

Contrastive Divergence in Gaussian Diffusions

Javier R. Movellan

movellan@mplab.ucsd.edu

*Institute for Neural Computation, University of California San Diego, La Jolla,
CA 92093-0515, U.S.A.*

This letter presents an analysis of the contrastive divergence (CD) learning algorithm when applied to continuous-time linear stochastic neural networks. For this case, powerful techniques exist that allow a detailed analysis of the behavior of CD. The analysis shows that CD converges to maximum likelihood solutions only when the network structure is such that it can match the first moments of the desired distribution. Otherwise, CD can converge to solutions arbitrarily different from the log-likelihood solutions, or they can even diverge. This result suggests the need to improve our theoretical understanding of the conditions under which CD is expected to be well behaved and the conditions under which it may fail. In addition the results point to practical ideas on how to improve the performance of CD.

1 Introduction ---

CD is a recent learning rule found to work well in practice despite still unclear theoretical underpinnings (Hinton, 2002; Hinton & Salakhutdinov, 2006; Hyvärinen, 2006; MacKay, 2001; Carreira-Perpinan & Hinton, 2005; Roth & Black, 2005; Williams & Agakov, 2002; Yuille, 2004). This letter presents an analysis of CD in Gaussian diffusions—a linear, continuous-time, continuous-state version of RNNs. These networks are of interest for two reasons: (1) powerful analytical tools exist that allow comparing the behavior of CD to other algorithms, like MLE, and (2) many nonlinear systems of interest for which CD has proven useful have multiple attractors about which the systems behave locally like gaussian diffusions. Thus, the analysis of the gaussian diffusion case may provide clues for a better understanding of CD in more general conditions. The analysis presented here shows that convergence of CD is guaranteed if the first moment of the Gaussian diffusion is at equilibrium. In this case, CD and MLE converge to the same solution; otherwise, CD may converge to arbitrarily different solutions from MLE or diverge altogether.

In this letter, we pursue a continuous-time formulation of CD that makes possible the use of stochastic calculus tools. The continuous-time case can be seen as the limit of the dynamics induced by the uncorrected discrete time Langevin MCMC method (Neal, 1996). In addition, it should be noted that CD is typically interpreted as a method for learning equilibrium distributions while here we also examine it as a method for learning finite time distributions.

Consider a stochastic process $X = \{X_t : t \in \mathcal{R}_+\}$ defined by the SDE,

$$dX_t = \theta(\gamma - X_t)dt + \sqrt{2\tau}dB_t, \quad (1.1)$$

$$X_0 \sim \mathcal{N}(\mu_0, \sigma_0), \quad (1.2)$$

Here we interpret the process as a neural network, where θ is a symmetric p. d. matrix of synaptic connections, γ is a fixed vector of synaptic biases that determine the mean of the equilibrium distribution, $\tau > 0$ is a fixed parameter that controls the degree of noise in the network, and dB_t is a Brownian motion differential. The solution to this equation is (Movellan, 2006b; Oksendal, 1992):

$$X_t = e^{-t\theta} \left(X_0 + (e^{t\theta} - I)\gamma + \sqrt{2\tau} \int_0^t e^{s\theta} dB_s \right). \quad (1.3)$$

where I is the identity matrix. Thus, $X_t \sim \mathcal{N}(\mu_t, \sigma_t)$ (Movellan, 2006b):

$$\mu_t \stackrel{\text{def}}{=} \mathbb{E}[X_t] = e^{-t\theta} \mu_0 + (I - e^{-t\theta})\gamma, \quad (1.4)$$

$$\sigma_t \stackrel{\text{def}}{=} \text{Cov}[X_t] = \tau\theta^{-1} + (\sigma_0 - \tau\theta^{-1})e^{-2t\theta}. \quad (1.5)$$

At equilibrium, the mean and covariance take the following form,

$$\mu_\infty \stackrel{\text{def}}{=} \lim_{t \rightarrow \infty} \mu_t = \gamma, \quad (1.6)$$

$$\sigma_\infty \stackrel{\text{def}}{=} \lim_{t \rightarrow \infty} \sigma_t = \tau\theta^{-1}, \quad (1.7)$$

and therefore

$$\mu_t = \mu_\infty + e^{-t\theta}(\mu_0 - \mu_\infty) = \mu_0 + (I - e^{-t\theta})(\mu_0 - \mu_\infty), \quad (1.8)$$

$$\sigma_t = \sigma_\infty + e^{-2t\theta}(\sigma_0 - \sigma_\infty) = \sigma_0 + (I - e^{-2t\theta})(\sigma_0 - \sigma_\infty). \quad (1.9)$$

Express the distribution of X_t in the following Boltzmann form,

$$p(x_t) \propto e^{\phi_t(x_t)}, \quad (1.10)$$

$$\phi_t(x) \stackrel{\text{def}}{=} x'\sigma_t^{-1}\mu_t - \frac{1}{2}x'\sigma_t^{-1}x, \quad (1.11)$$

where $-\phi_t$ is the potential at time t .

