover, for larger  $\nu$ , the performance is the same for both cases: both converge to the performance of MLE.

---

4. See Appendix B.1 for the calculations in the special case of orthogonal mixing matrices.

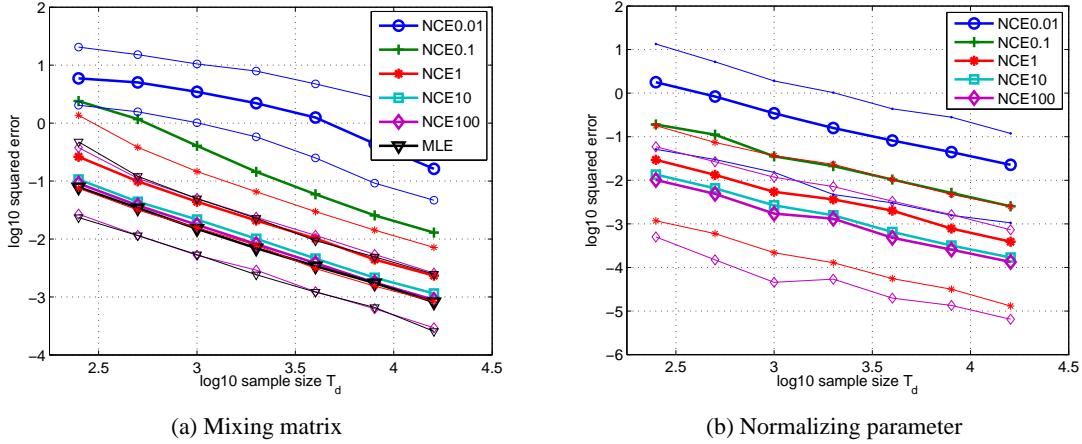


Figure 3: Validation of the theory of NCE: Estimation errors for an ICA model with four sources. Figures (a) and (b) show the mse for the mixing matrix  $\mathbf{B}$  and the normalizing parameter  $c$ , respectively. The performance of NCE approaches the performance of MLE(black triangles) as the ratio  $\nu = T_n/T_d$  increases: the case of  $\nu = 0.01$  is shown with blue circles,  $\nu = 0.1$  with green crosses,  $\nu = 1$  with red asterisks,  $\nu = 10$  with light blue squares, and  $\nu = 100$  with magenta diamonds. The thicker curves are the median of the performance for 500 random precision matrices with condition number smaller than ten. The finer curves show the 0.9 and 0.1 quantiles of the logarithm of the squared estimation error. To increase readability of the plots, the quantiles for  $\nu = 0.1$  and  $\nu = 10$  are not shown.

### 3.3 Scaling Properties

We use the ICA model from the previous subsection to study the behavior of NCE as the dimension  $n$  of the data increases. As before, we estimate the parameters by maximizing  $J_T(\boldsymbol{\theta})$  in Eq (8) with the nonlinear conjugate gradient algorithm of Rasmussen(2006). Again, we use a Gaussian with the same covariance structure as the data as noise distribution  $p_n$ .

The randomly chosen  $n \times n$  mixing matrices  $\mathbf{A}$  are restricted to be orthogonal. Orthogonality is only used to set up the estimation problem; in the estimation, the orthogonality property is not used. A reason for this restriction is that drawing mixing matrices at random as in the previous subsection leads more and more often to badly conditioned matrices as the dimension increases. Another reason is that the estimation error for orthogonal mixing matrices depends only on the dimension  $n$  and not on the particular mixing matrix chosen, see Appendix B.1 for a proof. Hence, this restriction allows us to isolate the effect of  $\dim n$  on the estimation accuracy.

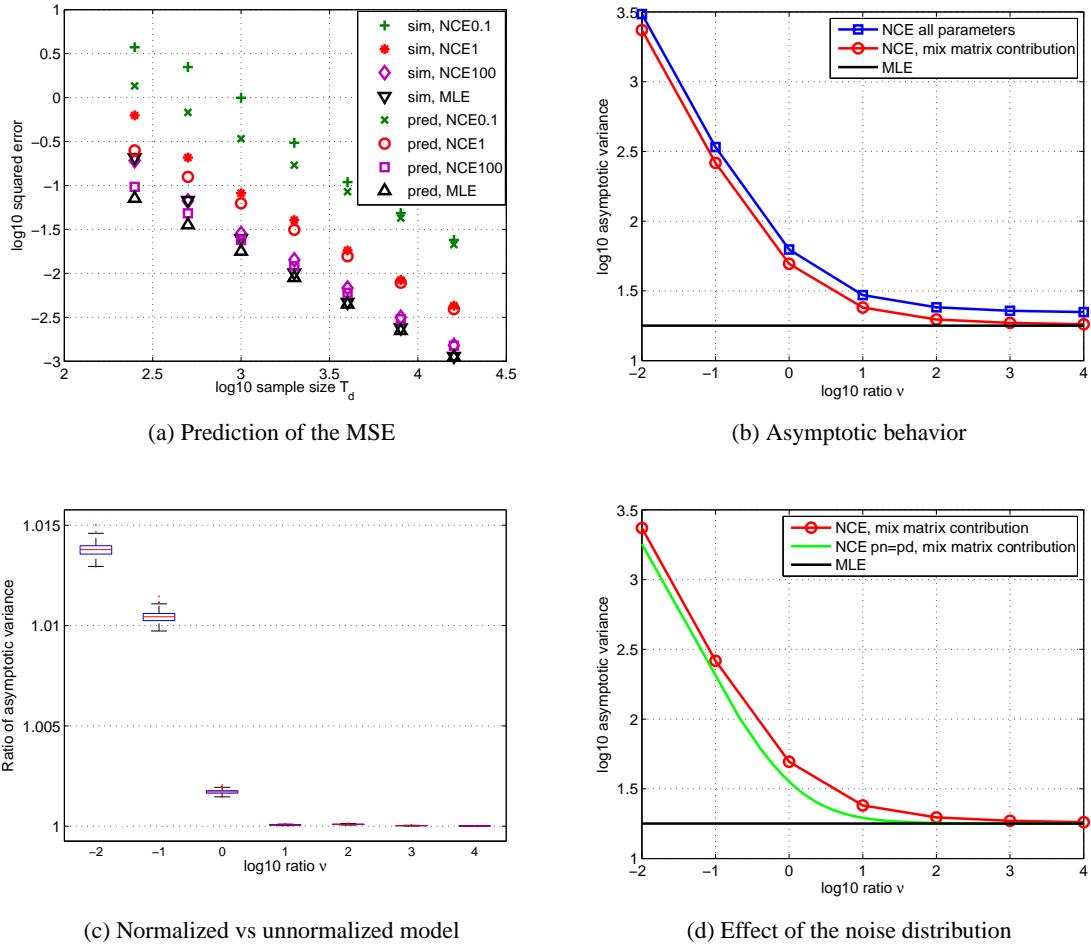


Figure 4: Validation of the theory of NCE: Estimation error for large sample sizes. Figure (a) shows that Corollary 4 correctly predicts the MSE for large samples sizes  $T_d$ . Figure (b) shows the asymptotic variance  $\text{tr } \Sigma$  as a function of  $\nu$ . Figure (c) shows a boxplot of the ratio between the asymptotic variance when the model is unnormalized and the asymptotic variance when the model is normalized. Figure (d) compares NCE with Gaussian noise to the hypothetical case where  $p_n$  equals the data distribution  $p_d$ . As in Figure 3, the curves in all figures but in Figure (c) are the median of the results for 500 random mixing matrices. The boxplot in Figure (c) shows the distribution for all the 500 matrices.

### 3.3.1 RESULTS

Figure 5(a) shows the asymptotic variance  $\text{tr } \Sigma_B$  related to the mixing matrix as a function of the dim  $n$ . NCE with  $\nu = T_n / T_d = 1$  is shown in red with asterisks as markers, MLE in black using triangles as markers. The markers show the theoretical prediction based on Corollary 4; the boxplots the simulation results for ten

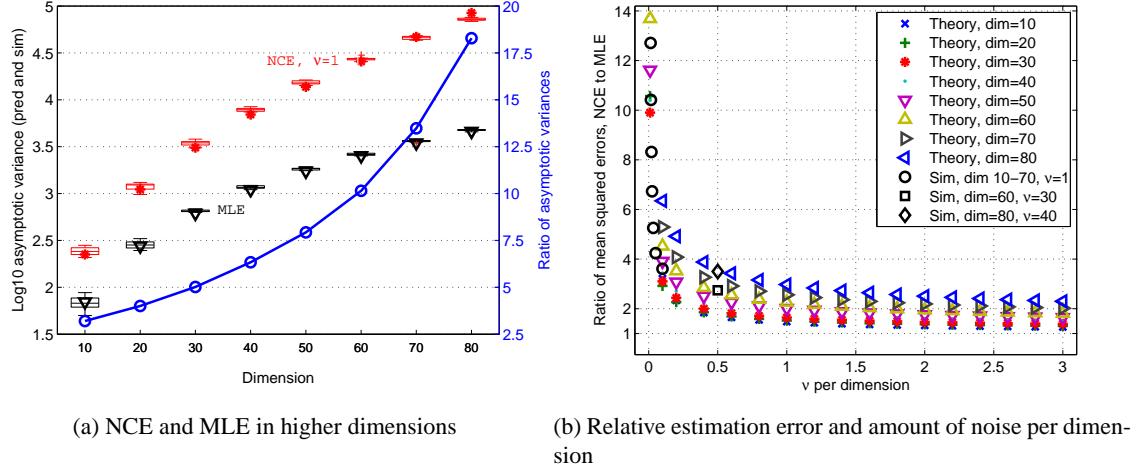


Figure 5: Investigating how NCE scales with the dimension of the data. Figure (a) shows the logarithm of the asymptotic variance for NCE ( $\nu = T_n/T_d = 1$ , in red) and MLE (in black). The boxplots show simulation results; the asterisks and triangles theoretical predictions for NCE and MLE, respectively. The same figure shows the ratio of the two asymptotic variances (blue circles, right scale). Figure (b) plots the ratio of the mean squared errors of the two estimators as a function of  $\nu$  per dimension  $n$ . The value of  $\nu$  needs to be increased as the dimensions increases; a linear increase leads to acceptable results.

random mixing matrices with  $T_d = 80000$ . The simulation results match the predictions well, which validates the theory of NCE in large dimensions.

Since the number of parameters increases with larger  $n$ , it is natural that  $\text{tr } \Sigma_B$  increases with  $n$ . However, for NCE, the increase is larger than for MLE. This is more clearly visible by considering the blue curve in Figure 5(a) (circles as markers, scale on the right axis). The curve shows the ratio between the asymptotic variance for noise-contrastive estimation and for MLE. By definition of the asymptotic variance, this ratio is equal to the ratio of the two estimation errors obtained with the two different methods. The ratio does not depend on the number of parameters and the sample size  $T_d$ . It is hence a suitable performance indicator to investigate how NCE scales with the dimension  $n$  of the data. The plot shows that for fixed  $\nu$ , the performance deteriorates as the dimension increases. In order to counteract this decline in performance, the parameter  $\nu$  needs to be increased as the dimension increases.

Figure 5(b) shows the ratio of the squared errors as a function of  $\nu/n$  where we varied  $n$  from ten to eighty dimensions as in Figure 5(a). Importantly, both theoretical results, where we numerically calculated the asymptotic variances, and simulation results show that for a reasonable performance in comparison to MLE,  $\nu$  does not need to be increased exponentially as the dimension  $n$  increases; a linear increase with, for instance,  $\nu \in [n/2, n]$  suffices to lead to estimation errors of about 2-4 times of those that are obtained by estimating normalized models with MLE.

## 4. Investigating the Trade-Off between Statistical and Computational Performance

We have seen that for large ratios  $\nu$  of noise sample size  $T_n$  to data sample size  $T_d$ , the estimation error for NCE behaves like the error in MLE. For large  $\nu$ , however, the computational load becomes also heavier because more noise samples need to be processed. There is thus a trade-off between statistical and computational performance. Such a trade-off exists also in other estimation methods for unnormalized models. In this section, we investigate the trade-off in NCE, and compare it to the trade-off in MCMLE (Geyer, 1994), contrastive divergence (Hinton, 2002) and persistent contrastive divergence<sup>5</sup> (Younes, 1989; Tielemans, 2008), as well as score matching (Hyvärinen et al., 2005).

In Section 4.1, we comment on the data which we use in the comparison. In Section 4.2, we review the different estimation methods with focus on the trade-off between statistical and computational performance. In Section 4.3, we point out the limitations of our comparison before presenting the simulation results in Section 4.4.

### 4.1 Data Used in the Comparison

For the comparison, we use artificial data which follows the ICA model in Equation (14) with the data log-pdf  $\ln p_d$  being given by Equation (15). We set the dimension  $n$  to ten and use  $T_d = 8000$  observations to estimate the parameters. In a first comparison, we assume Laplacian sources in the ICA model. The log-pdf  $\ln p_d$  is then specified with Equation (16). Note that this log-pdf has a sharp peak around zero where it is not continuously differentiable. In a second comparison, we use sources that follow the smoother logistic density. The nonlinearity  $f$  and the log normalizing constant  $c^*$  in Equation (15) are in that case

$$f(u) = -2 \ln \cosh \left( \frac{\pi}{2\sqrt{3}} u \right), \quad c^* = \ln |\det \mathbf{B}^*| + n \ln \left( \frac{\pi}{4\sqrt{3}} \right),$$

respectively. We are thus making the comparison for a relatively nonsmooth and smooth density. Both comparisons are based on 100 randomly chosen mixing matrices with condition number smaller than 10.

### 4.2 Estimation Methods Used in the Comparison

We introduce here briefly the different methods and comment on our implementation and choices of parameters.

#### 4.2.1 NCE

To estimate the parameters, we maximize  $J_T$  in Eq(8). We use here a Gaussian noise density  $p_n$  with a covariance matrix equal to the sample covariance of the data. As before,  $J_T$  is maximized using the nonlinear conjugate gradient method of Rasmussen (2006). To map out the trade-off between statistical and computational performance, we measured the estimation error and the time needed to optimize  $J_T$  for  $\nu \in \{1, 2, 5, 10, 20, 50, 100, 200, 400, 1000\}$ .

---

5. Persistent contrastive divergence= stochastic MLE.

#### 4.2.2 MC MLE

For normalized models, an estimate for the parameters  $\alpha$  can be obtained by choosing them such that the probability of the observed data is maximized. This is done by maximization of

$$J_{\text{MLE}}(\alpha) = \frac{1}{T_d} \sum_{t=1}^{T_d} \ln p_m^0(\mathbf{x}_t; \alpha) - \ln Z(\alpha). \quad (18)$$

If no analytical expression for the partition function  $Z(\alpha)$  is available, IS can be used to numerically approximate  $Z(\alpha)$  via its definition in Eq (2),

$$Z(\alpha) \approx \frac{1}{T_n} \sum_{t=1}^{T_n} \frac{p_m^0(\mathbf{n}_t; \alpha)}{p_{\text{IS}}(\mathbf{n}_t)}.$$

The  $\mathbf{n}_t$  are independent observations of “noise”  $\sim p_{\text{IS}}$ . Note that more sophisticated ways exist to numerically calculate the value of  $Z$  at a given  $\alpha$  (see Robert and Casella, 2004, in particular Chapter 3 and 4). The simple approach above leads to  $J_{\text{IS}}(\alpha)$  known as MCML (Geyer, 1994),

$$J_{\text{IS}}(\alpha) = \frac{1}{T_d} \sum_{t=1}^{T_d} \ln p_m^0(\mathbf{x}_t; \alpha) - \ln \left( \frac{1}{T_n} \sum_{t=1}^{T_n} \frac{p_m^0(\mathbf{n}_t; \alpha)}{p_{\text{IS}}(\mathbf{n}_t)} \right).$$

We maximized  $J_{\text{IS}}(\alpha)$  with the nonlinear conjugate gradient algorithm of Rasmussen (2006).

Like in NCE, there is a trade-off between statistical performance and running time: The larger  $T_n$  gets the better the approximation of the log-likelihood. Hence, the estimates become more accurate but the optimization of  $J_{\text{IS}}$  takes also more time. To map out the trade-off curve, we used the same values of  $T_n = \nu T_d$  as in NCE, and also the same noise distribution,  $p_{\text{IS}} = p_n$ .

#### 4.2.3 CONTRASTIVE DIVERGENCE

If  $J_{\text{MLE}}$  is maximized with a steepest ascent algorithm, the update rule for  $\alpha$  is

$$\alpha_{k+1} = \alpha_k + \mu_k \nabla_\alpha J_{\text{MLE}}(\alpha)|_{\alpha_k}, \quad (19)$$

where  $\mu_k$  is the step-size. For the calculation of  $\nabla_\alpha J_{\text{MLE}}$ , the gradient of the log partition function  $\ln Z(\alpha)$  is needed, see Eq (18). Above, IS was used to evaluate  $\ln Z(\alpha)$  and its gradient  $\nabla_\alpha \ln Z(\alpha)$ . The gradient of the log partition function:

$$\nabla_\alpha \ln Z(\alpha) = \frac{\nabla_\alpha Z(\alpha)}{Z(\alpha)} = \int \frac{p_m^0(\mathbf{n}; \alpha)}{Z(\alpha)} \nabla_\alpha \ln p_m^0(\mathbf{n}; \alpha) d\mathbf{n}. \quad (20)$$

If we had data  $\mathbf{n}_t$  at hand which follows the normalized model density  $p_m^0(\cdot; \alpha)/Z(\alpha)$ , the last equation could be evaluated by taking the sample average. The parameter vector  $\alpha$  could then be learned based on Eq (19). In general, sampling from the model density is, however, only possible by means of McMC methods. In CD (Hinton, 2002), to compute  $\alpha_{k+1}$ , Markov chains are started at the data points  $\mathbf{x}_t$  and stopped after a few Monte Carlo steps before they actually reach the stationary distribution  $p_m^0(\cdot; \alpha_k)/Z(\alpha_k)$ . The data points

$\mathbf{n}_t$  that are created in that way follow thus only approximately  $p_m^0(\cdot; \boldsymbol{\alpha}_k)/Z(\boldsymbol{\alpha}_k)$ . For every update of  $\boldsymbol{\alpha}$  the Markov chains are restarted from the  $\mathbf{x}_t$ . Note that this update rule for  $\boldsymbol{\alpha}$  is not directly optimizing a known objective function.

In our implementation, we used HMC (see for example Neal, 2010) with a rejection ratio of 10% for the sampling (like in Teh et al., 2004 ; Ranzato and Hinton , 2010 ). There are then four tuning parameters for CD: The number of Monte Carlo steps, the number of “leapfrog” steps in HMC, the choice of the step sizes  $\mu_k$ , as well as the number of data points  $\mathbf{x}_t$  and noise points  $\mathbf{n}_t$  used in each update step of  $\boldsymbol{\alpha}$ . The choice of the tuning parameters will affect the estimation error and the computation time . For our comparison here, we used contrastive divergence with one and three Monte Carlo steps (denoted by CD 1 and CD 3 in the figures below ), together with either three or twenty leapfrog steps . Ranzato and Hinton (2010 ) used CD 1 with twenty leapfrog steps (below denoted by CD 1 20), while Teh et al. (2004) used CD 1 30 to estimate unnormalized models from natural image data . For the  $\mu_k$ , we considered constant step sizes , as well as linearly and exponentially decaying step sizes.<sup>6</sup> For each update step, we chose an equal number of data and noise points. We considered the case of using all data in each update step, and the case of using minibatches of only 100 randomly chosen data points.

