

Energy Based Models Reading Group

James Allingham, Stratis Markou

1st December 2020



UNIVERSITY OF
CAMBRIDGE

Introduction

What are energy based models?

Suppose we have access to an **unnormalised** distribution \tilde{p}_θ

$$p_\theta(\mathbf{x}) = \frac{\tilde{p}_\theta(\mathbf{x})}{Z_\theta}. \quad (1)$$

We call the negative log of $\tilde{p}_\theta(\mathbf{x})$, its **energy function**

$$E_\theta(\mathbf{x}) = -\log \tilde{p}_\theta(\mathbf{x}). \quad (2)$$

Energy function equal to $-\log p_\theta(\mathbf{x})$ up to a constant independent of \mathbf{x}

$$p_\theta(\mathbf{x}) = \frac{e^{-E_\theta(\mathbf{x})}}{Z_\theta} \implies E_\theta(\mathbf{x}) = -\log p_\theta(\mathbf{x}) - \log Z_\theta, \quad (3)$$

Energy Based Models (EBMs)

We use the term Energy Based Model (EBM) for energy functions $E_\theta(\mathbf{x})$ where $Z_\theta = \int e^{-E_\theta(\mathbf{x})} d\mathbf{x}$ is not tractable.

Why Energy Based Models?

Energy-based models bring us **flexibility** in

- **Model Design:** EBMs place fewer restrictions on model design compared to other generative models (e.g. ARMs, VAEs, NFs).

We just need an $E_\theta : \mathcal{X} \rightarrow \mathbb{R}$, such that Z_θ is finite.

Not constrained to models with tractable likelihoods.

- **Problem Application:** The EBM framework is extremely general. If we can rephrase a problem as a scalar function, we can apply EBMs...

Example EBM Applications

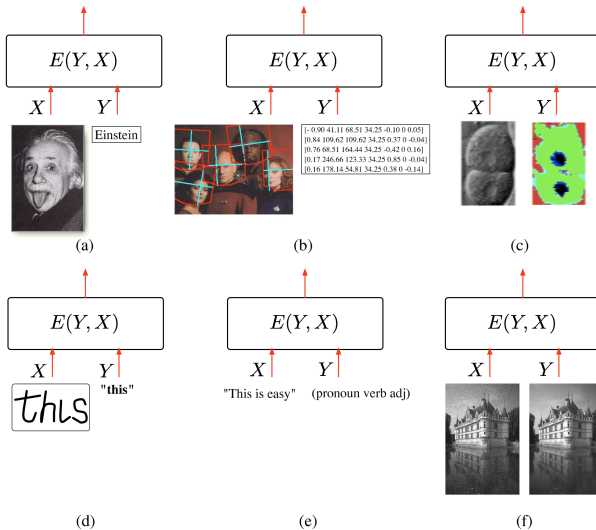


Figure 1: Examples of EBM applications [LeCun et al., 2006].

Old-school

An Example – Product of Experts [Hinton, 2002]

Product of N experts has the form

$$p_{\theta}(\mathbf{x}) = \frac{1}{Z_{\theta}} \prod_{n=1}^N p_{n,\theta}(\mathbf{x}) \iff \log p_{\theta}(\mathbf{x}) = \underbrace{\sum_{n=1}^N \log p_{n,\theta}(\mathbf{x})}_{-E_{\theta}(\mathbf{x})} - \log Z_{\theta}. \quad (4)$$

We want to train the model via maximum-likelihood

$$\theta^* = \arg \max_{\theta} \log p_{\theta}(\mathbf{x}), \quad (5)$$

but Z_{θ} is intractable for general $E_{\theta}(\mathbf{x})$. Take gradient w.r.t. θ

$$\nabla_{\theta} \log p_{\theta}(\mathbf{x}) = -E_{\theta}(\mathbf{x}) - \nabla_{\theta} \log Z_{\theta} \quad (6)$$

$$\nabla_{\theta} \log Z_{\theta} = -\nabla_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\theta}(\mathbf{x})} [\nabla_{\theta} E_{\theta}(\mathbf{x})] \quad (7)$$

Sampling $\mathbf{x} \sim p_{\theta}(\mathbf{x})$ is intractable.