2 ML and CD

The process X induces a family of distributions parameterized by t , θ , and γ . For now, we will treat the equilibrium mean γ as a fixed value and the connectivity matrix θ as an adaptive parameter. We will define learning as the process of finding values of θ under which the distribution of X_t approximates the distribution of a target rv ξ .

The method of ML calls for values of θ that maximize the likelihood function. Local maxima can be found by progressively changing θ in the direction of the log-likelihood gradient. For Boltzmann distributions, the log-likelihood gradient takes the following form (see the appendix, lemma 1),

$$\nabla_{\theta} \mathbb{E}[\log p_{X_t}(\xi)] = \mathbb{E}[\Psi_t(\xi)] - \mathbb{E}[\Psi_t(X_t)], \quad (2.1)$$

where $\Psi_t(x)$ is the unnormalized Fisher score function: $\Psi_t(x) \stackrel{\text{def}}{=} \nabla_{\theta} \phi_t(x)$.

CD was designed for situations in which the equilibrium potential $-\phi_{\infty} \stackrel{\text{def}}{=} \lim_{t \rightarrow \infty} -\phi_t$ is known but the finite time potentials are unknown. Rather than waiting for equilibrium conditions, CD operates with a finite $t > 0$ and progressively changes θ in the direction of **Hinton's CD statistic**:

$$H_t \stackrel{\text{def}}{=} \mathbb{E}[\Psi_{\infty}(\xi)] - \mathbb{E}[\Psi_{\infty}(X_t)]. \quad (2.3)$$

In gaussian diffusions, there are analytical expressions for the potentials at all times, thus allowing a direct comparison between ML and CD. It can be shown that the Fisher score function takes the following form (see the appendix, theorem 1),

$$\Psi_t(x) = x\mu'_t c_t + tx\sigma_t^{-1}e^{-t\theta}(\gamma - \mu_0)' - \frac{1}{2}xx'c_t, \quad (2.4)$$

where c_t is a p.d. matrix:

$$c_t \stackrel{\text{def}}{=} \tau \sigma_t^{-2} \left(\theta^{-2} (I - e^{-2t\theta}) - 2t \left(\theta^{-1} - \frac{1}{\tau} \sigma_0 \right) e^{-2t\theta} \right). \quad (2.5)$$

Thus, considering that $\lim_{t \rightarrow \infty} c_t = 1/\tau$, it follows that

$$\Psi_\infty(x) = \frac{1}{\tau} x \left(\gamma - \frac{1}{2} x \right)'. \quad (2.6)$$

Combining equations 2.1 and 2.4 gives us the gradient of the log likelihood function:

$$\begin{aligned} \nabla_\theta \mathbb{E}[\log p_{X_t}(\xi)] &= \frac{1}{2} (\mathbb{E}[X_t X_t'] - \mathbb{E}[\xi \xi']) c_t \\ &\quad + (\mathbb{E}[\xi] - \mu_t) \mu_t' c_t \\ &\quad + t \sigma_t^{-1} e^{-t\theta} [\mathbb{E}(\xi) - \mu_t] (\gamma - \mu_0)'. \end{aligned} \quad (2.7)$$

The gradient for the equilibrium distribution can be obtained by taking the limit as $t \rightarrow \infty$:

$$\begin{aligned} \nabla_\theta \mathbb{E}[\log p_{X_\infty}(\xi)] &= \frac{1}{2\tau} [\mathbb{E}(X_\infty X_\infty') - \mathbb{E}(\xi \xi')] \\ &\quad + \frac{1}{\tau} [\mathbb{E}(\xi) - \gamma] \gamma'. \end{aligned} \quad (2.8)$$

Combining equations 2.3 and 2.6 gives us Hinton's CD statistic:

$$H_t = \frac{1}{2\tau} [\mathbb{E}(X_t X_t') - \mathbb{E}(\xi \xi')] + \frac{1}{\tau} [\mathbb{E}(\xi) - \mu_t] \gamma'. \quad (2.9)$$

Note

$$\nabla_\theta \mathbb{E}[\log p_{X_t}(\xi)] = \tau H_t c_t + R_t, \quad (2.10)$$

where the residual term R_t is defined as follows:

$$R_t \stackrel{\text{def}}{=} t \sigma_t^{-1} e^{-t\theta} [\mathbb{E}(\xi) - \mu_t] (\gamma - \mu_0)' + (\mathbb{E}[\xi] - \mu_t) (\mu_t - \gamma)' c_t. \quad (2.11)$$