We selected the step size  $\mu_k$  and the number of data points used in each update by means of preliminary simulations on five data sets. We limited ourselves to CD with one Monte Carlo and three leapfrog steps (CD 1 3). For both Laplacian and logistic sources , using mini - batches with an exponential decaying step size gave the best results . The results are reported below in Section 4.4. The use of minibatches led to faster estimation results without affecting their accu - racy . Exponentially decaying step sizes are advocated by the theory of stochastic approximation ; in some cases, however, linear decay was found to be more appropriate (Tieleman , 2008, Section 4.5). For Laplacian sources, the initial step size  $\mu_0$  was 0.005; for logistic sources, it was  $\mu_0=0.01$ . Note that in this selection of the tuning parameters, we used the true parameters to compute the estimation error. Clearly , this cannot be done in real applications since the true parameter values are not known. The choice of the tuning parameters must then solely be based on experience, as well as trial and error.

#### 4.2.4 PERSISTENT CONTRASTIVE DIVERGENCE

As CD, PCD(Younes , 1989 ; Tieleman , 2008) uses the update rule in Eq (19) together with an approximative evaluation of the integral in Eq(20) to learn the parameters  $\boldsymbol{\alpha}$ . The integral is also computed based on McMC sampling . Unlike CD, the Markov chains are not restarted at the data points  $\mathbf{x}_t$ . For the computation of  $\boldsymbol{\alpha}_{k+1}$ , the Markov chains are initialized with the samples  $\mathbf{n}_t$  that were obtained in the previous iteration by running Markov chains  $\rightarrow p_m^0(\cdot; \boldsymbol{\alpha}_{k-1})/Z(\boldsymbol{\alpha}_{k-1})$ . As in CD, the Markov chains are only run for a short time and stopped before having actually converged.

Since PCD differs from CD only by the initialization of the Markov chains, it has the same tuning parameters . As in CD, we used preliminary simulations to select suitable parameters : again , exponentially decaying step sizes  $\mu_k$  together with minibatches of size 100 gave the best performance . The preliminary simulations yielded also the same initial step sizes  $\mu_0$  as in CD. It turned out, however,

---

6. Linear decay:  $\mu_k = \mu_0(1 - k/maxIteration)$ , exponential decay:  $\mu_k = \mu_0 C/(C + k)$  with  $C = 5000$ .

that the number of leapfrog steps in PCD needs to be larger than in CD: using, for example, only three leapfrog steps as in CD resulted in a poor performance in terms of estimation accuracy. For the results reported below in Section 4.4, we used 20 and 40 leapfrog steps, together with one and three Monte Carlo steps.

#### 4.2.5 SCORE MATCHING

In score matching (Hyvärinen, 2005), the parameter vector  $\alpha$  is estimated by minimization of the cost function  $J_{\text{SM}}$ ,

$$J_{\text{SM}}(\alpha) = \frac{1}{T_d} \sum_{t=1}^{T_d} \sum_{i=1}^n \frac{1}{2} \Psi_i^2(\mathbf{x}_t; \alpha) + \Psi'_i(\mathbf{x}_t; \alpha).$$

The term  $\Psi_i(\mathbf{x}; \alpha)$  is the derivative of the unnormalized model wrt  $\mathbf{x}(i)$ , the  $i$ -th element of  $\mathbf{x}$ ,

$$\Psi_i(\mathbf{x}; \alpha) = \frac{\partial \ln p_m^0(\mathbf{x}; \alpha)}{\partial \mathbf{x}(i)}.$$

The term  $\Psi'_i(\mathbf{x}; \alpha)$  denotes the derivative of  $\Psi_i(\mathbf{x}; \alpha)$  with respect to  $\mathbf{x}(i)$ . The presence of this derivative may make the objective function and its gradient algebraically rather complicated if a sophisticated model is estimated. For the ICA model with Laplacian sources,  $\Psi_i(\mathbf{x}; \alpha)$  equals

$$\Psi_i(\mathbf{x}; \alpha) = \sum_{j=1}^n -\sqrt{2} \text{sign}(\mathbf{b}_j \mathbf{x}) B_{ji} \quad (21)$$

which is not smooth enough to be used in score matching. Using the smooth approximation  $\text{sign}(u) \approx \tanh(10u)$  is a way to obtain a smooth enough  $\Psi_i(\mathbf{x}; \alpha)$  and  $\Psi'_i(\mathbf{x}; \alpha)$ . The optimization of  $J_{\text{SM}}$  is done by the nonlinear conjugate gradient algorithm of Rasmussen (2006). Note that, unlike the estimation methods considered above, score matching does not have a tuning parameter which controls the trade-off between statistical and computational performance. Moreover, score matching does not rely on sampling.

### 4.3 Limitations of the Comparison

For all considered methods but CD/PCD, the algorithm which is used to optimize the given objectives can be rather freely chosen. This choice will influence the trade-off between statistical and computational performance. Here, we use the optimization algorithm by Rasmussen (2006). Our results below show thus the trade-off of the different estimation methods in combination with this particular optimization algorithm. With this optimization algorithm, we used for each update all data. The algorithm is not suitable for stochastic optimization with minibatches (see for example Schraudolph and Graepel, 2002). Optimization based on mini-batches may well lead not only for (persistent) contrastive divergence to gains in speed but also for the other estimation methods, including nce.

It is well known that a Gaussian as noise (proposal) distribution is not the optimal choice for IS if the data has heavy tails (Wasserman, 2004, Chapter 24). Gaussian noise is not the optimal choice for NCE either. The presented results should thus not be considered as a general comparison of the two estimation methods perse. Importantly, however, the chosen setup allows one to assess how NCE behaves when the data has heavier tails than the noise, which is often the case in practice.

Finally, the reader may want to keep in mind that for other kinds of data, in particular also in very high dimensions, differences may occur.

## 4.4 Results

We first compare NCE with the methods for which we use the same optimization algorithm, that is MCMLE and score matching. Then, we compare it with CD/PCD.

### 4.4.1 COMPARISON WITH MC MLE AND SCORE MATCHING

Figure 6 shows the comparison of (NCE, red squares), MCML (IS, blue circles) and score matching (SM, black triangles). The left panels show the simulation results in form of “result points” where the x -coordinate represents the time till the algorithm converged and the y -coordinate the estimation error at convergence. Convergence in the employed nonlinear conjugate gradient algorithm by Rasmussen (2006) means that the line search procedure failed twice in a row to meet the strong Wolfe-Powell conditions (Sun and Yuan , 2006 , Chapter 2.5.2). For score matching , 100 result points corresponding to 100 different random mixing matrices are shown in each figure. For NCE and MCml , we used ten different values of  $\nu$  so that for these methods , each figure shows 1000 result points . The panels on the right present the simulation result in a more schematic way . For NCE and MCml , the different ellipses represent the outcomes for different values of  $\nu$ . Each ellipse contains 90% of the result points . We can see that increasing  $\nu$  reduces the estimation error but it also increases the running time. For score matching, there is no such trade-off.

Figure 6(a) shows that for Laplacian sources, NCE outperforms the other methods in terms of the trade-off between statistical and computational performance . The large estimation error of score matching is likely to be due to the smooth approximation of the sign function in Eq(21). The figure also shows that NCE handles noise that has lighter tails than the data more gracefully than MCMLE. The reason is that the nonlinearity  $h(\mathbf{u}; \boldsymbol{\theta})$  in the objective function in Eq (8) is bounded even if data and noise distribution do not match well (see also Pihlaja et al., 2010).

For logistic sources , shown in Figure 6(b), nce and MCml perform equally . Score matching reaches its level of accuracy about 20 times faster than the other methods . NCE and MCml can, however , have a higher estimation accuracy than score matching if  $\nu$  is large enough . Score matching can thus be considered to have a built -in trade - off between estimation performance and computation time : Computations are fast but the speed comes at the cost of not being able to reach an estimation accuracy as high as, for instance, NCE.

### 4.4.2 COMPARISON WITH CONTRASTIVE AND PERSISTENT CONTRASTIVE DIVERGENCE

Since CD/PCD do not have an objective function and given the randomness that is introduced by the minibatches, it is difficult to choose a reliable stopping criterion. Hence, we did not impose any stopping criterion but the maximal number of iterations. The two algorithms had always converged before this maximal number of iterations was reached in the sense that the estimation error did not visibly decrease any more.

We base our comparison on the estimation error as a function of the running time of the algorithm. This makes the comparison independent from the stopping criterion that is used in NCE. For NCE, the parameter  $\nu$  controls the trade-off between computational and statistical performance ; for CD/PCD, it is the number of leapfrog steps and the number of Markov steps taken in each update . We compiled a trade -off curve for each of the one hundred estimation problems by taking at any time point the minimum estimation error over the various estimation errors that are obtained for different values of the trade -off parameters.<sup>7</sup> Figure 7 shows an example for NCE and CD. The distribution of the trade-off curves is shown in Figure 8. For large running times, the distribution of the estimation error is for all estimation methods similar to the one for MLE. For shorter running times, NCE is seen to have for Laplacian sources a better trade-off than the other methods. For logistic sources, however, the situation is reversed.

#### 4.4.3 SUMMARY

The foregoing simulation results and discussion suggest that all estimation methods trade, in one form or the other, estimation accuracy against computation speed. In terms of this trade-off, NCE is particularly well suited for the estimation of data distributions with heavy tails . In case of thin tails , NCE performs similarly to MCml , and CD/PCD has a better trade-off. If the data distribution is particularly smooth and the model algebraically not too complicated , score matching may , depending on the required estimation accuracy, be the best option.

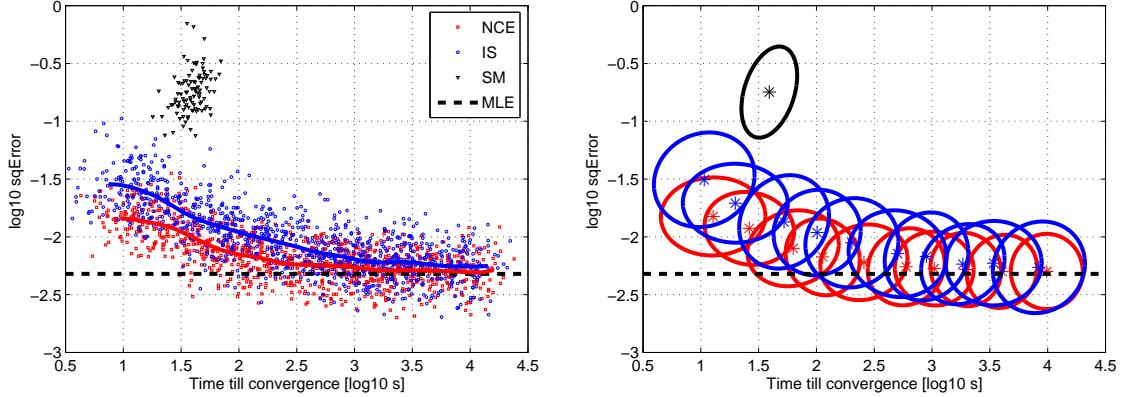
### 5. Simulations with Natural Images

In this section, we estimate with our new estimation method models of natural images. In the theory of NCE, we have assumed that all variables can be observed . NCE can thus not be used for models with latent variables which cannot be integrated out analytically . Such models occur for example in the work by Olshausen and Field (1996), Hyvarinen et al. (2001a), Karklin and Lewicki (2005), Lucke and Sahani (2008) and Osindero and Hinton (2008). We are here considering models which avoid latent variables . Recent models which are related to the models that we are considering here can be found in the work by Osindero et al. (2006), Koster and Hyvarinen (2010) and Ranzato and Hinton (2010). For a comprehensive introduction to natural image statistics, see for example the textbook by Hyvarinen et al. (2009).

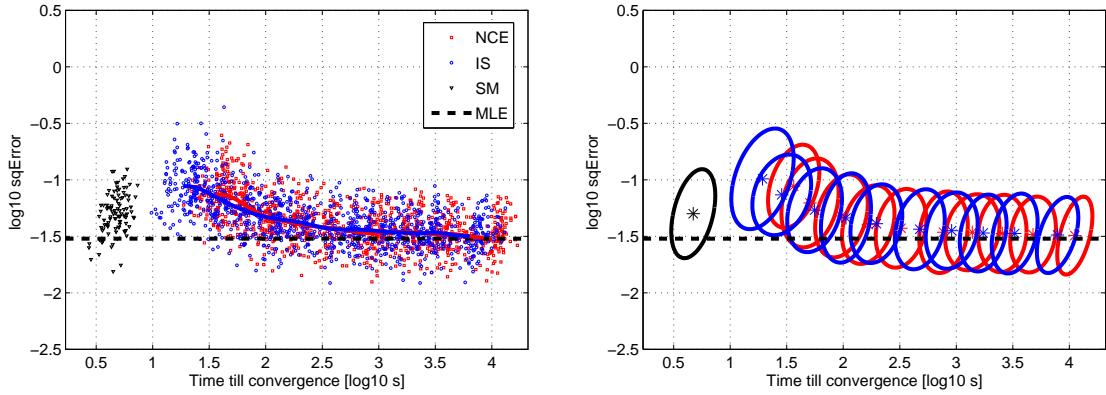
The presented models will consist of two processing layers, like in a multilayer neural network. The output of the network for a given input image gives the value of the model-pdf at that image. Because of the two processing layers, we call the models “two-layer models”.

We start with giving some preliminaries in Section 5.1. In Section 5.2, we present the settings of nce. In Section 5.3, we properly define the two-layer model and estimate a version with more than 50000 parameters . In Section 5.4, we present an extension of the model where the learned output nonlinearity of the network belongs to the flexible family of splines . The different models are compared in Section 5.5.

7. A comparison of CD and PCD for different settings can be found in Appendix C.1.



(a) Sources following a Laplacian density



(b) Sources following a logistic density

Figure 6: Trade-off between statistical and computational performance for (NCE, red squares), MCML (IS, blue circles) and score matching (SM, black triangles). Each point represents the result of one simulation. Performing local linear kernel smoothing regression on the result points yields the thick curves. For NCE and MCML, the ten ellipses represent the outcomes for the ten different values of  $\nu \in \{1, 2, 5, 10, 20, 50, 100, 200, 400, 1000\}$ . The ellipses were obtained by fitting a Gaussian to the distribution of the result points, each one contains 90% of the results points for a given  $\nu$ . The asterisks mark their center. For an ICA model with Laplacian sources, NCE has the best trade-off between statistical and computational performance. For logistic sources, NCE and IS perform equally well. For medium estimation accuracy, score matching outperforms the other two estimation methods.

## 5.1 Data, Preprocessing and Modeling Goal

Our basic data are a random sample of  $25\text{px} \times 25\text{px}$  image patches that we extracted from a subset of van Hateren's image database (van Hateren and van der Schaaf, 1998). The images in the subset showed wildlife scenes only. The sample size  $T_d$  is 160000.

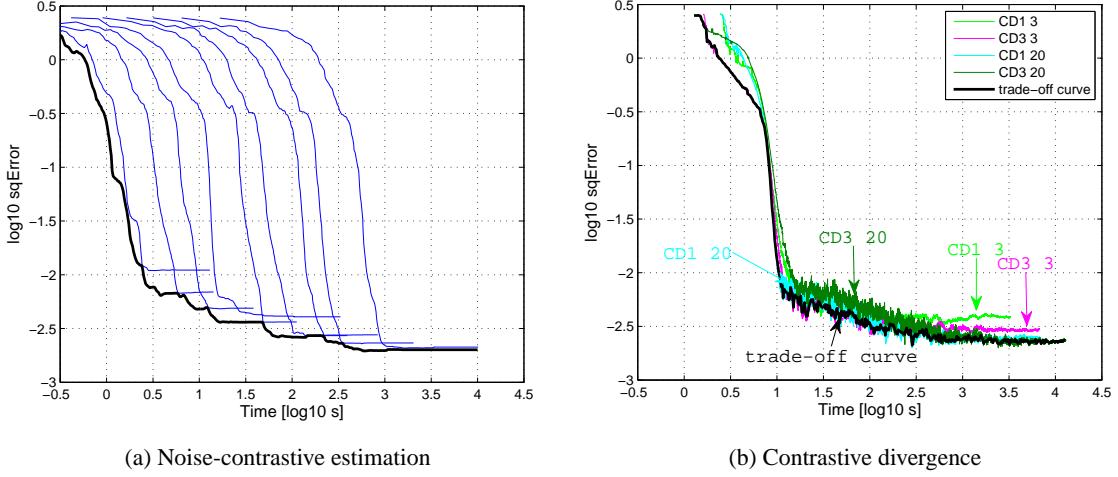


Figure 7: Example of a trade-off curve for NCE and CD.

- (a) The different curves in blue show the estimation error which is obtained for the various values of  $\nu$ . The thicker curve in black shows the trade-off curve. It is obtained by taking at any time point the minimum estimation error. (b) The trade-off curve, shown in black, is similarly obtained by taking the minimum over the estimation errors which are obtained with different settings of CD.

As preprocessing, we removed from each image patch its average value (local mean, DC component), whitened the data and reduced the dimension from  $d = 25 \cdot 25 = 625$  to  $n = 160$ . This retains 93% of the variance of the image patches. After dimension reduction, we additionally centered each data point and rescaled it to unit variance. In order to avoid division by small numbers, we avoided taking small variance patches. This gave our data  $X = (\mathbf{x}_1, \dots, \mathbf{x}_{T_d})$ . Because of the centering and rescaling, each data point  $\mathbf{x}_t$  satisfies

$$\sum_{k=1}^n \mathbf{x}_t(k) = 0, \quad \frac{1}{n-1} \sum_{k=1}^n \mathbf{x}_t(k)^2 = 1. \quad (22)$$

This means that each data point lies on the surface of a  $n - 1$  dimensional sphere  $\mathbb{S}$ .

This kind of preprocessing is a form of luminance and contrast gain control which aim at canceling out the effects of the lighting conditions (see for example Hyvärinen et al., 2009, Chapter 9, where also the statistical effects of such a preprocessing are analyzed). Centering and rescaling to unit variance has also been used in image quality assessment in order to access the structural component of an image, which is related to the reflectance of the depicted objects (Wang et al., 2004, in particular Section III.B). By modeling the data  $X$ , we are thus modeling the structure in the image patches.

Given a data point  $\mathbf{x}_t$ , we can reconstruct the original (vectorized) image patch via

$$\mathbf{i}_t = \mathbf{V}^- \mathbf{x}_t, \quad \mathbf{V}^- = \mathbf{E} \mathbf{D}^{1/2}, \quad (23)$$

where  $\mathbf{E}$  is the  $d \times n$  matrix formed by the leading  $n$  eigenvectors of the covariance matrix of the image patches. The diagonal  $n \times n$  matrix  $\mathbf{D}$  contains the corresponding eigenvalues. The matrix

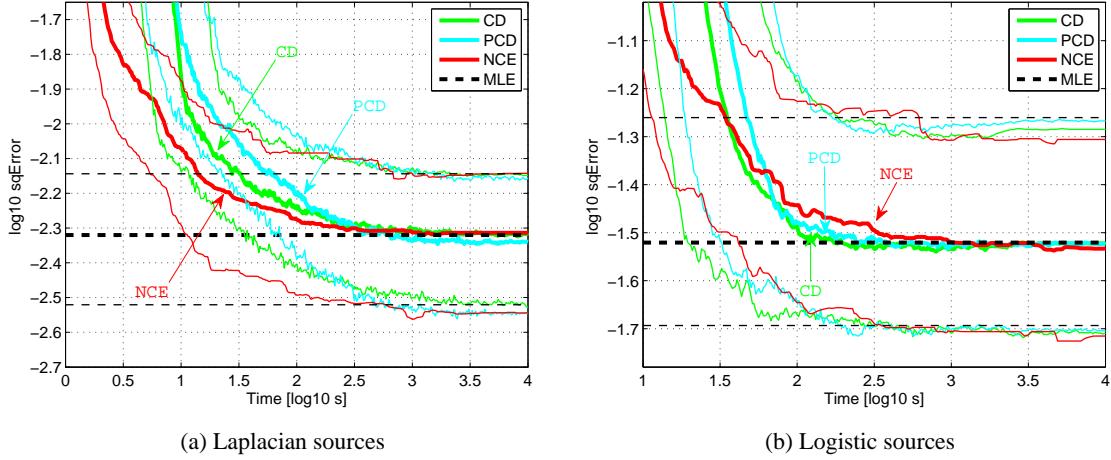


Figure 8: Distribution of the trade-off curves for CD(green), PCD(cyan), and NCE(red). The distribution of the estimation error for MLE is shown in black. The thick curves show the median, the finer curves the 0.9 and 0.1 quantiles.