Training PoEs by Contrastive Divergence

Contrastive divergence is a cancellation trick plus an approximation.

Let $p^0(\mathbf{x}) = p_D(\mathbf{x})$ be the true data distribution.

$$\theta^* = \arg \max_{\theta} \log p_{\theta}(\mathbf{x}) \iff \theta^* = \arg \min_{\theta} KL(p^0 || p_{\theta}) \quad (8)$$

Let $p_{\theta}^t(\mathbf{x})$ be the distribution of the data after t steps of MCMC

$$\mathbf{x} \sim p_{\theta}^t(\mathbf{x}) \iff \mathbf{x} \sim \text{MCMC}(\text{target} = p_{\theta}, \text{init} = \mathbf{x}_0), \mathbf{x}_0 \sim p^0(\mathbf{x}). \quad (9)$$

Note that $p_{\theta}^t \rightarrow p_{\theta}$ as $t \rightarrow \infty$, so $p_{\theta}^{\infty} = p_{\theta}$.

Idea: Run MCMC for a few iterations ($t = 1$), minimise

$$\Delta KL = KL(p^0 || p_{\theta}^{\infty}) - KL(p_{\theta}^t || p_{\theta}^{\infty}) \quad (10)$$

Motivation: Pesky term from Z_{θ} cancels.

Training PoEs by Contrastive Divergence

$$\Delta KL = \underbrace{KL(p^0 || p_\theta^\infty)}_{\textcircled{1}} - \underbrace{KL(p_\theta^t || p_\theta^\infty)}_{\textcircled{2}}. \quad (11)$$

Property 1: p_θ^t is always closer to p_θ^∞ than p^0 is to p_θ^∞ .

$$KL(p_\theta^t || p_\theta^\infty) \leq KL(p^0 || p_\theta^\infty) \implies \Delta KL \geq 0 \quad (12)$$

Property 2: If p^0 is equal to p_θ^∞ , so is p_θ^t

$$p_\theta^\infty = p^0 \implies p_\theta^t = p^0, \Delta KL = 0 \quad (13)$$

Property 3: Running $t \rightarrow \infty$ recovers maximum likelihood

$$\Delta KL \rightarrow KL(p^0 || p_\theta^\infty), \text{ as } t \rightarrow \infty. \quad (14)$$

Intuition: ΔKL encourages $\textcircled{1}$ p_θ^∞ close to p^0 and $\textcircled{2}$ p_θ^t far from p^0 . Since p_θ^t starts at p^0 , $\textcircled{2}$ encourages the chain to not wander from p^0 .

Training PoEs by Contrastive Divergence

Cancellation: Let's see how the Z_θ term cancels

$$\nabla_\theta \Delta KL = \nabla_\theta [KL(p^0 || p_\theta^\infty) - KL(p_\theta^t || p_\theta^\infty)] \quad (15)$$

$$= \int \left[\cancel{\frac{dp^0}{d\theta} \frac{\delta}{\delta p^0} KL(p^0 || p_\theta^\infty)} + \frac{dp_\theta^\infty}{d\theta} \frac{\delta}{\delta p_\theta^\infty} KL(p^0 || p_\theta^\infty) + \right. \\ \left. - \frac{dp_\theta^t}{d\theta} \frac{\delta}{\delta p_\theta^t} KL(p_\theta^t || p_\theta^\infty) - \frac{dp_\theta^\infty}{d\theta} \frac{\delta}{\delta p_\theta^\infty} KL(p_\theta^t || p_\theta^\infty) \right] dx \quad (16)$$

$$\frac{dp_\theta^\infty}{d\theta} \frac{\delta}{\delta p_\theta^\infty} KL(p^0 || p_\theta^\infty) = -p^0 \frac{d \log p_\theta^\infty}{d\theta} \quad (17)$$

$$\frac{dp_\theta^\infty}{d\theta} \frac{\delta}{\delta p_\theta^\infty} KL(p_\theta^t || p_\theta^\infty) = -p_\theta^t \frac{d \log p_\theta^\infty}{d\theta} \quad (18)$$

Now we also have that

$$\frac{d \log p_\theta^\infty}{d\theta} = -\frac{dE_\theta}{d\theta} - \frac{d \log Z_\theta}{d\theta}. \quad (19)$$

The terms involving Z_θ cancel!