Note

$$\lim_{t \rightarrow \infty} H_t = \nabla_\theta \mathbb{E}[\log p_{X_t}(\xi)], \quad (2.12)$$

Thus, in the limit Hinton's statistic becomes the gradient of the log likelihood. Hinton (2002) derived the H_t statistic as an approximation to the gradient of the difference between two K-L divergences: $\mathbb{D}(\xi, X_\infty) - \mathbb{D}(X_t, X_\infty)$. It can be shown (see the appendix, theorem 2) that¹

$$\nabla_\theta (\mathbb{D}(\xi, X_\infty)) - \mathbb{D}(X_t, X_\infty) = -(H_t + \tilde{R}_t), \quad (2.13)$$

where the residual \tilde{R}_t is a covariance statistic:

$$\tilde{R}_t \stackrel{\text{def}}{=} \text{Cov}[\phi_t(X_t) - \phi_\infty(X_t), \Psi_t(X_t)]. \quad (2.14)$$

Hinton (2002) proposed that this residual may be ignored in practice, resulting in the CD learning rule: $\Delta\theta \propto H_t$.

We are now ready to examine four learning rules:

- ML_t : MLE for the finite time process,

$$\Delta\theta \propto \nabla_\theta \mathbb{E}[\log p_{X_t}(\xi)] = \tau H_t c_t + R_t. \quad (2.15)$$

- ML_∞ : MLE for the process at stochastic equilibrium,

$$\Delta\theta \propto \nabla_\theta \mathbb{E}[\log p_{X_\infty}(\xi)] = H_\infty. \quad (2.16)$$

- ECD : Exact CD,

$$\Delta\theta \propto \nabla_\theta (\mathbb{D}(\xi, X_\infty)) - \mathbb{D}(X_t, X_\infty) = H_t + \tilde{R}_t. \quad (2.17)$$

- CD : $\Delta\theta \propto H_t$. (2.18)

First, note that as $t \rightarrow \infty$ (i.e., if we let the network settle to equilibrium), R_t and \tilde{R}_t vanish, and the four rules converge to the same solution (see the appendix, remark 1). A more interesting question is what happens when t is finite and obviously not enough time has been given for the network to achieve stochastic equilibrium. In this case, the learning rules may converge to different solutions. In fact, when large values of $\mu_0 - \gamma$ are chosen, the residual term R_t in equation 2.11 can be made arbitrarily large to the point that CD may not converge at all or may converge to solutions arbitrarily different from ML_t and ML_∞ . However, there are cases of interest in which the learning rules converge to the same results:

- **Case 1:** $\mu_\infty = \mu_0$. It follows that $\mu_t = \mu_\infty = \gamma$. Thus, the residual term R_t in equation 2.11 vanishes, and the gradient of the log-likelihood equals Hinton's CD statistic H_t times the positive definite matrix c_t . Thus, in this case, CD and ML_t converge to the same solution.

¹This result holds for more general processes, not just gaussian diffusions.

- **Case 2:** $\mu_t = \mathbb{E}(\xi)$. In this case the first moment of the desired distribution has already been learned. Note the residual term R_t in equation 2.11 also vanishes, and thus CD and ML_t converge to the same estimate.
- **Case 3:** This case combines case 1 and case 2: $\mu_\infty = \mu_0$ and $\mu_t = \mathbb{E}(\xi)$. Under these conditions (see the appendix, remark 1),

$$H_t = H_\infty(I - e^{-2t\theta}). \quad (2.19)$$

Since $I - e^{-2t\theta}$ is a positive definite matrix and H_∞ is proportional to the gradient of ML_∞ , it follows that CD , ML_t , and ML_∞ have positive inner products with each other and converge to the same solution.

2.1 Summary of Results. The analysis reveals the importance of initializing the network so that the first moment of the states is at equilibrium. If the first moment is not at equilibrium, then CD may converge to solutions arbitrarily different from ML solutions or diverge altogether. If at equilibrium, then CD and ML_t converge to the same solution. If, in addition, $\mu_0 = \mathbb{E}[\xi]$, then CD , ML_t , and ML_∞ converge to the same solution. There currently is nothing in the theory of CD to explain why it converges when the first moment is at equilibrium but may diverge otherwise.

3 Learning the Equilibrium Means

So far we have treated the equilibrium mean, γ , as a fixed vector. This was purposely done to establish that there are conditions under which CD may not converge. In this section, we study what happens if we treat the connectivity matrix θ as a fixed parameter and the bias parameter γ as adaptive. In this case, it can be shown that

$$H_t = \frac{1}{\tau} \theta (\mathbb{E}[\xi] - \mu_t), \quad (3.1)$$

$$\nabla_\gamma \mathbb{E}[\log p_{X_t}(\xi)] = \frac{1}{\tau} \theta (I - e^{-2t\theta})(\mathbb{E}[\xi] - \mu_t) = (I - e^{-2t\theta})H_t. \quad (3.2)$$

Thus, since $I - e^{-2t\theta}$ is a positive definite matrix, when applied to the bias parameter γ , both CD and ML_t have positive inner products with each other. In addition, they converge when the first moments of the desired and obtained distributions are matched.