$\mathbf{V}^-$  defined above is the pseudoinverse of the whitening matrix  $\mathbf{V} = \mathbf{D}^{-1/2}\mathbf{E}^T$ . Since the column vectors of  $\mathbf{V}^-$  form a basis for a  $n$  dimensional subspace of  $\mathbb{R}^d$ ,  $\mathbf{x}$  is the coordinate vector of  $\mathbf{i}$  with respect to that basis. The dimension reduction implies that the reconstruction cannot be perfect; the reconstruction can also only be performed up to the scale and average value of the patch because of the luminance and contrast gain control. Figure 9(a) shows examples of natural image patches after extraction from the data base; Figure 9(b) shows the corresponding reconstructions  $\mathbf{i}$ . Since all image patches in Figure 9 were rescaled to use the full colormap, the effects of luminance and contrast gain control are not visible. The effect of the dimension reduction is low-pass filtering.

## 5.2 Settings for NCE

Matlab code for the simulations is available from the authors' homepage so that our description here will not be exhaustive. All the models considered in the next subsections are estimated with noise-contrastive estimation. We learn the parameters by optimization of the objective  $J_T$  in Equation (8). The two-layer models are estimated by first estimating one-layer models. The learned parameters are used as initial values for the first layer in the estimation of the complete two-layer model. The second layer is initialized to small random values.

For the cnd  $p_n$ , we take a uniform distribution on the surface of the  $n - 1$  dimensional sphere  $\mathbb{S}$  on which  $\mathbf{x}$  is defined.<sup>8</sup> Examples of image patches with coordinates following  $p_n$  are shown in Figure 9(c). Samples from  $p_n$  can easily be created by sampling from  $N(0, 1)$ , followed by centering and rescaling such that Eq (22) holds. Since  $p_n$  is a constant, the log-ratio  $G(\cdot; \theta)$  in Eq (4) is up to an additive constant equal to

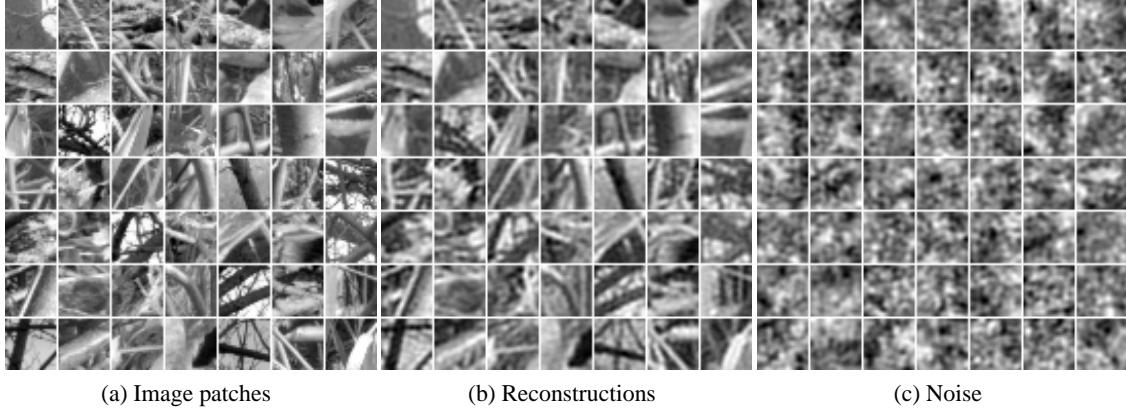


Figure 9: (a) Natural image patches of size  $25\text{px} \times 25\text{px}$ . (b) Reconstructed image patches after pre-processing. These are examples of the image patches denoted by  $\mathbf{i}$  in Equation (23) with coordinate vectors  $\mathbf{x} \in \mathbb{R}^{160}$ . (c) Noise images which are obtained via Equation (23) if the coordinates are uniformly distributed on the sphere  $\mathbb{S}$ . Comparison with Figure (b) shows that the coordinate vectors  $\mathbf{x}$  for natural images are clearly not uniformly distributed on the sphere. In the next subsections, we model their distribution.

$$\ln p_m(\cdot; \boldsymbol{\theta}),$$

$$G(\cdot; \boldsymbol{\theta}) = \ln p_m(\cdot; \boldsymbol{\theta}) + \text{constant}.$$

As pointed out in Section 2.2,  $\boldsymbol{\theta}$  evolves in the maximization of  $J_T$  such that  $G(\mathbf{u}; \hat{\boldsymbol{\theta}}_T)$  is as large as possible for  $\mathbf{u} \in X$  (natural images) but as small as possible for  $\mathbf{u} \in Y$  (noise). For uniform noise, the same must thus also hold for  $\ln p_m(\mathbf{u}; \hat{\boldsymbol{\theta}}_T)$ . This observation will be a useful guiding tool for the interpretation of the models below.

The factor  $\nu = T_n/T_d$  was set to 10. We found that an iterative optimization procedure where we separate the data into subsets and optimize  $J_T$  for increasingly larger values of  $\nu$  reduced computation time. The optimization for each  $\nu$  is done with the nonlinear conjugate gradient method of Rasmussen (2006). The size of the subsets is rather large, for example 80000 in the simulation of the next subsection.<sup>9</sup> A more detailed discussion of this optimization procedure can be found in Appendix C.2.

### 5.3 Two-Layer Model with Thresholding Nonlinearities

The first model that we consider is

$$\ln p_m(\mathbf{x}; \boldsymbol{\theta}) = \sum_{k=1}^n f(y_k; a_k, b_k) + c, \quad y_k = \sum_{i=1}^n Q_{ki} (\mathbf{w}_i^T \mathbf{x})^2, \quad (24)$$

where  $f$  is a smooth, compressive thresholding function that is parameterized by  $a_k$  and  $b_k$ . See Figure 10 for details regarding the parameterization and the formula for  $f$ . The parameters  $\boldsymbol{\theta}$  of

---

9. As pointed out in Section 4.3, the used nonlinear conjugate gradient algorithm is not suitable for stochastic optimization with small minibatches.

the model are the second-layer weights  $Q_{ki} \geq 0$ , the first-layer weights  $\mathbf{w}_i \in \mathbb{R}^n$ , the normalizing parameter  $c \in \mathbb{R}$ , as well as  $a_k > 0$  and  $b_k \in \mathbb{R}$  for the nonlinearity  $f$ . The definition of  $y_k$  shows that multiplying  $Q_{ki}$  by a factor  $\gamma_i^2$  and  $\mathbf{w}_i$  at the same time by the factor  $1/\gamma_i$  does not change the value of  $y_k$ . There is thus some ambiguity in the parameterization which could be resolved by imposing a norm constraint either on the  $\mathbf{w}_i$  or on the columns of the matrix  $\mathbf{Q}$  formed by the weights  $Q_{ki}$ . It turned out that for the estimation of the model such constraints were not necessary. For the visualization and interpretation of the results, we chose  $\gamma_i$  such that all the  $\mathbf{w}_i$  had norm one.

The motivation for the thresholding property of  $f$  is that, in line with Section 5.2,  $\ln p_m(\cdot; \boldsymbol{\theta})$  can easily be made large for natural images and small for noise. The  $y_k$  must just be above the thresholds for natural image input and below for noise. This occurs when the vectors  $\mathbf{w}_i$  detect features (regularities) in the input which are specific to natural images, and when, in turn, the second-layer weights  $Q_{ki}$  detect characteristic regularities in the squared first-layer feature outputs  $\mathbf{w}_i^T \mathbf{x}$ . The squaring implements the assumption that the regularities in  $\mathbf{x}$  and  $-\mathbf{x}$  are the same so that the pdf of  $\mathbf{x}$  should be an even function of the  $\mathbf{w}_i^T \mathbf{x}$ . Another property of the nonlinearity is its compressive log-like behavior for inputs above the threshold. The motivation for this is to “counteract” the squaring in the computation of  $y_k$ . The compression of large values of  $y_k$  leads to numerical robustness in the computation of  $\ln p_m$ .

A model like the one in Eq (24) has been studied before by Osindero et al. (2006) and Koster and Hyvärinen (2010). There are, however, a number of differences. The main difference is that in our case  $\mathbf{x}$  lies on a sphere while in the cited work,  $\mathbf{x}$  was defined in the whole space  $\mathbb{R}^n$ . This difference allows us to use nonlinearities that do not decay asymptotically to  $-\infty$  which is necessary if  $\mathbf{x}$  is defined in  $\mathbb{R}^n$ . A smaller difference is that we do not need to impose norm constraints to facilitate the learning of the parameters.

### 5.3.1 RESULTS

For the visualization of the first-layer feature detectors  $\mathbf{w}_i$ , note that the inner product  $\mathbf{w}_i^T \mathbf{x}$  equals  $(\mathbf{w}_i^T \mathbf{V}) \mathbf{i} = \mathbf{w}_i^T \mathbf{i}$ . The  $\mathbf{w}_i \in \mathbb{R}^n$  are coordinate vectors wrt the basis given by the columns of  $\mathbf{V}^-$ , see Section 5.1, while the  $\mathbf{w}_i \in \mathbb{R}^d$  are the coordinate vectors with respect to the pixel basis. The latter vectors can thus be visualized as images. This is done in Figure 11(a). Another way to visualize the first-layer feature detectors  $\mathbf{w}_i$  is to show the images which yield the largest feature output while satisfying the constraints in Eq(22). These optimal stimuli are proportional to  $\mathbf{V}^- (\mathbf{w}_i - \langle \mathbf{w}_i \rangle)$ , where  $\langle \mathbf{w}_i \rangle \in \mathbb{R}$  is the average value of the elements in the vector  $\mathbf{w}_i$ , see Appendix B.2 for a proof. The optimal stimuli are shown in Figure 11(b). Both visualizations show that the first layer computes “Gabor-like” features, which is in line with previous research on natural image statistics.

Figure 12 shows a random selection of the learned second-layer weights  $Q_{ik}$ . Figure 12(a) shows that the weights are extremely sparse. The optimization started with the weights being randomly assigned to small values, with the optimization most of them shrunk to zero; few selected ones, however, increased in magnitude. Note that this result was obtained without any norm constraints on  $\mathbf{Q}$ . From Figure 12(b), we see that the learned second-layer weights  $Q_{ik}$  are such that they combine first-layer features of similar orientation, which are centered at nearby locations (“complex cells”). The same figure shows also a condensed representation of the feature detectors using icons. This form of visualization is used in Figure 13 to visualize all the second-layer feature detectors.

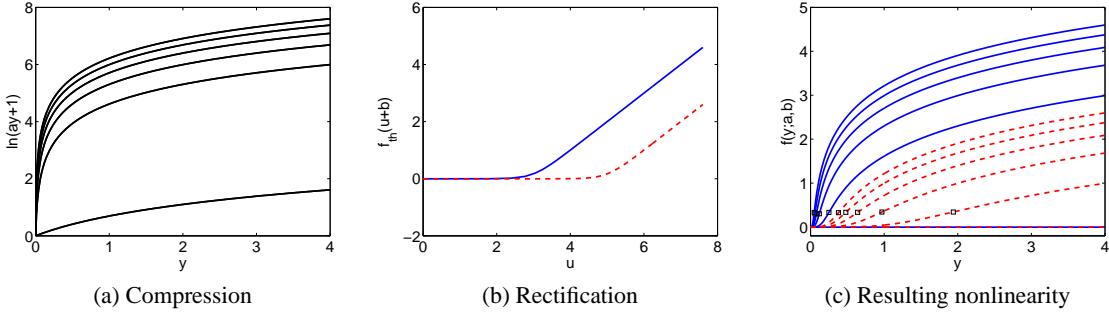


Figure 10: Two-layer model with thresholding nonlinearities. The family of nonlinearities used in the modeling is  $f(y; a, b) = f_{\text{th}}(\ln(ay + 1) + b)$ ,  $y \geq 0$ . The parameterized function is composed of a compressive nonlinearity  $\ln(ay + 1)$ , shown in Figure (a), and a smooth rectification function  $f_{\text{th}}(u + b)$  shown in Figure (b). Figure (c) shows examples of  $f(y; a, b)$  for different values of  $a$  and  $b$ . Parameter  $b$  sets the threshold, and parameter  $a$  controls the steepness of the function. Since the scale of the weights in Equation (24) is not restrained, the parameters  $a_k$  do not need to be learned explicitly. After learning, they can be identified by dividing  $y_k$  in Equation (24) by  $a_k$  so that its expectation is one for natural images. The formula for the thresholding function is  $f_{\text{th}}(u) = 0.25 \ln(\cosh(2u)) + 0.5u + 0.17$ . The curves shown in blue are for  $b = -3$  and  $a \in \{1, 50, 100, 200, \dots, 500\}$ . For the dashed curves in red,  $b = -5$ . The small squares in Figure (c) indicate where  $f$  changes from convex to concave.

Figure 14(a) shows the learned nonlinearities  $f(\cdot; a_k, b_k)$ . Note that we incorporated the learned normalizing parameter  $c$  as an offset  $c/n$  for each nonlinearity. The learned thresholding is similar for feature outputs of mid- and high-frequency feature detectors (black, solid curves). For the feature detectors tuned to low frequencies, the thresholds tend to be smaller (green, dashed curves). The nonlinearities in black are convex for arguments  $y$  smaller than two (see red rectangle in the figure). That is, they show a squashing behavior for  $y < 2$ . Looking at the distribution of the second-layer outputs  $y_k$  in Figure 14(b), we see that it is more likely that noise rather than natural images was the input when the second-layer feature outputs  $y_k$  are approximately between 0.5 and 2. In this regime, the squashing nonlinearities map thus more often the noise input to small values than natural images so that  $\ln p_m(\mathbf{u}; \hat{\theta}_T)$  tends to be larger when input  $\mathbf{u}$  is a natural image than when it is noise (see Section 5.2). One could, however, think that the thresholding nonlinearities are suboptimal because they ignore the fact that natural images lead, compared to the noise, rather often to  $y_k$  which are close to zero, see Figure 14(b). An optimal nonlinearity should, unlike the thresholding nonlinearities, assign a large value to both large and small  $y_k$  while mapping intermediate values of  $y_k$  to small numbers. The next subsection shows that such kinds of mappings emerge naturally when splines are used to learn the nonlinearities from the data.

#### 5.4 Two-Layer Model with Spline Nonlinearities

In the previous subsection, the family of nonlinearities  $f$  in Eq (24) was rather limited. Here, we look for  $f$  in the larger family of cubic splines where we consider the location of the knots to be fixed

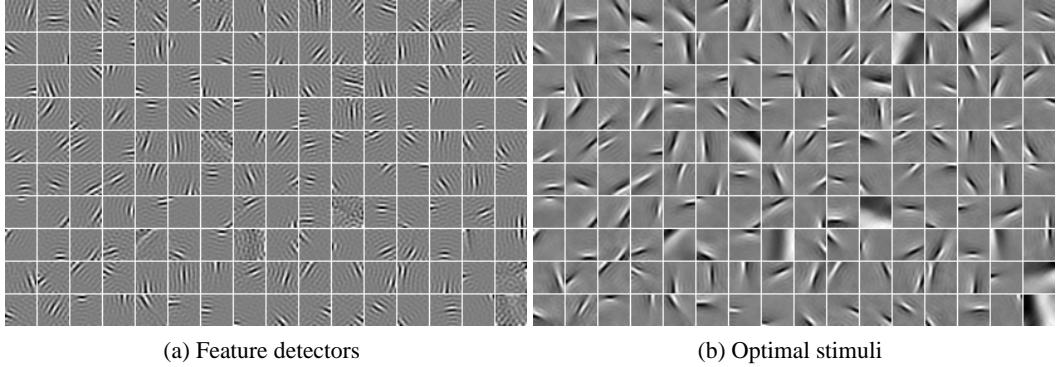


Figure 11: Two-layer model with thresholding nonlinearities: Visualization of the learned first-layer feature detectors  $\mathbf{w}_i$ . (a) The feature detectors in the pixel basis. (b) The corresponding optimal stimuli. The feature detectors in the first layer are “Gabor-like” (localized, oriented, bandpass). Comparison of the two figures shows that feature detectors which appear noisy in the pixel basis are tuned to low-frequency input.

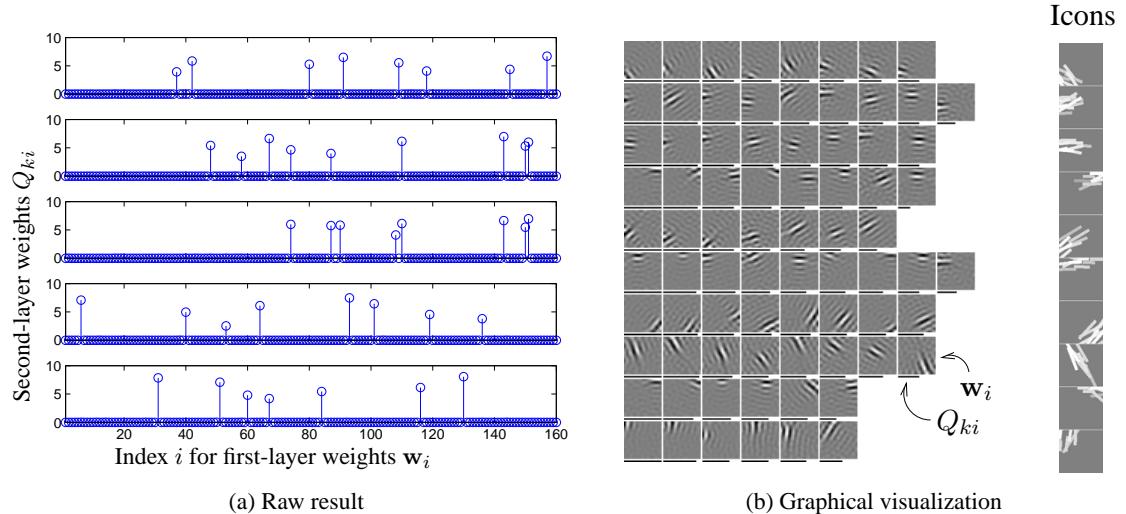


Figure 12: Two-layer model with thresholding nonlinearities: Random selection of second layer units. (a) Second-layer weights  $Q_{ki}$  for five different  $k$  (five different rows of the matrix  $\mathbf{Q}$ ) are shown. The weights are extremely sparse so that in the sum  $\sum_{i=1}^n Q_{ki}(\mathbf{w}_i^T \mathbf{x})^2$  only few selected squared first-layer outputs are added together. (b) Every row shows one second-layer feature detector. The first-layer feature detectors  $\mathbf{w}_i$  are shown as image patches like in Figure 11, and the black bar under each patch indicates the strength  $Q_{ki}$  by which a certain  $\mathbf{w}_i$  is pooled by the  $k$ -th second-layer feature detector. The numerical values  $Q_{ki}$  for the first five rows are shown in Figure (a). The right-most column shows a condensed visualization. The icons were created by representing each first-layer feature by a bar of the same orientation and similar length as the feature, and then superimposing them with weights given by  $Q_{ki}$ .