Training PoEs by Contrastive Divergence

The cancellation leaves us with

$$\nabla_{\theta} \Delta KL = \int \left[p^0 \frac{dE_{\theta}}{d\theta} - p_{\theta}^t \frac{dE_{\theta}}{d\theta} - \frac{dp_{\theta}^t}{d\theta} \frac{\delta}{\delta p_{\theta}^t} KL(p_{\theta}^t || p_{\theta}^{\infty}) \right] d\mathbf{x}$$

The first term can be estimated by setting \mathbf{x} equal to the data.

The second term can be estimated with simple Monte Carlo.

The last term is still tricky. Hinton [2002] ignores it!

$$\nabla_{\theta} \Delta KL \approx \mathbb{E}_{\mathbf{x} \sim p^0(\mathbf{x})} \left[\frac{dE_{\theta}(\mathbf{x})}{d\theta} \right] - \mathbb{E}_{\mathbf{x} \sim p_{\theta}^t(\mathbf{x})} \left[\frac{dE_{\theta}(\mathbf{x})}{d\theta} \right]. \quad (20)$$

Empirically shows that this update also reduces the ignored term.

(Restricted) Boltzmann Machines

Boltzmann Machines [BM; Hinton et al., 1986], early example of EBMs.

$$E_{\theta}(\mathbf{x}, \mathbf{h}) = \sum_{i \in X, j \in H} x_i h_j w_{ij} + \sum_{i \in X} x_i \theta_i + \sum_{i \in H} h_i \theta_i + \sum_{i < j \in X} x_i x_j w_{ij} + \sum_{i < j \in H} h_i h_j w_{ij}. \quad (21)$$

where $x_i, h_i \in \{0, 1\}$. Smolensky [1986] introduced restricted BMs.

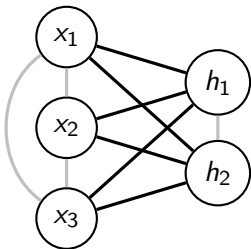


Figure 2: Graphical model of a (restricted) Boltzmann Machine.

(Restricted) Boltzmann Machines

Z_{θ} is analytic but computationally intractable, $2^{|\mathcal{X}|+|\mathcal{H}|}$ states.

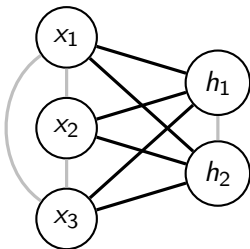


Figure 3: Graphical model of a (restricted) Boltzmann Machine.

Freund and Haussler [1994]: RBMs are PoEs \implies still intractable.

Hinton [2002] introduced Contrastive Divergence (CD) to train RBMs.

(Restricted) Boltzmann Machines

$$\nabla_{\theta} \Delta KL \approx \mathbb{E}_{\mathbf{x} \sim p^0(\mathbf{x})} \left[\frac{dE_{\theta}(\mathbf{x})}{d\theta} \right] - \mathbb{E}_{\mathbf{x} \sim p_{\theta}^t(\mathbf{x})} \left[\frac{dE_{\theta}(\mathbf{x})}{d\theta} \right]$$

$$\mathbf{x}^0 \sim p^0(\mathbf{x})$$

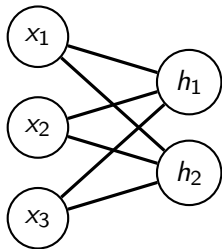


Figure 4: Gibbs sampling in an RBM for contrastive divergence.

(Restricted) Boltzmann Machines

$$\nabla_{\theta} \Delta KL \approx \mathbb{E}_{\mathbf{x} \sim p^0(\mathbf{x})} \left[\frac{dE_{\theta}(\mathbf{x})}{d\theta} \right] - \mathbb{E}_{\mathbf{x} \sim p_{\theta}^t(\mathbf{x})} \left[\frac{dE_{\theta}(\mathbf{x})}{d\theta} \right]$$

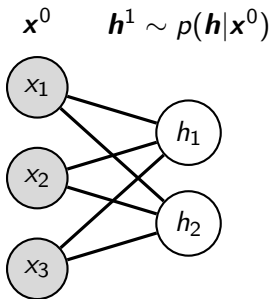


Figure 5: Gibbs sampling in an RBM for contrastive divergence.