4 The Partially Observable Case

In many cases of interest, the state vector X_t can be divided into a vector of observable units Y_t and a vector of hidden units Z_t : $X_t = (Y_t, Z_t)$. The goal in this case is for the observable units to approximate the distribution of the target vector ξ . The expectation maximization algorithm (EM) reduces

the partially observable case to the fully observable case (Dempster, Laird, & Rubin, 1977). EM operates in an iterative manner. At iteration k , we are given a fixed parameter $\theta^{(k)}$ and a target value ξ for the observable units. The goal then becomes to learn the fully observable joint distribution of $X_t = (\xi, Z_t^{(k)})$, where $Z_t^{(k)}$ is the distribution of samples of hidden states given observable state ξ and parameter $\theta^{(k)}$. The parameter $\theta^{(k+1)}$ that optimizes this joint likelihood becomes the starting point for the next iteration. Thus, the results obtained for the fully observable case generalize to the partially observable case.

5 Conclusion

We analyzed the behavior of CD in gaussian diffusion processes. We showed that in this case, CD converges to maximum likelihood solutions if the first moment of the state distribution is at equilibrium; otherwise, CD may diverge. There is nothing in the current theory of CD that would explain the difference in behavior between these two cases. In gaussian diffusion processes, once the first-order moments of the desired distribution have been matched, the CD learning rule achieves positive inner products with the log-likelihood gradients. The nonlinear systems for which CD has proven useful have potential functions with multiple attractors, around which the systems may behave like gaussian diffusions. This may help explain why CD works well in such systems. This view of CD suggests techniques to improve its performance. For example, since the residual term R_t vanishes when the first moment of the state distribution is at equilibrium, a two-stage process could be used: on each learning trial, the system can be run using zero temperature deterministic dynamics, thus allowing it to quickly find the equilibrium mean, followed by the stochastic dynamics to estimate the H_t statistic. In addition, H_t could be estimated more efficiently using deterministic sampling methods, like the unscented transform (Julier, Uhlmann, & Durrant-Whyte, 1995).

Appendix: Derivations

A.1 Notational Conventions. The appendix assumes the processes defined in the main body of the letter. Unless otherwise stated, capital letters are used for random variables, lowercase letters for specific values taken by random variables, and Greek letters for fixed parameters. The operators \mathbb{E} and \mathbb{D} stand, respectively, for expected value and Kullback-Leibler divergence. \mathcal{R} is the set of real numbers. We leave implicit the properties of the probability space in which the random variables are defined. To simplify the notation, we identify probability functions by their arguments, and when it does not lead to confusion, we leave implicit dependencies on network

parameters—for example,

$$p_t(x) \equiv p_{X_t(\theta)}(x), \quad (\text{A.1})$$

$$X_t \equiv X_t(\theta), \quad (\text{A.2})$$

$$\phi(x) \equiv \phi(x, \theta). \quad (\text{A.3})$$

Lemma 1. *Let ξ be a target rv, θ be a random parameter, and X be a random variable with a Boltzmann distribution:*

$$p_X(u \mid \theta) \stackrel{\text{def}}{=} p(X = u \mid \theta) = \frac{1}{Z} e^{\phi(u)}, \quad (\text{A.4})$$

Then (A.5)

$$\nabla_{\theta} \mathbb{E}[\log p_X(\xi) \mid \theta] = \mathbb{E}[\Psi(\xi) \mid \theta] - \mathbb{E}[\Psi(X) \mid \theta], \quad (\text{A.6})$$

where Ψ is the unnormalized Fisher score function

$$\Psi(x, \theta) \stackrel{\text{def}}{=} \nabla_{\theta} \phi(x, \theta). \quad (\text{A.7})$$

Lemma 2. *Let $\Sigma_t : \rightarrow \mathcal{R}^n \times \mathcal{R}^n$, a matrix function of a matrix $\mathcal{R}^n \times \mathcal{R}^n$*

$$\Sigma_t(\theta) \stackrel{\text{def}}{=} \tau \left(\theta^{-1} + \left(\frac{1}{\tau} \sigma_0 - \theta^{-1} \right) e^{-2t\theta} \right), \text{ for } \theta \in \mathcal{R}^n \times \mathcal{R}^n. \quad (\text{A.9})$$