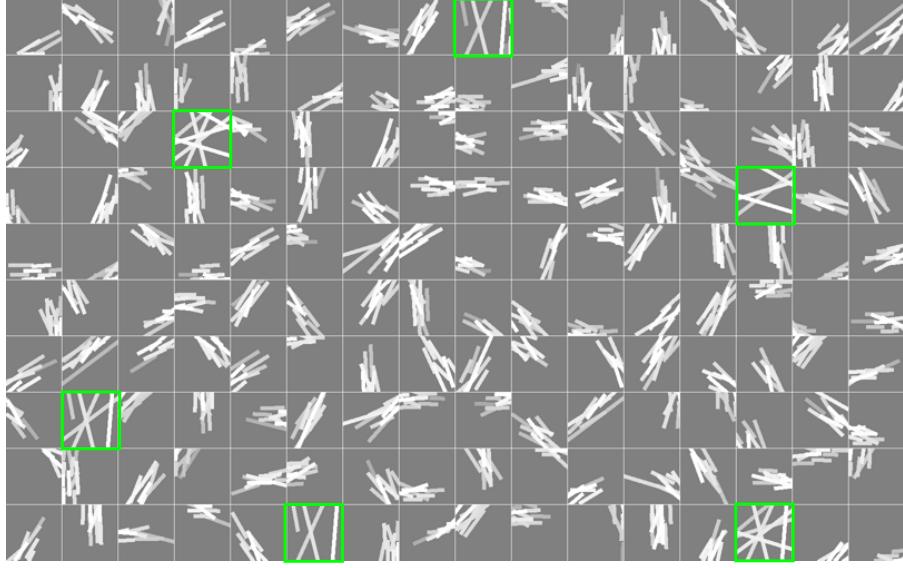


Figure 13: Two-layer model with thresholding nonlinearities: Visualization of the first- and second-layer feature detectors with icons. In the second layer, first-layer features of similar orientations are pooled together. See Figure 12 for details of how the icons were created. The feature detectors marked with a green frame are tuned to low frequencies.

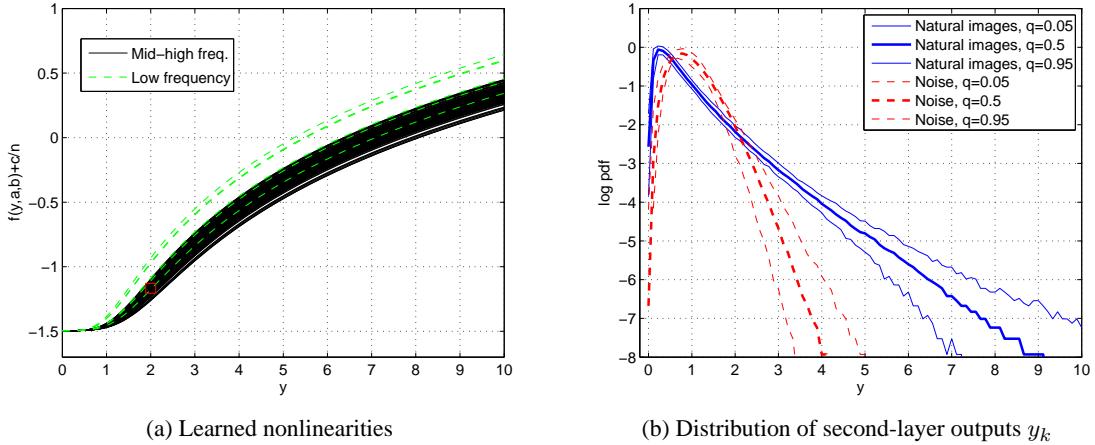


Figure 14: Two-layer model with thresholding nonlinearities: Learned nonlinearities and interpretation. Natural images tend to have larger second-layer outputs  $y_k$  than noise input since the two processing layers, visualized in Figures 11 to 13, detect structure inherent to natural images. Thresholding the  $y_k$  provides a way to assign to natural images large values in the model-pdf and to noise small values. In Figure (a), the nonlinearities acting on pooled low-frequency feature detectors are shown in green (dashed lines), those for medium and high frequency feature detectors in black (solid lines). The bold curves in Figure (b) show the median, the other curves the 5% and 95% quantiles. The solid curves in blue relate to natural images, the dashed curves in red to noise. As explained in Figure 10, the  $y_k$  have expectation one for natural images.

(regression splines represented with B-spline basis functions, Hastie et al., 2009, Chapter 5).

The model that we consider here is

$$\ln p_m(\mathbf{x}; \boldsymbol{\theta}) = \sum_{k=1}^n f(y_k; a_1, a_2, \dots) + c, \quad y_k = \sum_{i=1}^n Q_{ki} (\mathbf{w}_i^T \mathbf{x})^2. \quad (25)$$

The difference between this and the model of the previous subsection is that the output nonlinearity  $f$  is a cubic spline. Part of the parameters  $\boldsymbol{\theta}$  are thus as previously the  $\mathbf{w}_i \in \mathbb{R}^n$ ,  $Q_{ki} \geq 0$ , and  $c \in \mathbb{R}$ . Additional parameters are the  $a_i \in \mathbb{R}$  which are the coefficients of the B-spline basis functions of the cubic spline  $f$ . As before, we denote the matrix formed by the  $Q_{ki}$  by  $\mathbf{Q}$ .

For the modeling of the nonlinearity  $f$ , we must define its domain, which is the range of its arguments  $y_k$ . A way to control the range of  $y_k$  is to constrain the norm of the columns of  $\mathbf{Q}$  and also to constrain the vectors  $\mathbf{w}_k$  such that

$$\max_i \mathbb{E} \left\{ (\mathbf{w}_i^T \mathbf{x})^2 \right\} = 1, \quad (26)$$

where the expectation is taken over the natural images.

We estimated the model in Equation (25) by first estimating a spline-based one-layer model which is presented in Appendix C.3. In brief, in this model, we did not square the first-layer feature outputs  $\mathbf{w}_i^T \mathbf{x}$  and the matrix  $\mathbf{Q}$  was the identity. The arguments of the spline nonlinearity  $f$  were thus the feature outputs  $\mathbf{w}_i^T \mathbf{x}$  without additional processing. The learned nonlinearity is shown in Figure 16(a). In the following, we denote it by  $f_1$ . In Appendix C.3, we point out that the shape of  $f_1$  is closely related to the sparsity of the feature outputs when natural images are the input. Because  $f_1$  is an even function, and because of the squaring in the definition of  $y_k$ , we initialized  $f$  for the estimation of the two-layer model as  $f(u) = f_1(\sqrt{u})$ . This function is shown in Figure 16(b) (blue, dashes). The learned  $\mathbf{w}_i$  of the one-layer model were used as initial points for the estimation of the two-layer model. The  $Q_{ki}$  were randomly initialized to small values. It turned out that imposing Equation (26) was enough for the learning to work and no norm constraint for the columns of  $\mathbf{Q}$  was necessary. The results were very similar whether there were norm constraints or not. In the following, we report the results without any norm constraints.

#### 5.4.1 RESULTS

Figure 15 visualizes the learned parameters  $\mathbf{w}_i$  and  $Q_{ki}$  in the same way as in Figures 12 and 13 for the two-layer model with thresholding nonlinearities. The learned feature extraction stage is qualitatively very similar, up to two differences. The first difference is that many second-layer weights  $Q_{ki}$  shrank to zero: 66 out of 160 rows of the matrix  $\mathbf{Q}$  had so small values that we could omit them while accounting for 99.9% of the sum  $\sum_{ki} Q_{ki}$ . The second difference is that the pooling in the second layer is sometimes less sparse. In that case, the second layer still combines first-layer feature detectors of the same orientation but they are not all centered at the same location.

The learned nonlinearity  $f$  is shown in Figure 16(b) (black, solid). The nonlinearity from the one-layer model, shown in blue as a dashed curve, is altered so that small and large inputs are assigned to larger numbers while intermediate inputs are mapped to smaller numbers. Compared to the thresholding nonlinearities from the previous subsection, the learned nonlinearity has also for small inputs large outputs. Since the second-layer feature outputs  $y_k$  are sparser (that is, more often very small or large) for natural images than for the noise, the shape of the learned nonlinearity

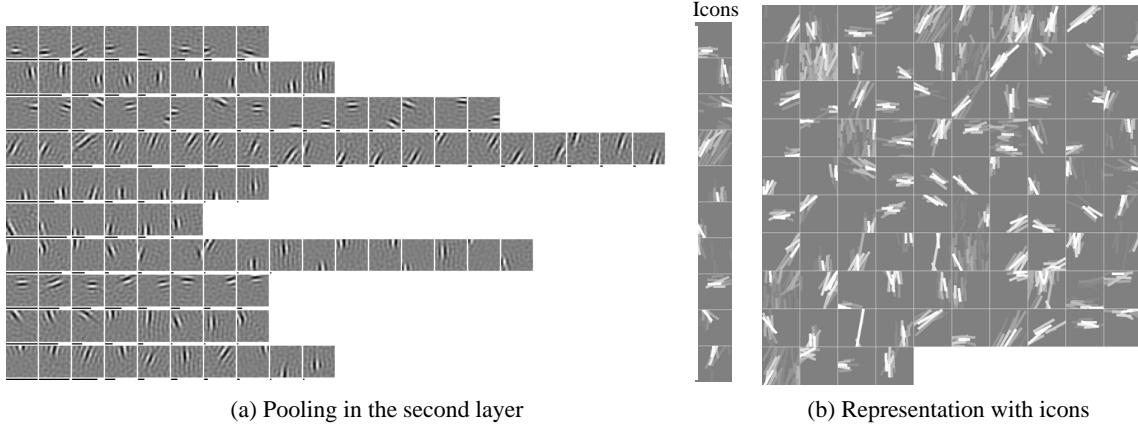


Figure 15: Two-layer model with spline nonlinearities. (a) Random selection of the learned second-layer units. (b) Representation of all the learned second-layer feature detectors as iconic images.

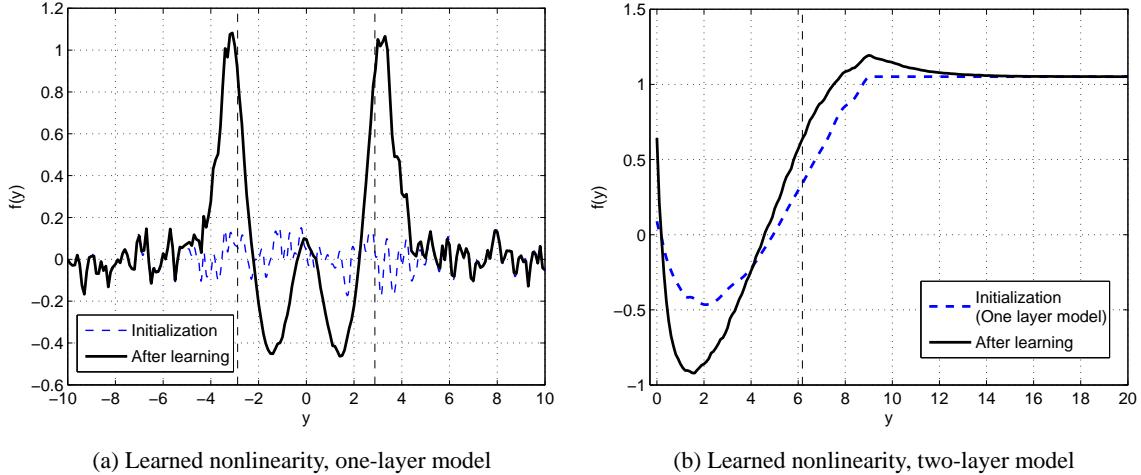


Figure 16: Two-layer model with spline nonlinearities. (a) Learned nonlinearity (black, solid) and its random initialization (blue, dashes) for the one-layer model. The learned nonlinearity is used as starting point in the learning of the two-layer model. (b) Learned nonlinearity (black, solid) and its initialization (blue, dashes) for the two-layer model. The dashed vertical lines indicate the 99% quantile for all the feature outputs for natural images. Due to the lack of training examples, the nonlinearities should not be considered valid beyond these lines.

implies that the estimated model assigns more often a higher probability density to natural images than to the noise.

## 5.5 Model Comparison

We have estimated models for natural images, both with thresholding nonlinearities and with splines. We make here a simple model comparison.

A quantitative comparison is done by calculating for ten validation sets the value of the objective function  $J_T$  of noise-contrastive estimation (see Equation (8) for the definition). The sample size of each validation set was  $T_v = 100000$ , and  $\nu$  was set to 10, as in the estimation of the models. For the same validation data, we also computed the performance measure  $\hat{L} = 1/T_v \sum_t \ln p_m(\mathbf{x}_t; \hat{\theta}_T)$ , which is an estimate for the rescaled log-likelihood, see Equation (13) in Section 3.1. As pointed out there,  $\hat{L}$  is only an estimate of the rescaled log-likelihood because  $\hat{c}$ , which is an element of the parameter vector  $\hat{\theta}_T$ , is used instead of the correct normalizing constant. Both  $J_T$  and the log-likelihood have the property that models which fit the data better have a higher score.

Comparing the structure of data points which are considered likely by the different models is a way to make a qualitative model comparison. Another approach would be to sample from the models, which we do in Appendix C.5. In order to get the likely points, we drew random samples that followed the noise distribution  $p_n$  (uniform on the sphere), and used them as initial points in the optimization of the various log-densities  $\ln p_m(\mathbf{x}; \hat{\theta}_T)$  with respect to  $\mathbf{x}$  under the constraint of Equation (22). We used the same initial points for all models and visualized the likely points  $\hat{\mathbf{x}}$  via Equation (23) as images  $\hat{\mathbf{i}} = \mathbf{V}^{-1} \hat{\mathbf{x}}$ .

The ICA model with Laplacian sources is a simple model for natural images. It has previously also been used to model natural images after they have been projected on a sphere (Hyvärinen et al., 2009, Chapter 9). The unnormalized model has been defined in Section 3.2 in Equation (17) and consists of one processing layer with the fixed nonlinearity  $f(u) = -\sqrt{2}|u|$ . We include it in our comparison and refer to it as one-layer model with ‘‘Laplacian nonlinearity’’.

### 5.5.1 RESULTS

Table 1 shows that the spline-based two-layer model of Section 5.4 gives, on average, the largest value of the objective function  $J_T$ , and also  $L_T$ . To investigate the merits of the spline output-nonlinearity, we fixed the feature extraction stage of the thresholding model in Section 5.3 and learned only the nonlinearity  $f$  using splines (for details, see Appendix C.4). The resulting model, labeled ‘‘refinement’’ in the table, performs nearly as good as the best model. The one-layer models with thresholding or Laplacian nonlinearities have the smallest objectives  $J_T$  and  $L_T$ . The two models achieve the objectives in different, complimentary ways. For the thresholding model, the absolute value of the feature outputs  $\mathbf{w}_i^T \mathbf{x}$  must be large to yield a large objective while for the model with the Laplacian nonlinearity  $f(\mathbf{w}_i^T \mathbf{x}) = -\sqrt{2}|\mathbf{w}_i^T \mathbf{x}|$ , the feature outputs must have small absolute values. The two models consider thus different aspects of the, for natural images, typically sparse feature outputs  $\mathbf{w}_i^T \mathbf{x}$ . The one-layer model with spline nonlinearity combines both aspects, see Figure 16(a), and yields also a higher score in the comparison. The same reason explains why spline-based two-layer models have higher scores than the two-layer model with the thresholding nonlinearity.

Figure 17 shows the likely data points from the various models  $p_m$ . The models with large objectives in Table 1 lead to image patches with particularly clear structure. The emergence of structure can be explained in terms of sparse coding since image patches which lead to sparse activations of the feature detectors are typically highly structured. Sparseness of the feature outputs

	One-layer model			Two-layer model		
	Thresholding	Laplacian	Spline	Thresholding	Refinement	Spline
$J_T$ , av	-1.871	-1.518	-1.062	-0.8739	-0.6248	<b>-0.6139</b>
$J_T$ , std	0.0022	0.0035	0.0030	0.0029	0.0030	0.0037
$L_T$ , av	-223.280	-222.714	-219,786	-220.739	-213.303	<b>-212.598</b>
$L_T$ , std	0.0029	0.0077	0.0137	0.0088	0.0282	0.0273

Table 1: Quantitative model comparison. The objective  $J_T$  of noise-contrastive estimation, see Equation (8), and the estimate  $\hat{L}$  of the (rescaled) log-likelihood, see Equation (13), are used to measure the performance. Larger values indicate better performance. The table gives the average (av) and the standard deviation (std) for ten validation sets. All models are defined on a sphere and learned with noise-contrastive estimation. The features for the one-layer models with thresholding and Laplacian nonlinearity are not shown in the paper. The “one-layer, thresholding” model is identical to the “two-layer, thresholding” model when the second layer is fixed to the identity matrix. With Laplacian nonlinearity we mean the function  $f(u) = -\sqrt{2}|u|$ . The “two-layer, thresholding” model has been presented in Section 5.3, and the “two-layer, spline” model in Section 5.4. The “one-layer, spline” and “two-layer, refinement” models are presented in the Appendix C.3 and C.4, respectively.

is facilitated by the nonlinearities in the models, and through the competition between the features by means of the sphere-constraint on the coordinates  $\mathbf{x}$ , as specified in Equation (22).

## 6. Conclusions

In this paper, we have considered the problem of estimating unnormalized statistical models for which the normalizing partition function cannot be computed in closed form. Such models cannot be estimated by maximization of the likelihood without resorting to numerical approximations which are often computationally expensive. The main contribution of the paper is a new estimation method for unnormalized models. A further contribution is made in the modeling of natural image statistics.

We have proven that our new estimation method, NCE, provides a consistent estimator for both normalized and unnormalized statistical models. The assumptions that must be fulfilled to have consistency are not stronger than the assumptions that are needed in MLE. We have further derived the asymptotic distribution of the estimation error which shows that, in the limit of arbitrarily many contrastive noise samples, the estimator performs like the maximum likelihood estimator. The new method has a very intuitive interpretation in terms of supervised learning : The estimation is performed by discriminating between the observed data and some artificially generated noise by means of logistic regression.

All theoretical results were illustrated and validated on artificial data where ground truth is known. We have also used artificial data to assess the balance between statistical and computational performance. In particular, we have compared the new estimation method to a number of other estimation methods for unnormalized models : Simulations suggest that NCE strikes a highly competitive trade-off. We have used the mse of the estimated param-

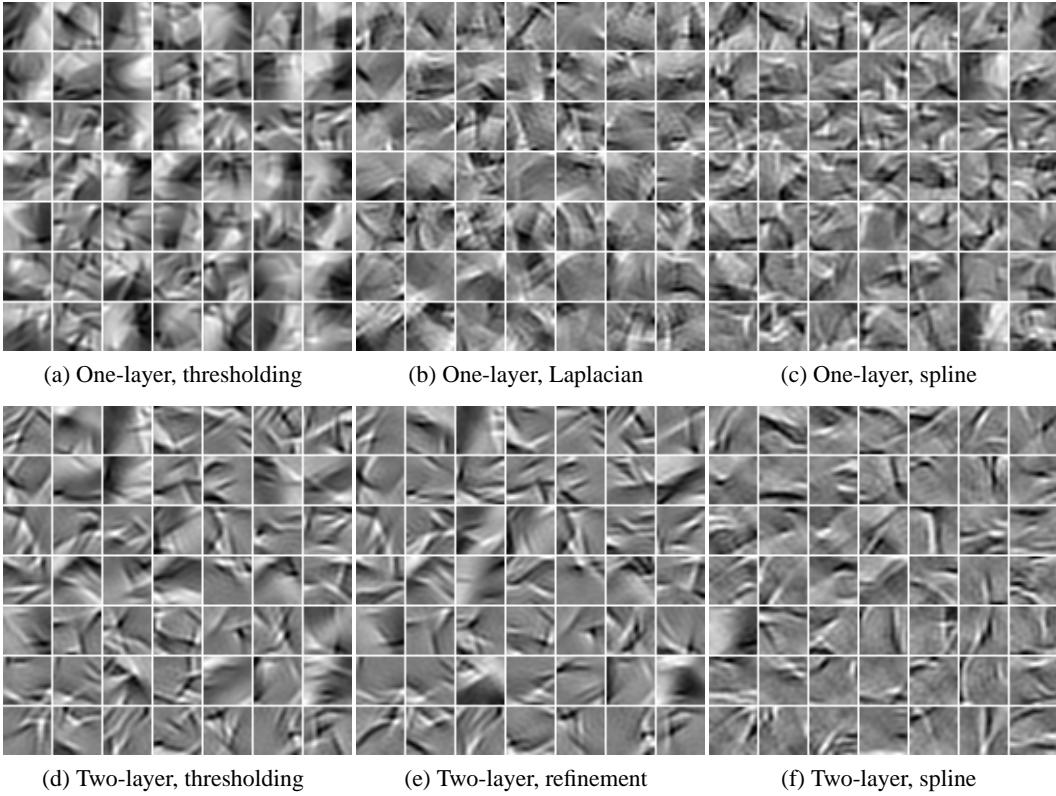


Figure 17: Likely points under the learned models for natural images. See caption of Table 1 for information on the models.

eters as statistical performance measure. It should be noted that this is only one possible criterion among many (see Hyvärinen, 2008, for a recently proposed alternative measure of performance).