RBM (generally PoEs): easy to sample $p(\mathbf{x}|\mathbf{h})$ and $p(\mathbf{h}|\mathbf{x})$.

(Restricted) Boltzmann Machines

$$\nabla_{\theta} \Delta KL \approx \mathbb{E}_{\mathbf{x} \sim p^0(\mathbf{x})} \left[\frac{dE_{\theta}(\mathbf{x})}{d\theta} \right] - \mathbb{E}_{\mathbf{x} \sim p_{\theta}^t(\mathbf{x})} \left[\frac{dE_{\theta}(\mathbf{x})}{d\theta} \right]$$

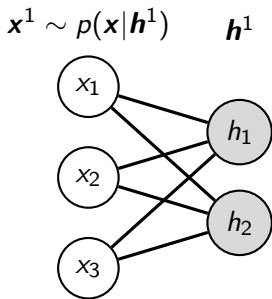


Figure 6: Gibbs sampling in an RBM for contrastive divergence.

RBM (generally PoEs): easy to sample $p(\mathbf{x}|\mathbf{h})$ and $p(\mathbf{h}|\mathbf{x})$.

Some results

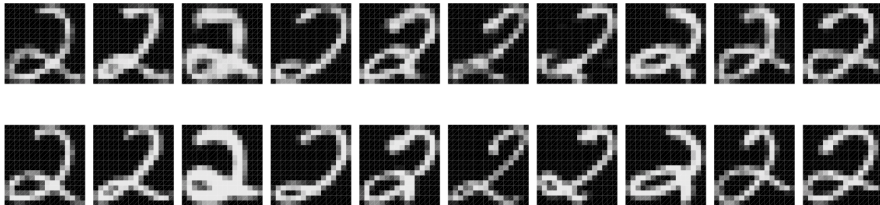


Figure 7: MNIST images (top) and their reconstructions (bottom) by an RBM.

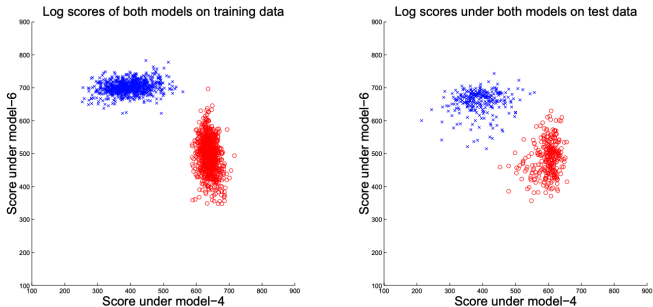


Figure 8: Energies of digits 4 and 6 under RBMs trained with digits 4 and 6.

Deep Boltzmann Machines

Deep BMs: Stack multiple RBMs [Salakhutdinov and Hinton, 2009].

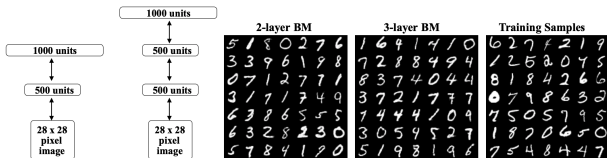


Figure 9: Training and generated data using a DBM on MNIST.

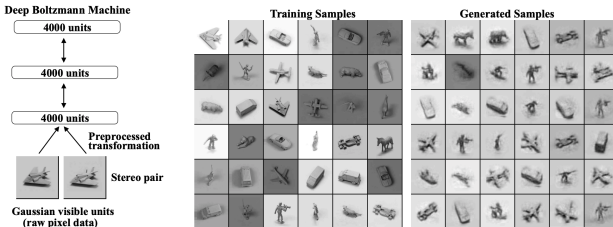


Figure 10: Training and generated data using a DBM on NORB.

Summary so far

- EBMs are a model class with intractable log-likelihoods (due to Z_θ).
- We can train EBMs by Contrastive Divergence
 - 1 Set up a Markov Chain for p_θ^t .
 - 2 Minimise $\Delta KL = KL(p^0 || p_\theta^\infty) - KL(p_\theta^t || p_\theta^\infty)$.
 - 3 Cancellation of Z_θ makes gradients tractable.
 - 4 MCMC particularly easy for PoEs (conditional independence).
- Can we leverage recent developments in Deep Learning for EBMs?
- Are there alternatives for training EBMs?