Let $\theta \in \mathcal{R}^n \times \mathcal{R}^n$ be a fixed symmetric invertible matrix and let $a \in \mathcal{R}^n \times \mathcal{R}^n$ be a fixed matrix. Let $\sigma_t \stackrel{\text{def}}{=} \Sigma_t(\theta)$, and let $\epsilon \in \mathcal{R}$. Then

$$\begin{aligned} & \frac{d}{d\epsilon} \Sigma_t^{-1}(\theta + \epsilon a) \\ &= \frac{1}{\tau} \sigma_t^{-1} \left(\theta^{-1} a \theta^{-1} (I - e^{-2t\theta}) - 2t \left(\theta^{-1} - \frac{1}{\tau} \sigma_0 \right) e^{-2t\theta} a \right) \sigma_t^{-1}. \end{aligned} \quad (\text{A.10})$$

Proof. To first order,

$$(\theta + \epsilon a)^{-1} \approx \theta^{-1} - \epsilon \theta^{-1} a \theta^{-1}, \quad (\text{A.11})$$

$$e^{-\epsilon 2ta} \approx I - \epsilon 2ta. \quad (\text{A.12})$$

Thus, to first order,

$$\begin{aligned} \Sigma_t^{-1}(\theta + \epsilon a) &\stackrel{\text{def}}{=} \frac{1}{\tau} \left((\theta + \epsilon a)^{-1} + \left(\frac{1}{\tau} \sigma_0 - (\theta + \epsilon a)^{-1} \right) e^{-2t(\theta + \epsilon a)} \right)^{-1} \\ &\approx \frac{1}{\tau} \left(\theta^{-1} - \epsilon \theta^{-1} a \theta^{-1} \right. \\ &\quad \left. + \left(\frac{1}{\tau} \sigma_0 - \theta^{-1} + \epsilon \theta^{-1} a \theta^{-1} \right) e^{-2t\theta} (I - \epsilon 2ta) \right)^{-1}. \quad (\text{A.13}) \end{aligned}$$

Separating out the constant, linear, and quadratic terms wrt ϵ ,

$$\begin{aligned} \Sigma_t^{-1}(\theta + \epsilon a) &\approx \frac{1}{\tau} \left(\theta^{-1} + \left(\frac{1}{\tau} \sigma_0 - \theta^{-1} \right) e^{-2t\theta} \right) \\ &\quad + \frac{\epsilon}{\tau} \left(-\theta^{-1} a \theta^{-1} (I - e^{-2t\theta}) - 2t \left(\frac{1}{\tau} \sigma_0 + \theta^{-1} \right) e^{-2t\theta} a \right) \\ &\quad + \frac{\epsilon^2}{\tau} 2t \theta^{-1} a \theta^{-1} e^{-2ta}. \quad (\text{A.14}) \end{aligned}$$

Using equation 1.9,

$$\frac{1}{\tau} \sigma_t = \theta^{-1} + \left(\frac{1}{\tau} \sigma_0 - \theta^{-1} \right) e^{-2t\theta}, \quad (\text{A.15})$$

and eliminating residual terms quadratic on ϵ , it follows that to first order,

$$\begin{aligned} \Sigma_t^{-1}(\theta + \epsilon a) &\approx \frac{1}{\tau} \left(\frac{1}{\tau} \sigma_t + \epsilon \left(-\theta^{-1} a \theta^{-1} (I - e^{-2t\theta}) \right. \right. \\ &\quad \left. \left. - 2t \left(\frac{1}{\tau} \sigma_0 - \theta^{-1} \right) e^{-2t\theta} a \right) \right)^{-1}. \quad (\text{A.16}) \end{aligned}$$

Using equation A.11,

$$\begin{aligned} \Sigma_t^{-1}(\theta + \epsilon a) &\approx \sigma_t^{-1} + \epsilon \tau \sigma_t^{-1} \left(\theta^{-1} a \theta^{-1} (I - e^{-2t\theta}) \right. \\ &\quad \left. - 2t \left(\theta^{-1} - \frac{1}{\tau} \sigma_0 \right) e^{-2t\theta} a \right) \sigma_t^{-1}. \quad (\text{A.17}) \end{aligned}$$

$$\begin{aligned}
\text{Thus, } \frac{d}{d\epsilon} \Sigma_t^{-1}(\theta + \epsilon a) &= \lim_{\epsilon \rightarrow 0} \frac{\Sigma(\theta + \epsilon a) - \Sigma(\theta)}{\epsilon} \\
&= \tau \sigma_t^{-1} \left(\theta^{-1} a \theta^{-1} (I - e^{-2t\theta}) - 2t \left(\theta^{-1} - \frac{1}{\tau} \sigma_0 \right) e^{-2t\theta} a \right) \sigma_t^{-1}. \quad (\text{A.18})
\end{aligned}$$

Lemma 3.