NCE as presented here extends the previous definition given by Gutmann and Hyvarinen (2010) since it allows for more noise samples than data points. We have also previously considered such a generalization (Pihlaja et al., 2010). Unlike in that preliminary version, our method here is asymptotically Fisher-efficient for all admissible noise densities when the number of noise samples becomes arbitrarily large. Pihlaja et al. (2010) has established links of NCE to IS which remain valid for this paper.

We applied NCE to the modeling of natural images. Besides validating the method on a large two-layer model, we have, as a new contribution to the understanding of natural image statistics, presented spline-based extensions: In previous models, the output nonlinearity in the pdf was hand-picked. Here, we have parameterized it as a spline and learned it from the data. The statistical models were all unnormalized and had several ten-thousands of parameters which demonstrates that our new method can handle demanding estimation problems.

## Acknowledgments

Yoshua Bengio, Ian Goodfellow, Pascal Vincent, Geoffrey Hinton, Nicolas Le Roux, Marc'Aurelio Ranzato, and Ilya Sutskever. This work was funded by the Centre-of-Excellence in Algorithmic Data Analysis and the Computational Sciences program, both of the Academy of Finland.

## Appendix A. Proofs of the Theorems

### A.1 Preliminaries

$$r_\nu(u) := \frac{1}{1 + \nu \exp(-u)},$$

which was introduced in Equation (6):

$$\begin{aligned} 1 - r_\nu(u) &= r_{\frac{1}{\nu}}(-u) \\ \frac{\partial r_\nu(u)}{\partial u} &= r_{\frac{1}{\nu}}(-u)r_\nu(u) \\ \frac{\partial}{\partial u} \ln r_\nu(u) &= r_{\frac{1}{\nu}}(-u) \\ \frac{\partial^2}{\partial u^2} \ln r_\nu(u) &= -r_{\frac{1}{\nu}}(-u)r_\nu(u) \\ \frac{\partial}{\partial u} \ln[1 - r_\nu(u)] &= -r_\nu(u) \\ \frac{\partial^2}{\partial u^2} \ln[1 - r_\nu(u)] &= -r_{\frac{1}{\nu}}(-u)r_\nu(u) \end{aligned}$$

The functions  $h(\mathbf{u}; \boldsymbol{\theta}) = r_\nu(G(\mathbf{u}; \boldsymbol{\theta}))$  and  $1 - h(\mathbf{u}; \boldsymbol{\theta}) = r_{\frac{1}{\nu}}(-G(\mathbf{u}; \boldsymbol{\theta}))$  are equal to

$$h(\mathbf{u}; \boldsymbol{\theta}) = \frac{p_m(\mathbf{u}; \boldsymbol{\theta})}{p_m(\mathbf{u}; \boldsymbol{\theta}) + \nu p_n(\mathbf{u})}, \quad 1 - h(\mathbf{u}; \boldsymbol{\theta}) = \frac{\nu p_n(\mathbf{u})}{p_m(\mathbf{u}; \boldsymbol{\theta}) + \nu p_n(\mathbf{u})}, \quad (27)$$

see Eq (3). ==>

$$\nu p_n(\mathbf{u})r_\nu(G(\mathbf{u}; \boldsymbol{\theta})) = \frac{\nu p_n(\mathbf{u})p_m(\mathbf{u}; \boldsymbol{\theta})}{p_m(\mathbf{u}; \boldsymbol{\theta}) + \nu p_n(\mathbf{u})}, \quad (28)$$

$$p_d(\mathbf{u})r_{\frac{1}{\nu}}(-G(\mathbf{u}; \boldsymbol{\theta})) = \frac{\nu p_n(\mathbf{u})p_d(\mathbf{u})}{p_m(\mathbf{u}; \boldsymbol{\theta}) + \nu p_n(\mathbf{u})}, \quad (29)$$

Taylor expansions

$$\begin{aligned}\ln r_\nu(u + \epsilon u_1 + \epsilon^2 u_2) &= \ln r_\nu(u) + \epsilon r_{\frac{1}{\nu}}(-u)u_1 + \\ &\quad \epsilon^2 \left[ r_{\frac{1}{\nu}}(-u)u_2 - \frac{1}{2} r_{\frac{1}{\nu}}(-u)r_\nu(u)u_1^2 \right] + \\ &\quad O(\epsilon^3),\end{aligned}\tag{30}$$

$$\begin{aligned}\ln [1 - r_\nu(u + \epsilon u_1 + \epsilon^2 u_2)] &= \ln[1 - r_\nu(u)] - \epsilon r_\nu(u)u_1 + \\ &\quad \epsilon^2 \left[ -r_\nu(u)u_2 - \frac{1}{2} r_{\frac{1}{\nu}}(-u)r_\nu(u)u_1^2 \right] + \\ &\quad O(\epsilon^3).\end{aligned}\tag{31}$$

## A.2 Proof of Theorem 1

For clarity of the proof, we state an important stepping stone as a lemma.

### A.2.1 LEMMA

**Lemma 8** For  $\epsilon > 0$  and  $\phi(\mathbf{x})$  a perturbation of the log-pdf  $f_m(\mathbf{x}) = \ln p_m(\mathbf{x})$ ,

$$\begin{aligned}\tilde{J}(f_m + \epsilon\phi) &= \tilde{J}(f_m) + \epsilon \int [p_d(\mathbf{u})r_{\frac{1}{\nu}}(-f_m(\mathbf{u}) + \ln p_n(\mathbf{u})) - \\ &\quad \boxed{\nu p_n(\mathbf{u})r_\nu(f_m(\mathbf{u}) - \ln p_n(\mathbf{u}))}] \phi(\mathbf{u}) d\mathbf{u} - \\ &\quad \frac{\epsilon^2}{2} \int r_{\frac{1}{\nu}}(-f_m(\mathbf{u}) + \ln p_n(\mathbf{u}))r_\nu(f_m(\mathbf{u}) - \ln p_n(\mathbf{u})) \\ &\quad (p_d(\mathbf{u}) + \nu p_n(\mathbf{u})) \phi(\mathbf{u})^2 d\mathbf{u} + O(\epsilon^3).\end{aligned}$$

Gateaux derivation

**Proof** making use of Eq (30) and (31) with  $u = f_m(\mathbf{x}) - \ln p_n(\mathbf{x})$ ,  $u_1 = \phi(\mathbf{x})$  and  $u_2 = 0$ .

■

### A.2.2 PROOF OF THE THEOREM

**Proof** A necessary condition for optimality is that in the expansion of  $J^\sim(f_m + \epsilon\phi)$ ,  $G$  derivation = 0 iff

$$p_d(\mathbf{u})r_{\frac{1}{\nu}}(-f_m(\mathbf{u}) + \ln p_n(\mathbf{u})) = \nu p_n(\mathbf{u})r_\nu(f_m(\mathbf{u}) - \ln p_n(\mathbf{u})).$$

With Eq (28) and (29), this implies that  $J^\sim$  has an extremum at  $p_m$  iff

$$\frac{\nu p_n(\mathbf{u})p_d(\mathbf{u})}{p_m(\mathbf{u}) + \nu p_n(\mathbf{u})} = \frac{\nu p_n(\mathbf{u})p_m(\mathbf{u})}{p_m(\mathbf{u}) + \nu p_n(\mathbf{u})}.$$

$\Leftrightarrow \nu > 0, p_m(\mathbf{u}) = p_d(\mathbf{u})$  where  $p_n(\mathbf{u}) \neq 0$ . At points where  $p_n(\mathbf{u}) = 0$ , the eq is trivially fulfilled. Hence,  $p_m = p_d$ , or  $f_m = \ln p_d$ , leads to an extremum of  $J^\sim$ .

Inserting  $f_m = \ln p_d$  into  $\tilde{J}$  in Lemma 8 leads to

$$\tilde{J}(\ln p_d + \epsilon\phi) = \tilde{J}(\ln p_d) - \frac{\epsilon^2}{2} \left\{ \int \frac{\nu p_n(\mathbf{u}) p_d(\mathbf{u})}{p_d(\mathbf{u}) + \nu p_n(\mathbf{u})} \phi(u)^2 d\mathbf{u} \right\} + O(\epsilon^3).$$

Since the term of order  $\epsilon^2$  is negative for all choices of  $\phi$ , the extremum is a maximum. The assumption that  $p_n(\mathbf{u}) > p_d(\mathbf{u})$  shows that  $f_m = \ln p_d$  is the only extremum.  $\blacksquare$

### A.3 Proof of Theorem 2 (Consistency)

#### A.3.1 LEMMATA

The Taylor expansions in Eq (30) and (31) are used to prove the following lemma which is like Lemma 8 for  $J^\sim$  but for the objective function  $J$  in Eq (10).

**Lemma 9** For  $\epsilon > 0$  and  $\boldsymbol{\varphi} \in \mathbb{R}^m$ ,

$$\begin{aligned} J(\boldsymbol{\theta} + \epsilon\boldsymbol{\varphi}) &= J(\boldsymbol{\theta}) + \epsilon \int u_1 [p_d(\mathbf{u})(1 - h(\mathbf{u}; \boldsymbol{\theta})) - \nu p_n(\mathbf{u})h(\mathbf{u}; \boldsymbol{\theta})] d\mathbf{u} + \\ &\quad \epsilon^2 \left\{ \int -\frac{1}{2}u_1^2(1 - h(\mathbf{u}; \boldsymbol{\theta}))h(\mathbf{u}; \boldsymbol{\theta})(p_d(\mathbf{u}) + \nu p_n(\mathbf{u})) d\mathbf{u} + \right. \\ &\quad \left. \int u_2(p_d(\mathbf{u})(1 - h(\mathbf{u}; \boldsymbol{\theta})) - \nu p_n(\mathbf{u})h(\mathbf{u}; \boldsymbol{\theta})) d\mathbf{u} \right\} + O(\epsilon^3), \end{aligned}$$

where

$$\begin{aligned} u_1 &= \boldsymbol{\varphi}^T \mathbf{g}(\mathbf{u}; \boldsymbol{\theta}), \\ u_2 &= \frac{1}{2} \boldsymbol{\varphi}^T \mathbf{H}_G(\mathbf{u}; \boldsymbol{\theta}) \boldsymbol{\varphi}. \end{aligned}$$

$\mathbf{g}(\mathbf{u}; \boldsymbol{\theta}) = \nabla G(\mathbf{u}; \boldsymbol{\theta})$ , and  $\mathbf{H}_G$  = the Hessian matrix of  $G(\mathbf{u}; \boldsymbol{\theta})$  wrt  $\boldsymbol{\theta}$ .

**Proof** With the definition of  $J$  in Eq (10),

$$\begin{aligned} J(\boldsymbol{\theta} + \epsilon\boldsymbol{\varphi}) &= \int \ln[r_\nu(G(\mathbf{u}; \boldsymbol{\theta} + \epsilon\boldsymbol{\varphi}))] p_d(\mathbf{u}) d\mathbf{u} + \\ &\quad \nu \int \ln[1 - r_\nu(G(\mathbf{u}; \boldsymbol{\theta} + \epsilon\boldsymbol{\varphi}))] p_n(\mathbf{u}) d\mathbf{u}. \end{aligned}$$

Developing  $G(\mathbf{u}; \boldsymbol{\theta} + \epsilon\boldsymbol{\varphi})$  till terms of order  $\epsilon^2$  yields

$$G(\mathbf{u}; \boldsymbol{\theta} + \epsilon\boldsymbol{\varphi}) = G(\mathbf{u}; \boldsymbol{\theta}) + \epsilon\boldsymbol{\varphi}^T \mathbf{g}(\mathbf{u}; \boldsymbol{\theta}) + \epsilon^2 \frac{1}{2} \boldsymbol{\varphi}^T \mathbf{H}_G(\mathbf{u}; \boldsymbol{\theta}) \boldsymbol{\varphi} + O(\epsilon^3).$$

Defining  $u_1$  and  $u_2$  as in the lemma, we obtain

$$\ln r_\nu(G(\mathbf{u}; \boldsymbol{\theta} + \epsilon\boldsymbol{\varphi})) = \ln r_\nu \left( G(\mathbf{u}; \boldsymbol{\theta}) + \epsilon u_1 + \epsilon^2 u_2 + O(\epsilon^3) \right).$$

Using now the Taylor expansions in Eq (30) and (31) for  $u = G(\mathbf{u}; \boldsymbol{\theta})$ ,

■

**Lemma 10** If  $p_n \gg p_d$  and

$$\mathcal{I}_\nu = \int \mathbf{g}(\mathbf{u}) \mathbf{g}(\mathbf{u})^T P_\nu(\mathbf{u}) p_d(\mathbf{u}) d\mathbf{u}$$

is full rank (PD), where

$$\begin{aligned} P_\nu(\mathbf{u}) &= \frac{\nu p_n(\mathbf{u})}{p_d(\mathbf{u}) + \nu p_n(\mathbf{u})}, \\ \mathbf{g}(\mathbf{u}) &= \nabla_{\boldsymbol{\theta}} \ln p_m(\mathbf{u}; \boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*}, \end{aligned}$$

then

$$J(\boldsymbol{\theta}^*) > J(\boldsymbol{\theta}) \quad \boldsymbol{\theta} \neq \boldsymbol{\theta}^*.$$

**Proof** A necessary condition for optimality is that in the expansion of  $J(\boldsymbol{\theta} + \epsilon \boldsymbol{\varphi})$  in Lemma 9, the term of order  $\epsilon = 0$  for any  $\boldsymbol{\varphi}$  iff

$$\begin{aligned} p_d(\mathbf{u})(1 - h(\mathbf{u}; \boldsymbol{\theta})) &= \nu p_n(\mathbf{u}) h(\mathbf{u}; \boldsymbol{\theta}), \\ \Leftrightarrow & \\ \frac{\nu p_n(\mathbf{u}) p_d(\mathbf{u})}{p_m(\mathbf{u}; \boldsymbol{\theta}) + \nu p_n(\mathbf{u})} &= \frac{\nu p_n(\mathbf{u}) p_m(\mathbf{u}; \boldsymbol{\theta})}{p_m(\mathbf{u}; \boldsymbol{\theta}) + \nu p_n(\mathbf{u})}, \end{aligned}$$

where we have used Eq (28) and (29) as in the proof for Lemma 8. The assumption that  $\nu > 0$  and  $p_d(\cdot) = p_m(\cdot; \boldsymbol{\theta}^*)$  together with the above equation  $\Rightarrow$  the term of order  $\epsilon = 0$  if  $\boldsymbol{\theta} = \boldsymbol{\theta}^*$ .

The objective function

$$\begin{aligned} J(\boldsymbol{\theta}^* + \epsilon \boldsymbol{\varphi}) &= J(\boldsymbol{\theta}^*) - \frac{\epsilon^2}{2} \int u_1^2 (1 - h(\mathbf{u}; \boldsymbol{\theta}^*)) h(\mathbf{u}; \boldsymbol{\theta}^*) \\ &\quad (p_d(\mathbf{u}) + \nu p_n(\mathbf{u})) d\mathbf{u} + O(\epsilon^3). \end{aligned}$$

The terms  $h(\mathbf{u}; \boldsymbol{\theta}^*)$  and  $1 - h(\mathbf{u}; \boldsymbol{\theta}^*)$  are with Equation (27)

$$h(\mathbf{u}; \boldsymbol{\theta}^*) = \frac{p_d(\mathbf{u})}{p_d(\mathbf{u}) + \nu p_n(\mathbf{u})}, \quad 1 - h(\mathbf{u}; \boldsymbol{\theta}^*) = \frac{\nu p_n(\mathbf{u})}{p_d(\mathbf{u}) + \nu p_n(\mathbf{u})}.$$

$\Rightarrow$

$$J(\boldsymbol{\theta}^* + \epsilon \boldsymbol{\varphi}) = J(\boldsymbol{\theta}^*) - \frac{\epsilon^2}{2} \boldsymbol{\varphi}^T \left[ \int \mathbf{g}(\mathbf{u}) \mathbf{g}(\mathbf{u})^T P_\nu(\mathbf{u}) p_d(\mathbf{u}) d\mathbf{u} \right] \boldsymbol{\varphi} + O(\epsilon^3)$$

by inserting the definition of  $u_1$  evaluated at  $\boldsymbol{\theta}^*$ . The term of order  $\epsilon^2$  defines the nature of the extremum at  $\boldsymbol{\theta}^*$ . If  $\mathcal{I}_\nu > 0$ ,  $J(\boldsymbol{\theta}^*)$  is a maximum.

By definition of  $J$ , for a suitable noise density  $p_n$ ,  $J$  attains a global maximum at  $\boldsymbol{\theta}^*$ .

■

### A.3.2 PROOF OF THE THEOREM

The proof of consistency goes along the same lines as MLE ( Wasserman, 2004, Chapter 9).

**Proof of consistency**, given  $\epsilon > 0$ ,  $P(||\hat{\theta}_T - \theta^*|| > \epsilon) \rightarrow 0$  as  $T_d \rightarrow \infty$ . In what follows, it is sometimes useful to make the underlying probability space explicit and write  $P(||\hat{\theta}_T - \theta^*|| > \epsilon)$ .  
 by Lemma 10,  $\hat{\theta}_T$  is a global maximum,  $||\hat{\theta} - \theta^*|| > \epsilon \implies$  there is a  $\delta(\epsilon)$  such that  $J(\hat{\theta}) < J(\theta^*) - \delta(\epsilon)$ . Hence,

$$\{||\hat{\theta}_T - \theta^*|| > \epsilon\} \subset \{J(\hat{\theta}_T) < J(\theta^*) - \delta(\epsilon)\}$$

and thus

$$P(||\hat{\theta}_T - \theta^*|| > \epsilon) < P(J(\hat{\theta}_T) < J(\theta^*) - \delta(\epsilon)). \quad (32)$$

Next, we investigate what happens to  $P(J(\hat{\theta}_T) < J(\theta^*) - \delta(\epsilon))$ ,  $T_d \rightarrow \infty$ .

Fact

$$\begin{aligned} J(\theta^*) - J(\hat{\theta}_T) &= J(\theta^*) - J_T(\theta^*) + J_T(\theta^*) - J(\hat{\theta}_T) \\ &\leq J(\theta^*) - J_T(\theta^*) + J_T(\hat{\theta}_T) - J(\hat{\theta}_T) \end{aligned}$$

as  $\hat{\theta}_T$  maximizes  $J_T$

$$\implies |J(\theta^*) - J(\hat{\theta}_T)| \leq 2 \sup_{\theta} |J(\theta) - J_T(\theta)|, \quad |$$

$$\implies P(|J(\theta^*) - J(\hat{\theta}_T)| > \delta) \leq P(2 \sup_{\theta} |J(\theta) - J_T(\theta)| > \delta).$$

Using the assumption that  $J_T(\theta)$  converges in probability uniformly over  $\theta$  to  $J(\theta)$ , we obtain that for large  $T_d$

$$P(|J(\theta^*) - J(\hat{\theta}_T)| > \delta) < \epsilon_2$$

As  $J(\theta^*) \geq J(\theta)$  for any  $\theta$ , we have

$$P(J(\hat{\theta}_T) < J(\theta^*) - \delta) < \epsilon_2$$

for any  $\epsilon_2 > 0$ .