Back to the Future

Modern Training – Contrastive Divergence I

We want to maximise the likelihood

$$p_{\theta}(\mathbf{x}) = \frac{\exp(-E_{\theta}(\mathbf{x}))}{Z_{\theta}} \quad (22)$$

but we can't compute the normalizing constant

$$Z_{\theta} = \int \exp(-E_{\theta}(\mathbf{x})) \, d\mathbf{x}. \quad (23)$$

However, it turns out that with a few tricks we can compute the gradient of the log-likelihood

$$\nabla_{\theta} \log p_{\theta}(\mathbf{x}) = -\nabla_{\theta} E_{\theta}(\mathbf{x}) - \nabla_{\theta} \log Z_{\theta}. \quad (24)$$

Modern Training – Contrastive Divergence II

The first term $\nabla_{\theta} E_{\theta}(\mathbf{x})$ is easy to compute with AD. But, the second term $\nabla_{\theta} \log Z_{\theta}$ is intractable to compute exactly. However, it can be approximated, with simple MC, as

$$\nabla_{\theta} \log Z_{\theta} = \mathbb{E}_{\mathbf{x} \sim p_{\theta}(\mathbf{x})} [-\nabla_{\theta} E_{\theta}(\mathbf{x})] \approx \frac{1}{N} \sum_n^N -\nabla_{\theta} E_{\theta}(\mathbf{x}_n), \quad \mathbf{x}_n \sim p_{\theta}(\mathbf{x}). \quad (25)$$

Thus, we are taking gradient steps in the direction

$$\nabla_{\theta} E_{\theta}(\mathbf{x}_{\text{train}}) - \nabla_{\theta} E_{\theta}(\mathbf{x}_{\text{sample}}). \quad (26)$$

Modern Training – Contrastive Divergence III

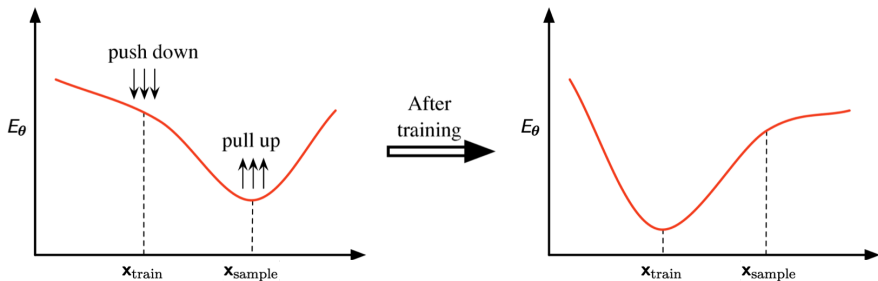


Figure 11: Taking steps in the direction of $\nabla_{\theta} E_{\theta}(\mathbf{x}_{\text{train}}) - \nabla_{\theta} E_{\theta}(\mathbf{x}_{\text{sample}})$. Adapted from [LeCun et al., 2006].

Modern Training – Contrastive Divergence IV

Alas, sampling from $p_{\theta}(\mathbf{x})$ is highly *non-trivial*. Thus, we must resort to further approximation. A common choice is to use Langevin MCMC

$$\mathbf{x}^{t+1} \leftarrow \mathbf{x}^t + \frac{\epsilon^2}{2} \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}^t) + \epsilon \mathbf{z}^t, \quad t = 0, 1, \dots, T - 1, \quad (27)$$

where $\mathbf{z}^t \sim \mathcal{N}(0, 1)$ and $\mathbf{x}^0 \sim p(\mathbf{x})$.

Note: $\nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}^t) = -\nabla_{\mathbf{x}} E_{\theta}(\mathbf{x}) - \cancel{\nabla_{\mathbf{x}} \log Z_{\theta}} = -\nabla_{\mathbf{x}} E_{\theta}(\mathbf{x})$.