$$\nabla_{\theta} x' \Sigma_t^{-1}(\theta) x = x x' c_t, \quad (\text{A.19})$$

where c_t is a $p.d.$ matrix:

$$c_t \stackrel{\text{def}}{=} \tau \sigma_t^{-2} \left(\theta^{-2} (I - e^{-2t\theta}) - 2t \left(\theta^{-1} - \frac{1}{\tau} \sigma_0 \right) e^{-2t\theta} \right). \quad (\text{A.20})$$

Proof. Using lemma 2 and considering the symmetry of the matrices at hand,

$$\begin{aligned}
\frac{\partial}{\partial \theta_{ij}} \Sigma_t^{-1}(\theta) &= \frac{d}{d\epsilon} \Sigma_t^{-1} \left(\theta + \epsilon \frac{1_i 1_j' + 1_j 1_i'}{2} \right) \\
&= \frac{1_i 1_j' + 1_j 1_i'}{2} \tau \sigma_t^{-2} \left(\theta^{-2} (I - e^{-2t\theta}) - 2t \left(\theta^{-1} - \frac{1}{\tau} \sigma_0 \right) e^{-2t\theta} \right) \\
&= \frac{1_i 1_j' + 1_j 1_i'}{2} c_t, \quad (\text{A.21})
\end{aligned}$$

where 1_i is a vector of Kröneckner delta terms $1_i \stackrel{\text{def}}{=} (\delta_{1,i}, \dots, \delta_{n,i})'$. Thus,

$$\frac{\partial}{\partial \theta_{ij}} x' \Sigma_t^{-1}(\theta) x = x' \frac{1_i 1_j' + 1_j 1_i'}{2} c_t x, \quad (\text{A.22})$$

$$\nabla_{\theta} x' \Sigma_t^{-1}(\theta) x = x x' c_t. \quad (\text{A.23})$$

We will now show that c_t is a positive definite matrix. First, note that c_t can be expressed in the following form:

$$c_t = 2\tau t \sigma_t^{-2} \frac{1}{\tau} \sigma_0 e^{-2t\theta} + \tau \sigma_t^{-2} \theta^{-2} ((I - e^{-2t\theta}) - 2t\theta e^{-2t\theta}). \quad (\text{A.24})$$

The first term is a positive definite matrix for $t > 0$. The second term has two factors: one is a positive definite matrix $\tau \sigma_t^{-2} \theta^{-2}$ and the other positive

definite for $t = 0$ and with a positive definite derivative with respect to time:

$$\frac{d}{dt}((I - e^{-2t\theta}) - 2t\theta e^{-2t\theta}) = 4t\theta^2 e^{-2t\theta}. \quad (\text{A.25})$$

Thus, c_t is positive definite for $t \geq 0$.

Lemma 4. Let $M_t : \mathcal{R}^n \times \mathcal{R}^n \rightarrow \mathcal{R}^n$ be a vector function of a matrix

$$M_t(\theta) \stackrel{\text{def}}{=} e^{-t\theta} \mu_0 + (I - e^{-t\theta}) \gamma. \quad (\text{A.26})$$

Let $\theta \in \mathcal{R}^n \times \mathcal{R}^n$ be a fixed symmetric invertible matrix, and let $\mu_t \stackrel{\text{def}}{=} M_t(\theta)$. Then

$$\nabla_{\theta} x' \Sigma_t^{-1}(\theta) M_t(\theta) = x \mu_t' c_t + t x \sigma_t^{-1} e^{-t\theta} (\gamma - \mu_0)'. \quad (\text{A.27})$$

with Σ_t, c_t as defined in the previous lemmas.

Proof.

$$\nabla_{\theta} x' \Sigma_t^{-1}(\theta) M_t(\theta) = \nabla_{\theta} x' \Sigma_t^{-1}(\theta) \mu_t + \nabla_{\theta} x' \sigma_t^{-1} M_t(\theta). \quad (\text{A.28})$$

Using the proof for lemma 3, it is easy to see that

$$\nabla_{\theta} x' \Sigma_t^{-1}(\theta) \mu_t = x \mu_t' c_t. \quad (\text{A.29})$$

Moreover, using standard matrix calculus rules (see Movellan, 2006a),

$$\begin{aligned} \nabla_{\theta} x' \sigma_t^{-1} M_t(\theta) &= \nabla_{\theta} x' \sigma_t^{-1} (e^{-t\theta} \mu_0 + (I - e^{-t\theta}) \gamma) \\ &= t x \sigma_t^{-1} e^{-t\theta} (\gamma - \mu_0)'. \end{aligned} \quad (\text{A.30})$$

Theorem 1.

$$\nabla_{\theta} \left(x' \sigma_t^{-1} \mu_t - \frac{1}{2} x' \sigma_t^{-1} x \right) = x \mu_t' c_t + t x \sigma_t^{-1} e^{-t\theta} (\gamma - \mu_0)' - \frac{1}{2} x x' c_t. \quad (\text{A.31})$$