■

### A.4 Proof of Theorem 3 (Asymptotic Normality)

#### A.4.1 LEMMATA

##### Lemma 11

$$0 = \nabla_{\theta} J_T(\theta^*) + \mathbf{H}_J(\theta^*)(\hat{\theta}_T - \theta^*) + O(||\hat{\theta}_T - \theta^*||^2)$$

where

$$\begin{aligned}\nabla_{\theta} J_T(\theta^*) &= \frac{1}{T_d} \sum_{t=1}^{T_d} (1 - h(\mathbf{x}_t; \theta^*)) \mathbf{g}(\mathbf{x}_t) - \nu \frac{1}{T_n} \sum_{t=1}^{T_n} h(\mathbf{y}_t; \theta^*) \mathbf{g}(\mathbf{y}_t), \\ \mathbf{H}_J(\theta^*) &= \frac{1}{T_d} \sum_{t=1}^{T_d} \left\{ -(1 - h(\mathbf{x}_t; \theta^*)) h(\mathbf{x}_t; \theta^*) \mathbf{g}(\mathbf{x}_t) \mathbf{g}(\mathbf{x}_t)^T + \right. \\ &\quad (1 - h(\mathbf{x}_t; \theta^*)) \mathbf{H}_G(\mathbf{x}_t; \theta^*) \} - \\ &\quad \nu \frac{1}{T_n} \sum_{t=1}^{T_n} \left\{ (1 - h(\mathbf{y}_t; \theta^*)) h(\mathbf{y}_t; \theta^*) \mathbf{g}(\mathbf{y}_t) \mathbf{g}(\mathbf{y}_t)^T + \right. \\ &\quad \left. h(\mathbf{y}_t; \theta^*) \mathbf{H}_G(\mathbf{y}_t; \theta^*) \right\}.\end{aligned}$$

**Proof** Using the chain rule, it follows from the relations in Section A.1 that

$$\begin{aligned}\nabla_{\theta} \ln h(\mathbf{x}_t; \theta) &= (1 - h(\mathbf{x}_t; \theta)) \mathbf{g}(\mathbf{x}_t; \theta) \\ \nabla_{\theta} \ln [1 - h(\mathbf{y}_t; \theta)] &= -h(\mathbf{y}_t; \theta) \mathbf{g}(\mathbf{y}_t; \theta).\end{aligned}$$

$\Rightarrow$

$$\boxed{\nabla_{\theta} J_T(\theta) = \frac{1}{T_d} \sum_{t=1}^{T_d} (1 - h(\mathbf{x}_t; \theta)) \mathbf{g}(\mathbf{x}; \theta) - \nu \frac{1}{T_n} \sum_{t=1}^{T_n} h(\mathbf{y}_t; \theta) \mathbf{g}(\mathbf{y}_t; \theta).}$$

As  $\hat{\theta}_T$  is the value of  $\theta$  which maximizes  $J_T(\theta)$ , must  $= 0$ .  $\Rightarrow$

$$0 = \nabla_{\theta} J_T(\theta^*) + \mathbf{H}_J(\theta^*)(\hat{\theta}_T - \theta^*) + O(||\hat{\theta}_T - \theta^*||^2). \quad |$$

The  $k$ -th row of the Hessian  $\mathbf{H}_J(\theta)$  is  $\nabla_{\theta} F_k(\theta)^T$  where  $F_k$  is the  $k$ -th element of the vector  $\nabla_{\theta} J_T$ . Denoting by  $g_k$  the  $k$ -th element of the score function  $\mathbf{g}$ , we have

$$\begin{aligned}\nabla_{\theta} F_k(\theta) &= \frac{1}{T_d} \sum_{t=1}^{T_d} \left\{ -\nabla_{\theta} h(\mathbf{x}_t; \theta) g_k(\mathbf{x}_t; \theta) + (1 - h(\mathbf{x}_t; \theta)) \nabla_{\theta} g_k(\mathbf{x}_t; \theta) \right\} \\ &\quad - \nu \frac{1}{T_n} \sum_{t=1}^{T_n} \left\{ \nabla_{\theta} h(\mathbf{y}_t; \theta) g_k(\mathbf{y}_t; \theta) + h(\mathbf{y}_t; \theta) \nabla_{\theta} g_k(\mathbf{x}_t; \theta) \right\}.\end{aligned}$$

Using the chain rule, it follows from the relations in Section A.1 that

$$\nabla_{\theta} h(\mathbf{u}; \boldsymbol{\theta}) = (1 - h(\mathbf{u}; \boldsymbol{\theta}))h(\mathbf{u}; \boldsymbol{\theta})\mathbf{g}(\mathbf{u}; \boldsymbol{\theta}).$$

Hence,

$$\begin{aligned} \nabla_{\theta} F_k(\boldsymbol{\theta}) &= \frac{1}{T_d} \sum_{t=1}^{T_d} \left\{ -(1 - h(\mathbf{x}_t; \boldsymbol{\theta}))h(\mathbf{x}_t; \boldsymbol{\theta})\mathbf{g}(\mathbf{x}_t; \boldsymbol{\theta})g_k(\mathbf{x}_t; \boldsymbol{\theta}) + \right. \\ &\quad \left. (1 - h(\mathbf{x}_t; \boldsymbol{\theta}))\nabla_{\theta} g_k(\mathbf{x}_t; \boldsymbol{\theta}) \right\} - \\ &\quad \nu \frac{1}{T_n} \sum_{t=1}^{T_n} \left\{ (1 - h(\mathbf{y}_t; \boldsymbol{\theta}))h(\mathbf{y}_t; \boldsymbol{\theta})\mathbf{g}(\mathbf{y}_t; \boldsymbol{\theta})g_k(\mathbf{y}_t; \boldsymbol{\theta}) + \right. \\ &\quad \left. h(\mathbf{y}_t; \boldsymbol{\theta})\nabla_{\theta} g_k(\mathbf{y}_t; \boldsymbol{\theta}) \right\}, \end{aligned}$$

■

**Lemma 12**  $\mathbf{H}_J(\boldsymbol{\theta}^*)$  converges in probability to  $-\mathcal{I}_{\nu}$ , the sample size  $T_d \rightarrow \infty$ .

**Proof**  $T_n = \nu T_d \rightarrow \infty$ . As the sample sizes become arbitrarily large, the sample averages become integration over the corresponding densities so that

$$\begin{aligned} \lim_{T_d \rightarrow \infty} \mathbf{H}_J(\boldsymbol{\theta}^*) &\xrightarrow{P} - \int (1 - h(\mathbf{u}; \boldsymbol{\theta}^*))h(\mathbf{u}; \boldsymbol{\theta}^*)\mathbf{g}(\mathbf{u})\mathbf{g}(\mathbf{u})^T(p_d(\mathbf{u}) + \nu p_n(\mathbf{u}))d\mathbf{u} + \\ &\quad \int ((1 - h(\mathbf{u}; \boldsymbol{\theta}^*))p_d(\mathbf{u}) - h(\mathbf{u}; \boldsymbol{\theta}^*)\nu p_n(\mathbf{u}))\mathbf{H}_G(\mathbf{u}; \boldsymbol{\theta}^*)d\mathbf{u}. \end{aligned}$$

With Eq (28) and (29), we have

$$\begin{aligned} (1 - h(\mathbf{u}; \boldsymbol{\theta}^*))p_d(\mathbf{u}) &= h(\mathbf{u}; \boldsymbol{\theta}^*)\nu p_n(\mathbf{u}), \\ (1 - h(\mathbf{u}; \boldsymbol{\theta}^*))h(\mathbf{u}; \boldsymbol{\theta}^*)(p_d(\mathbf{u}) + \nu p_n(\mathbf{u})) &= \frac{\nu p_n(\mathbf{u})p_d(\mathbf{u})}{p_d(\mathbf{u}) + \nu p_n(\mathbf{u})}. \end{aligned} \tag{33}$$

Hence,

$$\lim_{T_d \rightarrow \infty} \mathbf{H}_J(\boldsymbol{\theta}^*) \xrightarrow{P} - \int \frac{\nu p_n(\mathbf{u})p_d(\mathbf{u})}{p_d(\mathbf{u}) + \nu p_n(\mathbf{u})}\mathbf{g}(\mathbf{u})\mathbf{g}(\mathbf{u})^T d\mathbf{u},$$

which is  $-\mathcal{I}_{\nu}$ .

■

**Lemma 13**  $E \nabla_{\theta} J_T(\theta^*) = 0$ .

**Proof**

$$\begin{aligned} E \nabla_{\theta} J_T(\theta^*) &= \frac{1}{T_d} \sum_{t=1}^{T_d} E \mathbf{g}(\mathbf{x}_t) (1 - h(\mathbf{x}_t; \theta^*)) - \nu \frac{1}{T_n} \sum_{t=1}^{T_n} E \mathbf{g}(\mathbf{y}_t) h(\mathbf{y}_t; \theta^*) \\ &= E \mathbf{g}(\mathbf{x}) (1 - h(\mathbf{x}; \theta^*)) - \nu E \mathbf{g}(\mathbf{y}) h(\mathbf{y}; \theta^*) \\ &= \int \mathbf{g}(\mathbf{u}) (1 - h(\mathbf{u}; \theta^*)) p_d(\mathbf{u}) d\mathbf{u} - \nu \int \mathbf{g}(\mathbf{u}) h(\mathbf{u}; \theta^*) p_n(\mathbf{u}) d\mathbf{u}, \end{aligned}$$

from the i.i.d. assumption of the sample  $X$  and  $Y$ .  $\Rightarrow$

$$E \nabla_{\theta} J_T(\theta^*) = \int \mathbf{g}(\mathbf{u}) ((1 - h(\mathbf{u}; \theta^*)) p_d(\mathbf{u}) - h(\mathbf{u}; \theta^*) \nu p_n(\mathbf{u})) d\mathbf{u},$$

■

**Lemma 14**  $\text{Var} \nabla_{\theta} J_T(\theta^*) = \frac{1}{T_d} \left( \mathcal{I}_{\nu} - \left( 1 + \frac{1}{\nu} \right) E(P_{\nu} \mathbf{g}) E(P_{\nu} \mathbf{g})^T \right),$

where  $\mathcal{I}_{\nu}$ ,  $P_{\nu}$  and  $\mathbf{g}$  were defined in Lemma 10, and the expectation is taken over the data-pdf  $p_d$ .

**Proof** As  $E \nabla_{\theta} J_T(\theta^*) = 0$ , the variance is  $E \nabla_{\theta} J_T(\theta^*) \nabla_{\theta} J_T(\theta^*)^T$ . Multiplying out gives

$$\begin{aligned} \text{Var} \nabla_{\theta} J_T(\theta^*) &= \frac{1}{T_d^2} E \left[ \sum_{t=1}^{T_d} (1 - h(\mathbf{x}_t; \theta^*)) \mathbf{g}(\mathbf{x}_t) \sum_{t=1}^{T_d} (1 - h(\mathbf{x}_t; \theta^*)) \mathbf{g}(\mathbf{x}_t)^T \right] - \\ &\quad \frac{1}{T_d^2} E \left[ \sum_{t=1}^{T_d} (1 - h(\mathbf{x}_t; \theta^*)) \mathbf{g}(\mathbf{x}_t) \sum_{t=1}^{T_n} h(\mathbf{y}_t; \theta^*) \mathbf{g}(\mathbf{y}_t)^T \right] - \\ &\quad \frac{1}{T_d^2} E \left[ \sum_{t=1}^{T_n} h(\mathbf{y}_t; \theta^*) \mathbf{g}(\mathbf{y}_t) \sum_{t=1}^{T_d} (1 - h(\mathbf{x}_t; \theta^*)) \mathbf{g}(\mathbf{x}_t)^T \right] + \\ &\quad \frac{1}{T_d^2} E \left[ \sum_{t=1}^{T_n} h(\mathbf{y}_t; \theta^*) \mathbf{g}(\mathbf{y}_t) \sum_{t=1}^{T_n} h(\mathbf{y}_t; \theta^*) \mathbf{g}(\mathbf{y}_t)^T \right]. \end{aligned}$$

Since the samples are all independent from each other, we have

$$\begin{aligned}
\text{Var} \nabla_{\boldsymbol{\theta}} J_T(\boldsymbol{\theta}^*) &= \frac{1}{T_d^2} \sum_{t=1}^{T_d} \mathbb{E} \left[ (1 - h(\mathbf{x}_t; \boldsymbol{\theta}^*))^2 \mathbf{g}(\mathbf{x}_t) \mathbf{g}(\mathbf{x}_t)^T \right] + \\
&\quad \frac{1}{T_d^2} \sum_{\substack{t, \tau=1 \\ t \neq \tau}}^{T_d} \mathbb{E}[(1 - h(\mathbf{x}_t; \boldsymbol{\theta}^*)) \mathbf{g}(\mathbf{x}_t)] \mathbb{E}[(1 - h(\mathbf{x}_\tau; \boldsymbol{\theta}^*)) \mathbf{g}(\mathbf{x}_\tau)^T] - \\
&\quad \frac{1}{T_d^2} \sum_{t=1}^{T_d} \sum_{\tau=1}^{T_n} \mathbb{E}[(1 - h(\mathbf{x}_t; \boldsymbol{\theta}^*)) \mathbf{g}(\mathbf{x}_t)] \mathbb{E}[h(\mathbf{y}_\tau; \boldsymbol{\theta}^*) \mathbf{g}(\mathbf{y}_\tau)^T] - \\
&\quad \frac{1}{T_d^2} \sum_{t=1}^{T_n} \sum_{\tau=1}^{T_d} \mathbb{E}[h(\mathbf{y}_t; \boldsymbol{\theta}^*) \mathbf{g}(\mathbf{y}_t)] \mathbb{E}[(1 - h(\mathbf{x}_\tau; \boldsymbol{\theta}^*)) \mathbf{g}(\mathbf{x}_\tau)^T] + \\
&\quad \frac{1}{T_d^2} \sum_{\substack{t, \tau=1 \\ t \neq \tau}}^{T_n} \mathbb{E}[h(\mathbf{y}_t; \boldsymbol{\theta}^*) \mathbf{g}(\mathbf{y}_t)] \mathbb{E}[h(\mathbf{y}_\tau; \boldsymbol{\theta}^*) \mathbf{g}(\mathbf{y}_\tau)^T] + \\
&\quad \frac{1}{T_d^2} \sum_{t=1}^{T_n} \mathbb{E}[h(\mathbf{y}_t; \boldsymbol{\theta}^*)^2 \mathbf{g}(\mathbf{y}_t) \mathbf{g}(\mathbf{y}_t)^T].
\end{aligned}$$

As we assume that all  $\mathbf{x}_t$ , and also  $\mathbf{y}_t$ , are identically distributed, the above expression simplifies to

$$\begin{aligned}
\text{Var} \nabla_{\boldsymbol{\theta}} J_T(\boldsymbol{\theta}^*) &= \frac{1}{T_d} \int (1 - h(\mathbf{u}; \boldsymbol{\theta}^*))^2 \mathbf{g}(\mathbf{u}) \mathbf{g}(\mathbf{u})^T p_d(\mathbf{u}) d\mathbf{u} + \\
&\quad \frac{T_d^2 - T_d}{T_d^2} \mathbf{m}_x \mathbf{m}_x^T - \frac{T_d T_n}{T_d^2} \mathbf{m}_x \mathbf{m}_y^T - \\
&\quad \frac{T_d T_n}{T_d^2} \mathbf{m}_y \mathbf{m}_x^T + \frac{T_n^2 - T_n}{T_d^2} \mathbf{m}_y \mathbf{m}_y^T + \\
&\quad \frac{T_n}{T_d^2} \int h(\mathbf{u}; \boldsymbol{\theta}^*)^2 \mathbf{g}(\mathbf{u}) \mathbf{g}(\mathbf{u})^T p_n(\mathbf{u}) d\mathbf{u}, \tag{34}
\end{aligned}$$

where

$$\begin{aligned}
\mathbf{m}_x &= \int (1 - h(\mathbf{u}; \boldsymbol{\theta}^*)) \mathbf{g}(\mathbf{u}) p_d(\mathbf{u}) d\mathbf{u}, \\
\mathbf{m}_y &= \int h(\mathbf{u}; \boldsymbol{\theta}^*) \mathbf{g}(\mathbf{u}) p_n(\mathbf{u}) d\mathbf{u}.
\end{aligned}$$

Denoting by  $A$  the sum of the first and last line of Eq (34), we have

$$A = \frac{1}{T_d} \int \mathbf{g}(\mathbf{u}) \mathbf{g}(\mathbf{u})^T \left[ (1 - h(\mathbf{u}; \boldsymbol{\theta}^*))^2 p_d(\mathbf{u}) + h(\mathbf{u}; \boldsymbol{\theta}^*)^2 \nu p_n(\mathbf{u}) \right] d\mathbf{u}$$

since  $T_n = \nu T_d$ . Now, Eq (27) and  $p_m(\mathbf{u}; \boldsymbol{\theta}^*) = p_d(\mathbf{u}) \implies$

$$\begin{aligned}
(1 - h(\mathbf{u}; \boldsymbol{\theta}^*))^2 p_d(\mathbf{u}) + h(\mathbf{u}; \boldsymbol{\theta}^*)^2 \nu p_n(\mathbf{u}) &= \frac{\nu p_n(\mathbf{u}) p_d(\mathbf{u})}{p_d(\mathbf{u}) + \nu p_n(\mathbf{u})} \\
&= P_\nu p_d(\mathbf{u}),
\end{aligned}$$

$\implies$

$$\begin{aligned} A &= \frac{1}{T_d} \int \mathbf{g}(\mathbf{u}) \mathbf{g}(\mathbf{u})^T P_\nu p_d(\mathbf{u}) d\mathbf{u} \\ &= \frac{1}{T_d} \mathcal{I}_\nu. \end{aligned}$$

Denote by  $B$  the second line of Eq (34). Rearranging the terms, we have

$$\begin{aligned} B &= \mathbf{m}_x \int [(1 - h(\mathbf{u}; \boldsymbol{\theta}^*)) p_d(\mathbf{u}) - h(\mathbf{u}; \boldsymbol{\theta}^*) \nu p_n(\mathbf{u})] \mathbf{g}(\mathbf{u})^T d\mathbf{u} - \\ &\quad \frac{1}{T_d} \mathbf{m}_x \mathbf{m}_x^T. \end{aligned} \tag{35}$$

Again, Eq (27) and  $p_m(\mathbf{u}; \boldsymbol{\theta}^*) = p_d(\mathbf{u}) \implies$

$$\begin{aligned} (1 - h(\mathbf{u}; \boldsymbol{\theta}^*)) p_d(\mathbf{u}) &= h(\mathbf{u}; \boldsymbol{\theta}^*) \nu p_n(\mathbf{u}) \\ &= \frac{\nu p_n(\mathbf{u}) p_d(\mathbf{u})}{p_d(\mathbf{u}) + \nu p_n(\mathbf{u})} \\ &= P_\nu p_d(\mathbf{u}), \end{aligned}$$

so that the first line in Eq (35) is zero and

$$\mathbf{m}_x = \int P_\nu \mathbf{g}(\mathbf{u}) p_d(\mathbf{u}) d\mathbf{u}.$$

Thus

$$B = -\frac{1}{T_d} \int P_\nu \mathbf{g}(\mathbf{u}) p_d(\mathbf{u}) d\mathbf{u} \int P_\nu \mathbf{g}(\mathbf{u})^T p_d(\mathbf{u}) d\mathbf{u}.$$

Denote by  $C$  the third line of Eq (34). Rearranging the terms, we have with  $T_n = \nu T_d$

$$C = -\frac{\nu}{T_d} \mathbf{m}_y \mathbf{m}_y^T + \nu \mathbf{m}_y (\nu \mathbf{m}_y^T - \mathbf{m}_x^T).$$

The term  $\nu \mathbf{m}_y$  is with Equation (27) and  $p_m(\mathbf{u}; \boldsymbol{\theta}^*) = p_d(\mathbf{u})$

$$\nu \mathbf{m}_y = \int P_\nu \mathbf{g}(\mathbf{u}) p_d(\mathbf{u}) d\mathbf{u},$$

so that  $\nu \mathbf{m}_y = \mathbf{m}_x$ , and hence

$$\begin{aligned} C &= -\frac{1}{\nu T_d} (\nu \mathbf{m}_y) (\nu \mathbf{m}_y^T) \\ &= \frac{1}{\nu} B. \end{aligned}$$

All in all, the variance  $\text{Var} \nabla_{\boldsymbol{\theta}} J_T(\boldsymbol{\theta}^*)$  is thus

$$\begin{aligned} \text{Var} \nabla_{\boldsymbol{\theta}} J_T(\boldsymbol{\theta}^*) &= A + B + C \\ &= \frac{1}{T_d} \left( \mathcal{I}_\nu - \left( 1 + \frac{1}{\nu} \right) \mathbb{E}(P_\nu \mathbf{g}) \mathbb{E}(P_\nu \mathbf{g}^T) \right), \end{aligned}$$

where

$$\mathbb{E}(P_\nu \mathbf{g}) = \int P_\nu \mathbf{g}(\mathbf{u}) p_d(\mathbf{u}) d\mathbf{u}.$$

■

#### A.4.2 PROOF OF THE THEOREM

We are now ready to give the proof of Theorem 3.