When $\epsilon \rightarrow 0$ and $T \rightarrow \infty$, $\mathbf{x}^T \sim p_{\theta}(\mathbf{x})$.

However, we usually do not run the chain to convergence due to the high computational cost.

Improved CD for EBM's I

But wait! The CD gradient doesn't come from the LL, but rather

$$\Delta KL = KL(p^0 || p_\theta^\infty) - KL(p_\theta^t || p_\theta^\infty). \quad (28)$$

Taking gradients w.r.t this objective results in

$$-\nabla_\theta E_\theta(\mathbf{x}) + \mathbb{E}_{\mathbf{x} \sim p_\theta^t(\mathbf{x})} [\nabla_\theta E_\theta(\mathbf{x})] - \nabla_\theta p_\theta^t(\mathbf{x}) \nabla_{p_\theta^t(\mathbf{x})} KL(p_\theta^t || p_\theta^0) \quad (29)$$

where we have an additional **KL term** [Du et al., 2020].

CD training without the **KL term** doesn't minimize any scalar loss function [Sutskever and Tieleman, 2010].

Improved CD for EBM's II

Du et al. [2020] show that Adding the **KL term** is equivalent to adding a KL loss

$$\mathcal{L}_{KL} = \underbrace{\mathbb{E}_{p_{\theta}^t(\mathbf{x})} [\mathbf{x}]}_{\textcircled{1}} + \underbrace{\mathbb{E}_{p_{\theta}^t(\mathbf{x})} [\log p_{\theta}^t(\mathbf{x})]}_{\textcircled{2}}. \quad (30)$$

Estimation of the KL loss is fairly involved.

- ① requires differentiating through the MCMC sampling.
- ② is estimated via a nearest-neighbours approximation.

Improved CD for EBM's III

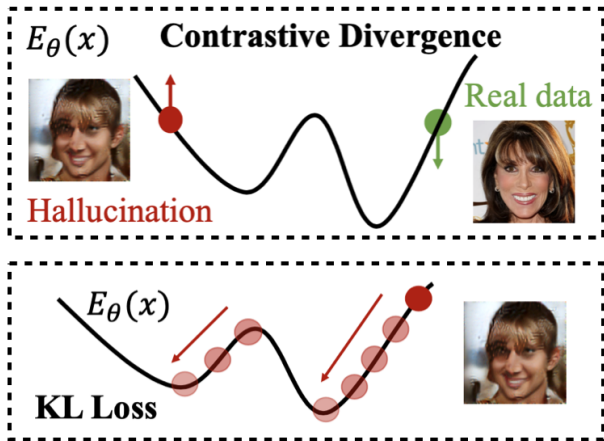


Figure 12: Intuition for the KL loss – regularization which prevents bad samples! Adapted from [Du et al., 2020].

Improved CD for EBM's IV

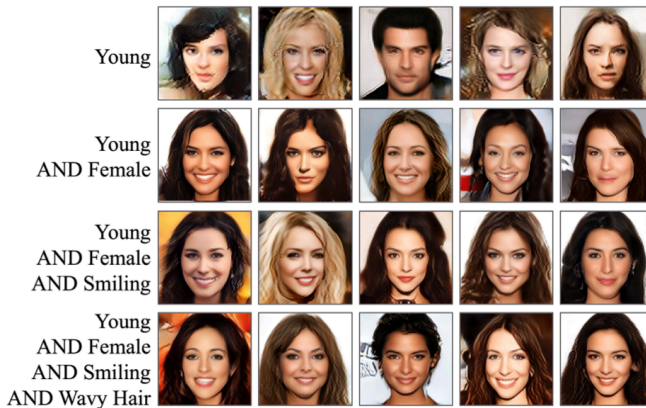


Figure 13: Samples from an EBM composed of individual conditional EBMs trained on the CelebA-HQ dataset using Improved CD. Adapted from [Du et al., 2020].