Proof. Direct consequence from the previous lemmas :-)

Lemma 5. Let θ be a random vector and $\{X_i : i = 1, 2\}$ be random variables such that

$$p_i(x | \theta) \stackrel{\text{def}}{=} p(X_i = x | \theta) = \frac{1}{Z_i(\theta)} e^{\phi_i(x, \theta)}, \quad \text{for } i = 1, 2, \quad (\text{A.32})$$

$$Z_i(\theta) \stackrel{\text{def}}{=} \int e^{\phi_i(x, \theta)} dx. \quad (\text{A.33})$$

Then

$$\begin{aligned}\nabla_{\theta} \mathbb{D}(X_i, X_j \mid \theta) &= \mathbb{E}[\Psi_j(X_j)] - \mathbb{E}[\Psi_j(X_i)] \\ &\quad + \text{Cov}[\phi_i(X_i) - \phi_j(X_i), \Psi_i(X_i)],\end{aligned}\tag{A.34}$$

where \mathbb{D} is the Kullback-Leibler divergence and $\Psi_i(x) \stackrel{\text{def}}{=} \nabla_{\theta} \phi_i(x, \theta)$.

Proof. To simplify the notation, we leave implicit the dependencies on θ . First note

$$\begin{aligned}\nabla_{\theta} p_i(x) \log p_j(x) &= p_i(x) \nabla_{\theta} \log p_j(x) \\ &\quad + \log p_j(x) p_i(x) \nabla_{\theta} \log p_i(x),\end{aligned}\tag{A.35}$$

and considering that,

$$\nabla_{\theta} \log p_i(x) = \Psi_i(x) - \mathbb{E}[\Psi_i(X_i)], \quad \text{for } i = 1, 2.\tag{A.36}$$

It follows that

$$\begin{aligned}\nabla_{\theta} p_i(x) \log p_j(x) &= p_i(x) (\Psi_j(x) - \mathbb{E}[\Psi_j(X_j)]) \\ &\quad + \log p_j(x) (p_i(x) (\Psi_i(x) - \mathbb{E}[\Psi_i(X_i)])).\end{aligned}\tag{A.37}$$

Thus,

$$\begin{aligned}&\nabla_{\theta} \int p_i(x) \log p_j(x) dx \\ &= \int p_i(x) \Psi_j(x) dx - \mathbb{E}[\Psi_j(X_j)] + \int p_i(x) \log p_j(x) \Psi_i(x) dx \\ &\quad - \int p_i(x) \log p_j(x) dx \mathbb{E}[\Psi_i(X_i)] \\ &= \mathbb{E}[\Psi_j(X_i)] - \mathbb{E}[\Psi_j(X_j)] \\ &\quad + \mathbb{E}[\log p_j(X_i) \Psi_i(X_i)] - \mathbb{E}[\log p_j(X_i)] \mathbb{E}[\Psi_i(X_i)] \\ &= \mathbb{E}[\Psi_j(X_i)] - \mathbb{E}[\Psi_j(X_j)] + \text{Cov}[\log p_j(X_i), \Psi_i(X_i)] \\ &= \mathbb{E}[\Psi_j(X_i)] - \mathbb{E}[\Psi_j(X_j)] + \text{Cov}[\phi_j(X_i), \Psi_i(X_i)];\end{aligned}\tag{A.38}$$

it follows that

$$\begin{aligned}\nabla_{\theta} \mathbb{D}(X_i, X_j) &= \nabla_{\theta} \int p_i(x) \log \frac{p_i(x)}{p_j(x)} dx \\ &= \mathbb{E}[\Psi_j(X_j)] - \mathbb{E}[\Psi_j(X_i)] + \text{Cov}[\phi_i(X_i), \Psi_i(X_i)] \\ &\quad - \text{Cov}[\phi_j(X_i), \Psi_i(X_i)].\end{aligned}\tag{A.39}$$

Theorem 2. Let θ be a random parameter vector and ξ a random vector independent of θ . Let $\{X_t : t \in \mathcal{R}_+\}$ be a collection of random variables with distribution

$$p_t(x_t | \theta) \propto e^{\phi_t(x, \theta)}. \quad (\text{A.40})$$

Then

$$\nabla_\theta (\mathbb{D}(\xi, X_\infty)) - \mathbb{D}(X_t, X_\infty) = -(H_t + \tilde{R}_t), \quad (\text{A.41})$$

where H_t is Hinton's CD statistic,

$$H_t \stackrel{\text{def}}{=} \mathbb{E}[\Psi_\infty(\xi)] - \mathbb{E}[\Psi_\infty(X_t)], \quad (\text{A.42})$$

and the residual \tilde{R}_t is a covariance statistic,

$$\tilde{R}_t \stackrel{\text{def}}{=} \text{Cov}[\phi_t(X_t) - \phi_\infty(X_t), \Psi_t(X_t)]. \quad (\text{A.43})$$