**Proof** Up to terms of order  $O(||\hat{\theta}_T - \theta^*||^2)$ , we have with Lemma 11

$$\sqrt{T_d}(\hat{\theta}_T - \theta^*) = -\mathbf{H}_J^{-1}\sqrt{T_d}\nabla_{\theta}J_T(\theta^*).$$

By Lemma 12,  $\mathbf{H}_J \xrightarrow{P} -\mathcal{I}_{\nu}$  for large sample sizes  $T_d$ . Using Lemma 13 and Lemma 14, we see that

$$\sqrt{T_d}\nabla_{\theta}J_T(\theta^*)$$

converges in distribution to a normal distribution of mean zero and covariance matrix

$$\mathcal{I}_{\nu} - \left(1 + \frac{1}{\nu}\right)\mathbf{E}(P_{\nu}\mathbf{g})\mathbf{E}(P_{\nu}\mathbf{g})^T,$$

which implies that  $\sqrt{T_d}(\hat{\theta}_T - \theta^*)$  converges in distribution to a normal distribution of mean zero and covariance matrix  $\Sigma$ ,

$$\Sigma = \mathcal{I}_{\nu}^{-1} - \left(1 + \frac{1}{\nu}\right)\mathcal{I}_{\nu}^{-1}\mathbf{E}(P_{\nu}\mathbf{g})\mathbf{E}(P_{\nu}\mathbf{g})^T\mathcal{I}_{\nu}^{-1}.$$

■

## Appendix B. Calculations

The following sections contain calculations needed in Section 3.3 and Section 5.3.

### B.1 Theory, Section 3.3: Asymptotic Variance for Orthogonal ICA Model

We calculate here the asymptotic covariance matrix of the estimation error for an orthogonal ICA model when a Gaussian distribution is used as noise distribution in NCE. This result is used to make the predictions about the estimation error in Section 3.3. The calculations show that the asymptotic variance does not depend on the mixing matrix but only on the dim of the data. Similar calculations can be used to show that this also holds for MLE.

A rv  $\mathbf{x}$  following an ICA model with orthogonal mixing matrix  $\mathbf{A} = (\mathbf{a}_1 \dots \mathbf{a}_n)$  has the distribution:

$$p_d(\mathbf{x}) = \frac{1}{Z} \prod_{i=1}^n f(\mathbf{a}_i^T \mathbf{x}),$$

where  $Z$  is the partition function. By orthogonality of  $\mathbf{A}$ ,

$$p_d(\mathbf{Ax}) = \frac{1}{Z} \prod_{i=1}^n f(x_i),$$

$= p_s(\mathbf{x})$  the distribution of the sources  $\mathbf{s}$  of the ICA model. Also by orthogonality of  $\mathbf{A}$ , the noise distribution  $p_n$  with the same covariance as  $\mathbf{x}$  is the std. In particular,  $p_n(\mathbf{Ax}) = p_n(\mathbf{x})$ .

For the calculation of the asymptotic variance, we need to compute the matrix  $\mathcal{I}_\nu$  which occurs in Theorem 2,  $\mathcal{I}_\nu = \int \mathbf{g}(\mathbf{u})\mathbf{g}(\mathbf{u})^T P_\nu(\mathbf{u}) p_d(\mathbf{u}) d\mathbf{u}$ . With the above data and noise distribution,  $P_\nu(\mathbf{u})$  has the property that

$$\begin{aligned} P_\nu(\mathbf{A}\mathbf{u}) &= \frac{\nu p_n(\mathbf{A}\mathbf{u})}{p_d(\mathbf{A}\mathbf{u}) + \nu p_n(\mathbf{A}\mathbf{u})} \\ &= \frac{\nu p_n(\mathbf{u})}{p_s(\mathbf{u}) + \nu p_n(\mathbf{u})}. \end{aligned}$$

Hence  $P_\nu(\mathbf{A}\mathbf{u})$  does not depend on  $\mathbf{A}$ . Below, we will denote  $P_\nu(\mathbf{A}\mathbf{u})$  by  $\tilde{P}_\nu(\mathbf{u})$ . For the ICA model, the vector  $\mathbf{g}(\mathbf{u})$  has the form

$$\mathbf{g}(\mathbf{u}) = (\mathbf{g}_1(\mathbf{u}), \dots, \mathbf{g}_n(\mathbf{u}), g_c(\mathbf{u}))^T$$

where  $\mathbf{g}_i(\mathbf{u}) = \nabla_{\mathbf{a}_i} \ln p_m(\mathbf{u}) = f'(\mathbf{a}_i^T \mathbf{u}) \mathbf{u}$  and  $g_c(\mathbf{u}) = \partial_c \ln p_m(\mathbf{u}) = 1$ . By orthogonality of  $\mathbf{A}$ , we have

$$\mathbf{g}_i(\mathbf{A}\mathbf{u}) = \mathbf{A}f'(u_i)\mathbf{u}.$$

We denote the vector  $f'(u_i)\mathbf{u}$  by  $\tilde{\mathbf{g}}_i(\mathbf{u})$  so that  $\mathbf{g}_i(\mathbf{A}\mathbf{u}) = \mathbf{A}\tilde{\mathbf{g}}_i(\mathbf{u})$ . Hence,

$$\mathbf{g}(\mathbf{A}\mathbf{u}) = \mathcal{A}(\tilde{\mathbf{g}}_1(\mathbf{u}), \dots, \tilde{\mathbf{g}}_n(\mathbf{u}), 1)^T$$

where  $\mathcal{A}$  is a block-diagonal matrix with  $n$  matrices  $\mathbf{A}$  on the diagonal and a single 1 in the  $(n+1)$ -th slot. As a shorthand, we will denote  $\mathbf{g}(\mathbf{A}\mathbf{u})$  by  $\mathcal{A}\tilde{\mathbf{g}}(\mathbf{u})$ .

With these preliminaries, using the change of variables  $\mathbf{u} = \mathbf{Av}$ ,

$$\begin{aligned} \mathcal{I}_\nu &= \int p_d(\mathbf{u})\mathbf{g}(\mathbf{u})\mathbf{g}(\mathbf{u})^T P_\nu(\mathbf{u}) d\mathbf{u} \\ &= \int p_s(\mathbf{v})\mathcal{A}\tilde{\mathbf{g}}(\mathbf{v})\tilde{\mathbf{g}}(\mathbf{v})^T \mathcal{A}^T \tilde{P}_\nu(\mathbf{v}) d\mathbf{v} \\ &= \mathcal{A}\tilde{\mathcal{I}}_\nu\mathcal{A}^T, \end{aligned}$$

where the matrix

$$\tilde{\mathcal{I}}_\nu = \int p_s(\mathbf{v})\tilde{\mathbf{g}}(\mathbf{v})\tilde{\mathbf{g}}(\mathbf{v})^T \tilde{P}_\nu(\mathbf{v}) d\mathbf{v}$$

does not depend on the mixing matrix  $\mathbf{A}$  but only on the distribution of the sources  $\mathbf{s}$ , the noise distribution  $p_n$ , and  $\nu$ . Moreover, by orthogonality of  $\mathbf{A}$ , the inverse of  $\mathcal{I}_\nu$  is given by

$$\mathcal{I}_\nu^{-1} = \mathcal{A}\tilde{\mathcal{I}}_\nu^{-1}\mathcal{A}^T.$$

The same reasoning shows that

$$\int p_d(\mathbf{u})P_\nu(\mathbf{u})\mathbf{g}(\mathbf{u}) d\mathbf{u} = \mathcal{A} \int p_s(\mathbf{v})\tilde{\mathbf{g}}(\mathbf{v})\tilde{P}_\nu(\mathbf{v}) d\mathbf{v},$$

which we will denote below by  $\mathcal{A}\tilde{\mathbf{m}}$ . Again,  $\tilde{\mathbf{m}}$  does not depend on  $\mathbf{A}$ . Hence, the asymptotic covariance matrix  $\Sigma$ ,

$$\Sigma = \mathcal{I}_\nu^{-1} - \left(1 + \frac{1}{\nu}\right) \mathcal{I}_\nu^{-1} \mathbb{E}(P_\nu \mathbf{g}) \mathbb{E}(P_\nu \mathbf{g})^T \mathcal{I}_\nu^{-1},$$

in Theorem 3 is for the ICA model with orthogonal mixing matrix  $\mathbf{A}$  given by

$$\Sigma_{\text{ortICA}} = \mathcal{A} \left[ \tilde{\mathcal{I}}_\nu^{-1} - \left( 1 + \frac{1}{\nu} \right) \tilde{\mathcal{I}}_\nu^{-1} \tilde{\mathbf{m}} \tilde{\mathbf{m}}^T \tilde{\mathcal{I}}_\nu^{-1} \right] \mathcal{A}^T.$$

The block matrix  $\mathcal{A}$  is orthogonal since  $\mathbf{A}$  is orthogonal. The asymptotic variance, that is the trace of  $\Sigma_{\text{ortICA}}$ , does hence not depend on  $\mathbf{A}$ .

## B.2 Natural Images, Section 5.3: Optimal Stimuli

We show here that the optimal stimulus, namely the image which yields the largest feature output for feature  $\mathbf{w}$  while satisfying the sphere constraints in Eq (22), is proportional to  $\mathbf{V}^-(\mathbf{w} - \langle \mathbf{w} \rangle)$ .

$\langle \mathbf{w} \rangle :=$  the average value of the elements  $\mathbf{w}$ .

Each coordinate vector  $\mathbf{x}$  defines an image  $\mathbf{i} = \mathbf{V}^-\mathbf{x}$ , see Eq (23). The optimal image is thus  $\mathbf{i}^* = \mathbf{V}^-\mathbf{x}^*$  where  $\mathbf{x}^*$  is the solution to the optimization problem

$$\max_{\mathbf{x}} \mathbf{w}^T \mathbf{x}$$

subject to  $\sum_{k=1}^n \mathbf{x}(k) = 0$  and  $1/(n-1) \sum_{k=1}^n \mathbf{x}(k)^2 = 1$ , which are the constraints in Equation (22). The Lagrangian associated with this constrained optimization problem is

$$L(\mathbf{x}, \lambda, \omega) = \mathbf{w}^T \mathbf{x} - \lambda \left( \frac{1}{n-1} \sum_{k=1}^n \mathbf{x}(k)^2 - 1 \right) - \omega \sum_{k=1}^n \mathbf{x}(k)$$

The maximizing  $\mathbf{x}^*$  is  $\mathbf{x}^* = (n-1)/(2\lambda)(\mathbf{w} - \omega)$ . Taking  $\omega$  such that the constraint  $\sum_{k=1}^n \mathbf{x}^*(k) = 0$  is fulfilled gives

$$\mathbf{x}^* = \frac{n-1}{2\lambda} (\mathbf{w} - \langle \mathbf{w} \rangle).$$

Hence, the optimal image  $\mathbf{i}^*$  is proportional to  $\mathbf{V}^-(\mathbf{w} - \langle \mathbf{w} \rangle)$ .

If we had a norm constraint on  $\mathbf{i}$  instead of the constraints in Eq (22), the Lagrangian would be

$$\tilde{L}(\mathbf{x}, \lambda) = \mathbf{w}^T \mathbf{x} - \lambda \left( \sum_{k=1}^n \mathbf{x}(k)^2 d_k - 1 \right)$$

where we have used that  $\mathbf{i}^T \mathbf{i} = \mathbf{x}^T \mathbf{V}^{-T} \mathbf{V}^- \mathbf{x} = \mathbf{x}^T \mathbf{D} \mathbf{x}$ . The  $n \times n$  matrix  $\mathbf{D}$  is diagonal with the eigenvalue  $d_k$  of the covariance matrix of the natural image patches as  $k$ -th element. The optimal  $\mathbf{x}$  would thus be  $\tilde{\mathbf{x}}^* = 1/(2\lambda)\mathbf{D}^{-1}\mathbf{w}$  so that the optimal image  $\tilde{\mathbf{i}}^*$  would be proportional to  $\mathbf{V}^-\mathbf{D}^{-1}\mathbf{w} = \mathbf{E}\mathbf{D}^{-1/2}\mathbf{w} = \mathbf{V}^T\mathbf{w}$ , for which we have used the notation  $\tilde{\mathbf{w}}$  in Section 5.3. Since the eigenvalues  $d_k$  fall off with the spatial frequency  $f$  (like  $1/f^2$ , see for example Hyvärinen et al., 2009, Chapter 5.6) the norm constraint on  $\mathbf{i}$  punishes low frequencies more heavily than the constraints in Equation (22). As a consequence, the  $\tilde{\mathbf{w}}$ , which are shown in Figure 11(a), are tuned to high frequencies while the optimal stimuli  $\mathbf{i}^*$ , shown in Figure 11(b), contain more low frequency components.

## Appendix C. Further Simulation Results

The following sections contain additional simulation results related to Section 4 and 5.

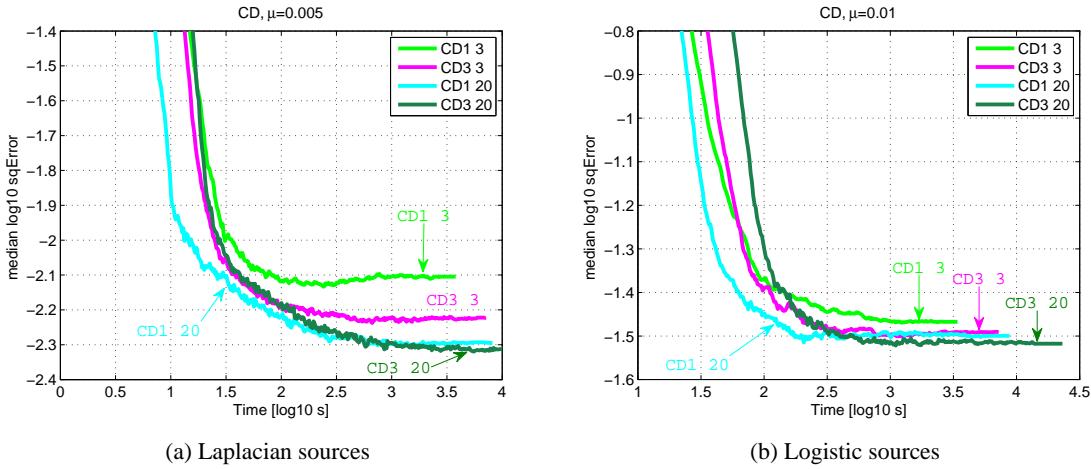


Figure 18: Trade-off between statistical and computational performance for CD. While the algorithms were running, measurements of the estimation error at a given time were made. The time variable indicates thus the time since the algorithm was started. Note the difference to Figure 6 where the time indicates the time-till-convergence. The plots show the median performance over the 100 estimation problems.  $CDx\,y$  refers to CD with  $x$  Monte Carlo steps, each using  $y$  leapfrog steps.

### C.1 Trade-Off, Section 4: Comparison of the Different Settings of Contrastive and Persistent Contrastive Divergence

We compare here the different settings of contrastive and persistent contrastive divergence. Since the two estimations methods do not have an objective function, and given the randomness that is introduced by the minibatches, choosing a reliable stopping criterion is difficult. Hence, we did not impose any stopping criterion but the maximal number of iterations. The algorithms had always converged before this maximal number of iterations was reached, in the sense that the estimation error did not visibly decrease any more. In real applications, where the true parameters are not known, assessing convergence based on the estimation error is, however, clearly not possible.

#### C.1.1 RESULTS

Figure 18 shows that for CD, using 20 leapfrog steps gives better results than using only three leapfrog steps. A trade-off between computation time and accuracy is visible: running the Markov chains for three Markov steps (CD3 20, in dark green) yields more accurate estimates than running them for one Markov step (CD1 20, in cyan) but the computations take also longer.

Figure 19 shows that for the tested schemes of PCD, using one Markov step together with 40 leapfrog steps (PCD1 40, in cyan) is the preferred choice for Laplacian sources; for logistic sources, it is PCD1 20 (shown in light green).

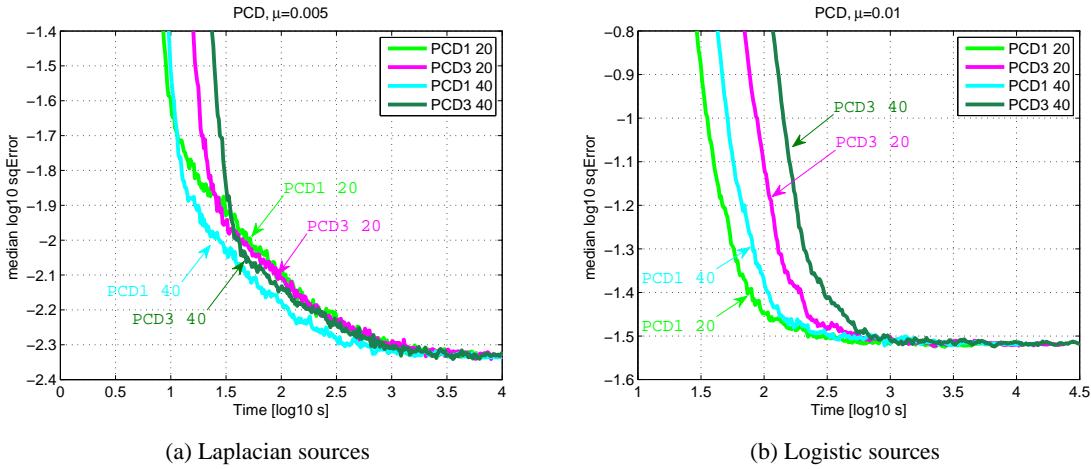


Figure 19: Trade-off between statistical and computational performance for PCD. The results are plotted in the same way as for CD in Figure 18.

## C.2 Natural Images, Section 5: Reducing Computation Time in the Optimization

The objective function  $J_T$  in Equation (8) is defined through an sample average. In an iterative optimization scheme, not all the data may be used to compute the average. The reason for using a smaller subset of the data can lie in memory considerations or in the desire to speed up the computations. We analyze here what statistical cost (reduction of estimation accuracy) such a optimization scheme implies. Furthermore, we show that optimizing  $J_T$  for increasingly larger values of  $\nu$  reduces computation time without affecting estimation accuracy. The presented results were obtained by using the the nonlinear conjugate gradient algorithm of Rasmussen (2006) for the optimization.