Modern Training – Score Matching

If $f(\mathbf{x})$ and $g(\mathbf{x})$ have equal first derivatives (a.k.a *score functions*), then $f(\mathbf{x}) \equiv g(\mathbf{x}) + \text{constant}$. When $f(\mathbf{x})$ and $g(\mathbf{x})$ are log PDFs, $f(\mathbf{x}) \equiv g(\mathbf{x})$. Thus, Hyvärinen and Dayan [2005] propose to learn an EBM by minimising

$$D_F(p_{\text{data}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) = \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x})\|^2 \right]. \quad (31)$$

- The expectation can be approximated with a simple MC estimator.
- The term $\nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) = -\nabla_{\mathbf{x}} \log E_{\theta}(\mathbf{x})$.
- Unfortunately, the term $\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$ is intractable.

However, using integration by parts, we can rewrite this as

$$D_F(p_{\text{data}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) = \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \left[\frac{1}{2} \sum_{i=1}^d \left(\frac{\partial E_{\theta}(\mathbf{x})}{\partial x_i} \right)^2 + \frac{\partial^2 E_{\theta}(\mathbf{x})}{(\partial x_i)^2} \right] + C. \quad (32)$$

Modern Training – Denoising Score Matching

Naive Score Matching has two potentially problematic requirements:

- 1 $p_{\text{data}}(\mathbf{x})$ is continuously differentiable and finite everywhere, and
- 2 the computation of expensive second-order gradients.

Problem 1 can be solved by adding noise to each data point $\tilde{\mathbf{x}} = \mathbf{x} + \epsilon$, resulting in a noisy data distribution $q(\tilde{\mathbf{x}}) = \int q(\tilde{\mathbf{x}} | \mathbf{x}) p_{\text{data}}(\mathbf{x}) d\mathbf{x}$.

Vincent [2011] solve problem 2 by showing:

$$\begin{aligned} D_F(q(\tilde{\mathbf{x}}) \parallel p_{\theta}(\tilde{\mathbf{x}})) &= \mathbb{E}_{q(\tilde{\mathbf{x}})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}} \log q(\tilde{\mathbf{x}}) - \nabla_{\mathbf{x}} \log p_{\theta}(\tilde{\mathbf{x}})\|_2^2 \right] \\ &= \mathbb{E}_{q(\mathbf{x}, \tilde{\mathbf{x}})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}} \log q(\tilde{\mathbf{x}}|\mathbf{x}) - \nabla_{\mathbf{x}} \log p_{\theta}(\tilde{\mathbf{x}})\|_2^2 \right] + C, \end{aligned} \tag{33}$$

thus avoiding any expensive second-order gradients.

New problems: Inconsistency. Trade-off between estimator variance and noise magnitude.

Modern Training – Sliced Score Matching

Sliced Score Matching minimizes the sliced Fisher divergence

$$D_{SF}(p_{\text{data}}||p_{\theta}) = \mathbb{E}_{p_{\text{data}}(\mathbf{x})}\mathbb{E}_{p(\mathbf{v})} \left[\frac{1}{2}(\mathbf{v}^T \nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x}) - \mathbf{v}^T \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}))^2 \right],$$

where $p(\mathbf{v})$ denotes a projection dist. such that $\mathbb{E}_{p(\mathbf{v})}[\mathbf{v}\mathbf{v}^T]$ is pos. def.

As before, we can use the chain rule to avoid $\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$:

$$\mathbb{E}_{p_{\text{data}}(\mathbf{x})}\mathbb{E}_{p(\mathbf{v})} \left[\frac{1}{2} \sum_{i=1}^d \left(\frac{\partial E_{\theta}(\mathbf{x})}{\partial x_i} v_i \right)^2 + \sum_{i=1}^d \sum_{j=1}^d \frac{\partial^2 E_{\theta}(\mathbf{x})}{\partial x_i \partial x_j} v_i v_j \right]. \quad (34)$$

However, unlike before, this form has a computational complexity of $\mathcal{O}(d)$ rather than $\mathcal{O}(d^2)$:

$$\sum_{i=1}^d \sum_{j=1}^d \frac{\partial^2 E_{\theta}(\mathbf{x})}{\partial x_i \partial x_j} v_i v_j = \sum_{i=1}^d \frac{\partial}{\partial x_i} \underbrace{\left(\sum_{j=1}^d \frac{\partial E_{\theta}(\mathbf{x})}{\partial x_j} v_j \right)}_{:=f(\mathbf{x})} v_i. \quad (35)$$

Modern Training – Noise Contrastive Estimation

Gutmann and Hyvärinen [2010] proposed an alternative training method.