Proof. Let

$$\phi_\xi(u, \theta) \stackrel{\text{def}}{=} \log p(\xi = u | \theta). \quad (\text{A.44})$$

Since ξ is independent of θ , the $p(\xi | \theta)$ is constant with respect to θ . Thus,

$$\Psi_\xi(u) \stackrel{\text{def}}{=} \nabla_\theta \phi(u, \theta) = 0. \quad (\text{A.45})$$

Using lemma 5, it follows that

$$\nabla_\theta \mathbb{D}(\xi, X_\infty) = \mathbb{E}[\Psi_\infty(X_\infty)] - \mathbb{E}[\Psi_\infty(\xi)] \quad (\text{A.46})$$

$$\begin{aligned} \nabla_\theta \mathbb{D}(X_t, X_\infty) &= \mathbb{E}[\Psi_\infty(X_\infty)] - \mathbb{E}[\Psi_\infty(X_t)] \\ &\quad + \text{Cov}[\phi_t(X_t) - \phi_\infty(X_t), \Psi_t(X_t)]. \end{aligned} \quad (\text{A.47})$$

Remark 1. Analysis of the Learning Cases. First use equation 1.4 to note that if $\mu_0 = \gamma$, then $\mu_t = \mu_0$ for all t . Since under case 3, $\mu_0 = \mathbb{E}[\xi]$ and $\gamma = \mu_0$, it follows that $\mu_t = \mathbb{E}[\xi]$. Moreover, using equation 1.9 and the fact that $\mu_t = \mu_\infty$, we get that

$$\begin{aligned} \mathbb{E}[X_t X'_t] - E[\xi] E[\xi]' &= \sigma_t = \mathbb{E}[X_\infty X'_\infty] - \mathbb{E}[\xi] \mathbb{E}[\xi]' \\ &\quad + (\mathbb{E}[\xi \xi'] - \mathbb{E}[X_\infty X'_\infty]) e^{-2t\theta}. \end{aligned}$$

Thus, using equation 2.9,

$$\begin{aligned}
 H_t &= \frac{1}{2\tau} (\mathbb{E}[X_t X_t'] - \mathbb{E}[\xi \xi']) \\
 &= \frac{1}{2\tau} (\mathbb{E}[X_\infty X_\infty'] - \mathbb{E}[\xi \xi'] + (\mathbb{E}[\xi \xi'] - \mathbb{E}[X_\infty X_\infty'])e^{-2t\theta}) \\
 &= \frac{1}{2\tau} (\mathbb{E}[X_\infty X_\infty'] - \mathbb{E}[\xi \xi'])(I - e^{-2t\theta}) \\
 &= H_\infty(I - e^{-2t\theta}).
 \end{aligned} \tag{A.48}$$

Thus, H_t equals H_∞ times a positive definite matrix. To see that \tilde{R}_t vanishes as t increases, use equation A.43 to note that

$$\tilde{R}_t = \nabla_\theta(\mathbb{D}(X_t, X_\infty) - \mathbb{D}(\xi, X_\infty)) - H_t. \tag{A.49}$$

Thus

$$\lim_{t \rightarrow \infty} \tilde{R}_t = -\nabla_\theta \mathbb{D}(\xi, X_\infty) - H_t = \nabla_\theta \mathbb{E}[\log p_{X_\infty}(\xi)] - H_t = 0. \tag{A.50}$$

Acknowledgments

National Science Foundation, ECCS grant 0622229.

References

- Carreira-Perpinan, M. A., & Hinton, G. E. (2005). On contrastive divergence learning. In *10th Int. Workshop on Artificial Intelligence and Statistics (AISTATS'2005)*. Available online at <http://www.gatsby.ucl.ac.uk/aistats/Alabst.htm>.
- Dempster, A. P., Laird, N. M., & Rubin, D. B. (1997). Maximum likelihood from incomplete data via the EM algorithm. *J. Roy. Stat. Soc.*, 39, 1–38.
- Hinton, G. E. (2002). Training products of experts by minimizing contrastive divergence. *Neural Computation*, 14(8), 1771–1800.
- Hinton, G. E., & Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. *Science*, 313(5786), 504–507.
- Hyvärinen, A. (2006). *Connections between score matching, contrastive divergence, and pseudolikelihood for continuous-valued variables* (Tech. Rep.). Available online at <http://www.cs.helsinki.fi/u/ahyvarin/papers/unsuper.shtml>. Helsinki: Helsinki University.
- Julier, S. J., Uhlmann, J. K., & Durrant-Whyte, H. F. (1995). A new approach for filtering nonlinear systems. In *Proceedings of the American Control Conference* (pp. 1628–1632). Piscataway, NJ: IEEE.