As working example, we consider the unnormalized Gaussian distribution of Section 3.1 for  $n = 40$ . Estimating the precision matrix and the normalizing parameter means estimating 821 parameters. We use  $T_d = 50000$ , and  $\nu = 10$ . We assume further that, for whatever reason, it is not feasible to work with all the data points at the same time but only with  $\tilde{T}_d = 25000$  samples (although for the present example, it is of course possible to use all the data).

### C.2.1 RESULTS

The lower black curve in Figure 20(a) shows the performance for the hypothetical situation where we could use all the data. The MSE reaches the level which Corollary 4 predicts (dashed horizontal line). This is the smallest error which can be obtained with nce for  $\nu = 10$  and  $T_d = 50000$ . The upper black curve in the same figure shows the MSE when only a fixed subset with  $\tilde{T}_d = 25000$  data points is used in the optimization. This clearly leads to less precise estimates. The performance can, however, be improved by randomly choosing a new subset of size  $\tilde{T}_d$  after two updates of the parameters (red curve). The improved performance comes, however, at the cost of slowing down convergence. If the resampling of the

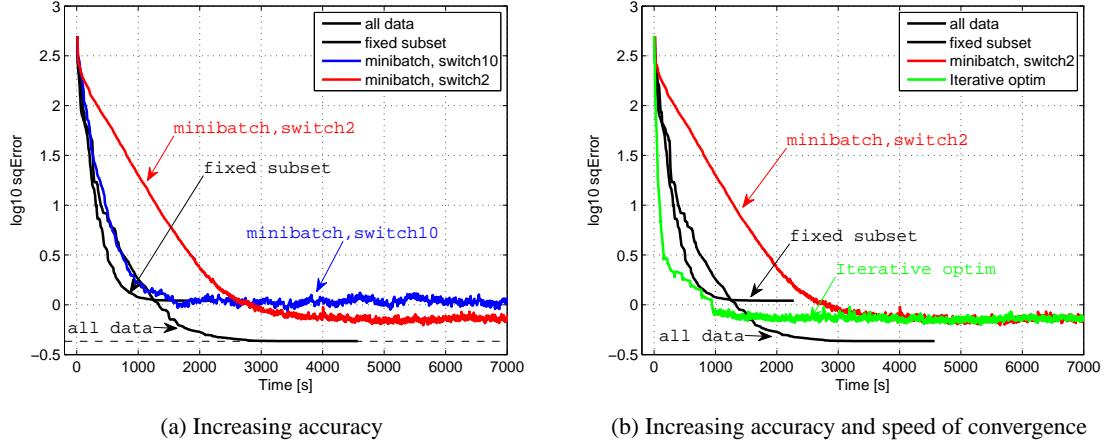


Figure 20: Analysis of the optimization strategy in Section 5. See Section C.2 for details.

subset is switched at a lower rate, for example, after 10 updates, the speed of convergence stays the same but the accuracy does not improve (blue curve).

Figure 20(b) shows the proposed optimization strategy, which we also use in Section 5 for the simulations with natural image data: We iteratively optimize  $J_T$  for increasingly larger values of  $\nu$ . Whenever we increase  $\nu$  to  $\nu + 1$ , we also take a new subset. When  $\nu$  reaches its maximal value, which is here  $\nu = 10$ , we switch the subset after two parameter updates. For the other values of  $\nu$ , we switch the subsets at a lower rate of 50 iterations. The results for this optimization strategy are shown in green (curve labelled “iterative optim”). It speeds up convergence while achieving the same precision as in the optimization with resampled subsets of size  $\tilde{T}_d$  alone (red curve in Figures (a) and (b)). By resampling new subsets, all the data are actually used in the optimization. However, the estimation accuracy is clearly worse than when all the data are used at once (as in the lower black curve). Hence, there is room for improvement in the way the optimization is performed.

### C.3 Natural Images, Section 5.4: Details for the Spline-Based One-Layer Model

The one-layer model that we consider here is

$$\ln p_m(\mathbf{x}; \boldsymbol{\theta}) = \sum_{k=1}^n f(\mathbf{w}_k^T \mathbf{x}; a_1, a_2, \dots) + c,$$

where the nonlinearity  $f$  is a cubic spline. While the two-layer models in Section 5.3 and Section 5.4 were hardcoded to assign the same value to  $\mathbf{x}$  and  $-\mathbf{x}$ , here, no symmetry assumption is made. The parameters are the feature weights  $\mathbf{w}_k \in \mathbb{R}^n$ ,  $c \in \mathbb{R}$  for the normalization of the pdf, as well as the  $a_i \in \mathbb{R}$  for the parameterization of the nonlinearity  $f$ . For the modeling of the nonlinearity, its domain needs to be defined. Its domain is related to the range of its arguments  $\mathbf{w}_k^T \mathbf{x}$ . To avoid ambiguities in the model specification, we constrain the vectors as in Equation (26). Defining  $f$  as a cubic spline on the whole real line is impossible since the number of parameters  $a_i$  would become intractable. With the constraint in Equation (26), it is enough to define  $f$  only on the interval  $[-10, 10]$  as a cubic spline. For that, we use a knot sequence with an equal spacing of 0.1. Outside

the interval, we define  $f$  to stay constant. With these specifications, we can write  $f$  in terms of B-spline basis functions with 203 coefficients  $a_1, \dots, a_{203}$ .

### C.3.1 RESULTS

The learned features are “Gabor-like” (results not shown). We observed, however, a smaller number of feature detectors that are tuned to low frequencies. Figure 16(a) in Section 5.4 shows the learned nonlinearity  $f$  (black solid curve) and the random initialization (blue dashed curve). The dashed vertical lines indicate the interval where 99% of the feature outputs occur for natural image input. The learned nonlinearity should thus only be considered valid on that interval. The nonlinearity has two striking properties: First, it is an even function. Note that no such constraint was imposed, so the symmetry of the nonlinearity is due to the symmetry in the natural images. This result validates the symmetry assumption inherent in the two-layer models. It also updates a previous result of ours where we have searched for  $f$  in a more restrictive space of functions and no symmetric nonlinearity emerged (Gutmann and Hyvärinen, 2009). Second,  $f$  is not monotonic. The shape of  $f$  is closely related to the sparsity of the feature outputs  $\mathbf{w}_k^T \mathbf{x}$ . Since the absolute values of the feature outputs are often very large or very small in natural images,  $f$  tends to map natural images to larger numbers than the noise input. This means that the model assigns more often a higher probability density to natural images than to the noise.

## C.4 Natural Images, Section 5.5: Refinement of the Thresholding Model

We are taking here a simple approach to the estimation of a two-layer model with spline nonlinearity  $f$ : We leave the feature extraction layers that were obtained for the thresholding model in Section 5.3 fixed, and learn only the cubic spline  $f$ . The model is thus

$$\ln p_m(\mathbf{x}; \boldsymbol{\theta}) = \sum_{k=1}^n f(y_k; a_1, a_2, \dots) + c, \quad y_k = \sum_{i=1}^n Q_{ki} (\mathbf{w}_i^T \mathbf{x})^2,$$

where the vector  $\boldsymbol{\theta}$  contains the parameters  $a_i$  for  $f$  and the normalizing parameter  $c$ . The knots of the spline are set to have an equal spacing of 0.1 on the interval [0 20]. Outside that interval, we define  $f$  to stay constant. With that specification, we can write  $f$  in terms of 203 B-spline basis functions. The parameter vector  $\boldsymbol{\theta} \in \mathbb{R}^{204}$  contains then the 203 coefficients for the basis functions and the parameter  $c$ .

### C.4.1 RESULTS

Figure 21(a) shows the learned nonlinearity (black solid curve) and its random initialization (blue dashed curve). The dashed vertical line around  $y = 4$  indicates the border of validity of the nonlinearity since 99% of the  $y_k$  fall, for natural image input, to the left of the dashed line. The salient property of the emerging nonlinearity is the “dip” after zero which makes  $f$  non-monotonic, as the nonlinearity which emerged in Section 5.4. Figure 21(b) shows the effective nonlinearities  $f_k$  when the different scales of the second layer outputs  $y_k$  and the normalizing parameter  $c$  are taken into account, as we have done in Figure 14(a). We calculated the scale  $\sigma_k$  by taking the average value of  $y_k$  over the natural images. The different scales  $\sigma_k$  then define different nonlinearities. Incorporating the normalizing parameter  $c$  into the nonlinearity, we obtain the set of effective nonlinearities  $f_k(y)$ ,

$$f_k(y) = f(\sigma_k y) + c/n, \quad k = 1, \dots, n. \quad (36)$$

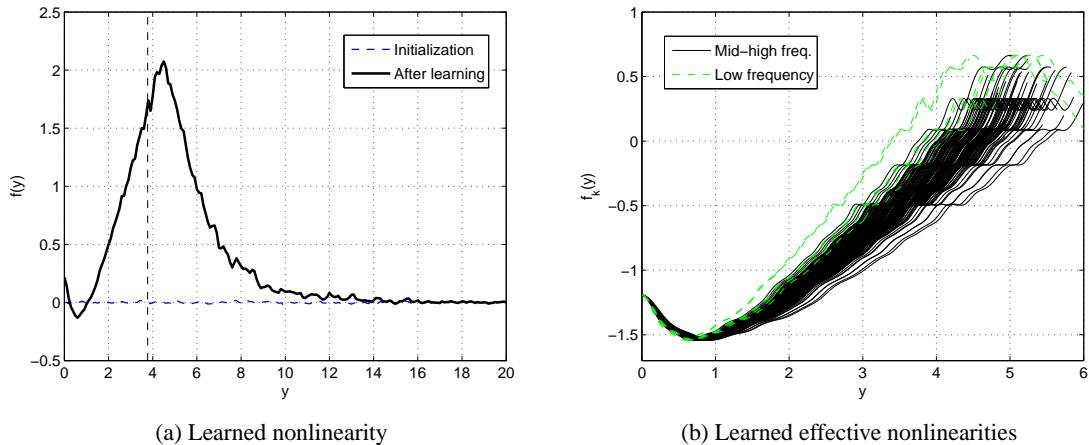


Figure 21: Refinement of the thresholding model of Section 5.3. Only the nonlinearity was learned, the features were kept fixed. The features are shown in Figures 11 to 13. (a) Learned spline (black solid curve) and the initialization (blue dashed curve). The dashed vertical line indicates the border of validity of the learned nonlinearity since 99% of the  $y_k$  fall, for natural image input, to the left of it. (b) The different scales of the  $y_k$  give rise to a set of effective nonlinearities  $f_k$ , as defined in Equation (36). Nonlinearities acting on low-frequency feature detectors are shown in green (dashed lines), the others in black (solid lines), as in Figure 14(a).

For the nonlinearities  $f_k$ , the dip occurs between zero and two. Inspection of Figure 14(b) shows that the optimal nonlinearities  $f_k$  take, unlike the thresholding nonlinearities, the distribution of the second-layer outputs  $y_k$  fully into account. The region where the dip occurs is just the region where noise input is more likely than natural image input. This means that the model is assigning more often a higher probability density to natural images than to the noise.

## C.5 Natural Images, Section 5.5: Samples from the Different Models

In Figure 17, we compared images which are considered likely by the different models. In Figure 22, we show samples that we drew from the models using Markov chains (HMC). Since the models are defined on a sphere, we constrained the Hamiltonian dynamics by projecting the states after each leapfrog step back onto the sphere. The number of leapfrog steps was set to 100, and the rejection rate to 0.35 (Neal, 2010, Section 4.4, p.30). The top row shows the most likely samples while the bottom row show the least likely ones. The least likely samples appear similar for all models. For the more probable ones, however, the two-layer models lead to more structured samples than the one-layer models.

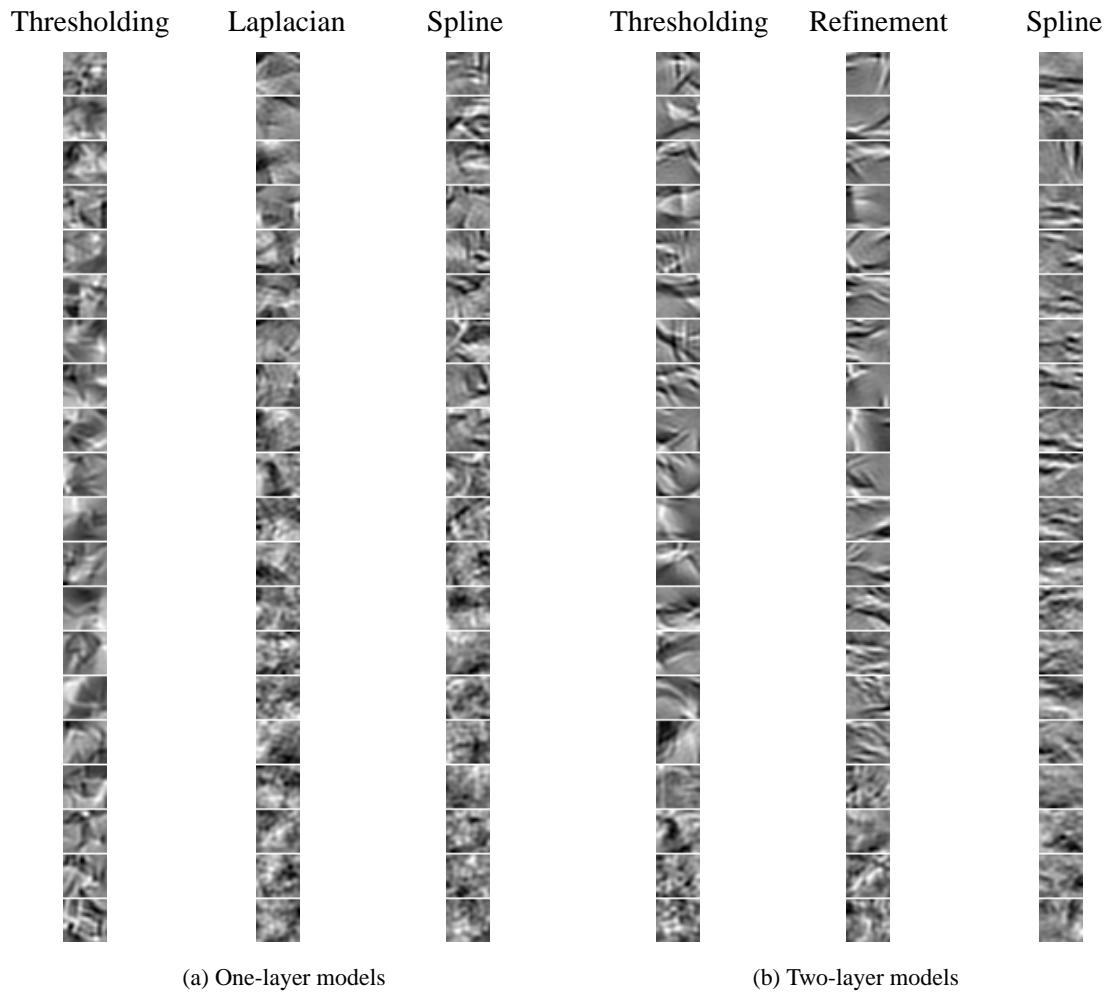


Figure 22: Sampling from the learned models of natural images. Figure (a) shows samples from the one-layer models, Figure (b) shows samples from the two-layer models. The samples are sorted so that the top ones are the most likely ones while those at the bottom are the least probable ones. See caption of Table 1 in Section 5.5 for information on the models used. Samples of the training data and the noise are shown in Figure 9 in Section 5.1.

## References

- C. Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, 1995.
- C.J. Geyer. On the convergence of Monte Carlo maximum likelihood calculations. *Journal of the Royal Statistical Society, Series B (Methodological)*, 56(1):261–274, 1994.
- M. Gutmann and A. Hyvärinen. Learning features by contrasting natural images with noise. In *Proceedings of the 19th International Conference on Artificial Neural Networks (ICANN)*, volume 5769 of *Lecture Notes in Computer Science*, pages 623–632. Springer Berlin / Heidelberg, 2009.
- M. Gutmann and A. Hyvärinen. Noise-contrastive estimation: A new estimation principle for un-normalized statistical models. In *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 9 of *JMLR W&CP*, pages 297–304, 2010.
- T. Hastie, R. Tibshirani, and J.H. Friedman. *The Elements of Statistical Learning*. Springer, 2009.
- G. Hinton. Training products of experts by minimizing contrastive divergence. *Neural Computation*, 14(8):1771–1800, 2002.
- A. Hyvärinen. Estimation of non-normalized statistical models using score matching. *Journal of Machine Learning Research*, 6:695–709, 2005.
- A. Hyvärinen. Optimal approximation of signal priors. *Neural Computation*, 20:3087–3110, 2008.
- A. Hyvärinen, P.O. Hoyer, and M. Inki. Topographic independent component analysis. *Neural Computation*, 13(7):1527–1558, 2001a.
- A. Hyvärinen, J. Karhunen, and E. Oja. *Independent Component Analysis*. Wiley-Interscience, 2001b.
- A. Hyvärinen, J. Hurri, and P.O. Hoyer. *Natural Image Statistics*. Springer, 2009.
- Y. Karklin and M. Lewicki. A hierarchical Bayesian model for learning nonlinear statistical regularities in nonstationary natural signals. *Neural Computation*, 17:397–423, 2005.
- D. Koller and N. Friedman. *Probabilistic Graphical Models*. MIT Press, 2009.
- U. Köster and A. Hyvärinen. A two-layer model of natural stimuli estimated with score matching. *Neural Computation*, 22(9):2308–2333, 2010.
- J. Lücke and M. Sahani. Maximal causes for non-linear component extraction. *Journal of Machine Learning Research*, 9:1227–1267, 2008.
- R.M. Neal. *Handbook of Markov Chain Monte Carlo*, chapter MCMC using Hamiltonian Dynamics. Chapman & Hall /CRC Press, 2010.
- B.A. Olshausen and D.J. Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583):607–609, 1996.

- S. Osindero and G. Hinton. Modeling image patches with a directed hierarchy of Markov random fields. In *Advances in Neural Information Processing Systems 20*, pages 1121–1128. MIT Press, 2008.
- S. Osindero, M. Welling, and G. E. Hinton. Topographic product models applied to natural scene statistics. *Neural Computation*, 18 (2):381–414, 2006.
- M. Pihlaja, M. Gutmann, and A. Hyvärinen. A family of computationally efficient and simple estimators for unnormalized statistical models. In *Proceedings of the 26th Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 442–449. AUAI Press, 2010.
- M.A. Ranzato and G. Hinton. Modeling pixel means and covariances using factorized third-order Boltzmann machines. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2551–2558, 2010.
- C.E. Rasmussen. Conjugate gradient algorithm, Matlab code version 2006-09-08. Downloaded from <http://learning.eng.cam.ac.uk/carl/code/minimize/minimize.m>. 2006.
- C.P. Robert and G. Casella. *Monte Carlo Statistical Methods*. Springer, 2nd edition, 2004.
- N.N. Schraudolph and T. Graepel. Towards stochastic conjugate gradient methods. In *Proceedings of the 9th International Conference on Neural Information Processing (ICONIP)*, volume 2, pages 853–856, 2002.
- W. Sun and Y. Yuan. *Optimization Theory and Methods: Nonlinear Programming*. Springer, 2006.
- Y. Teh, M. Welling, S. Osindero, and G. Hinton. Energy-based models for sparse overcomplete representations. *Journal of Machine Learning Research*, 4:1235–1260, 2004.
- T. Tieleman. Training restricted Boltzmann machines using approximations to the likelihood gradient. In *Proceedings of the 25th International Conference on Machine Learning*, pages 1064–1071, 2008.
- J. H. van Hateren and A. van der Schaaf. Independent component filters of natural images compared with simple cells in primary visual cortex. *Proceedings of the Royal Society of London. Series B: Biological Sciences*, 265(1394):359–366, 1998.
- Z. Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.
- L. Wasserman. *All of Statistics*. Springer, 2004.
- L. Younes. Parametric inference for imperfectly observed Gibbsian fields. *Probability Theory and Related Fields*, 82(4):625–645, 1989.