Define a known and tractable reference distribution $p_r(\mathbf{x})$.

Treat Z_θ as a trainable variable.

$$\log p_\theta = \log \tilde{p}_\theta(\mathbf{x}) + C. \quad (36)$$

Draw \mathbf{x} from $p_r(\mathbf{x})$ or from $p_D(\mathbf{x})$ (denoted $y = 0, 1$ respectively).

Use EBM to set up a classifier which distinguishes $y = 0, 1$

$$p(y = 0|\mathbf{x}, \theta) = \frac{p_r(\mathbf{x})}{p_\theta(\mathbf{x}) + p_r(\mathbf{x})}, \quad p(y = 1|\mathbf{x}, \theta) = \frac{p_\theta(\mathbf{x})}{p_\theta(\mathbf{x}) + p_r(\mathbf{x})}, \quad (37)$$

where $y = 0, 1$ mean \mathbf{x} drawn from $p_r(\mathbf{x})$ and from $p_D(\mathbf{x})$ respectively.

Modern Training – Noise Contrastive Estimation

Draw $\mathbf{x}_1, \dots, \mathbf{x}_N \sim p_r(\mathbf{x})$ and $\mathbf{x}_{N+1}, \dots, \mathbf{x}_{2N} \sim p_D(\mathbf{x})$. Minimise

$$\mathcal{L}_{NCE} = \sum_{n=1}^N \log p(y_n = 1 | \mathbf{x}_n, \boldsymbol{\theta}) + \log p(y_{N+n} = 0 | \mathbf{x}_{N+n}, \boldsymbol{\theta}), \quad (38)$$

binary cross entropy loss for classifying samples.

Theorem (informal) Gutmann and Hyvärinen [2010]

If there exists θ^* such that $p_{\theta^*} = p_D$, then in the limit $N \rightarrow \infty$ we have $\theta \rightarrow \theta^*$. Further, θ^* is a unique global optimum.

Observation: Objective automatically takes care of Z_θ

$$p(y = 0 | \mathbf{x}, \boldsymbol{\theta}) = \frac{p_r(\mathbf{x})}{p_\theta(\mathbf{x}) + p_r(\mathbf{x})}, p(y = 1 | \mathbf{x}, \boldsymbol{\theta}) = \frac{p_\theta(\mathbf{x})}{p_\theta(\mathbf{x}) + p_r(\mathbf{x})}, \quad (39)$$

Intuition: Trainable variable C cannot go to neither $-\infty$ nor ∞ .

Conditional Noise Contrastive Estimation

Challenge with NCE: How to choose the reference distribution $p_r(\mathbf{x})$?

Ceylan and Gutmann [2018] propose $p_r(\mathbf{x}) = \int q(\mathbf{x}|\mathbf{x}')p_D(\mathbf{x}')d\mathbf{x}'$ and

$$D(\mathbf{x}, \mathbf{x}') = p(y = 1|\mathbf{x}, \mathbf{x}', \theta) = \frac{q(\mathbf{x}'|\mathbf{x})p_\theta(\mathbf{x})}{q(\mathbf{x}'|\mathbf{x})p_\theta(\mathbf{x}) + q(\mathbf{x}|\mathbf{x}')p_\theta(\mathbf{x}')}, \quad (40)$$

where a typical choice for q is

$$q(\mathbf{x}'|\mathbf{x}) = \mathcal{N}(\mathbf{x}'; \mathbf{x}, \sigma^2 \mathbf{I}). \quad (41)$$

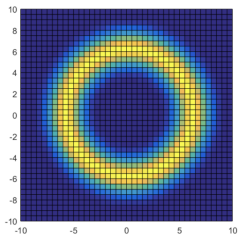
Observation: Equation (40) invariant to scaling p_θ , and Z_θ cancels.

Similarly to NCE, set up classification task and minimise

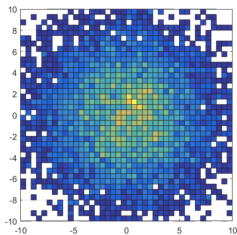
$$\mathcal{L}_{CNCE} = - \sum_{n=1}^N [\log D(\mathbf{x}_n, \mathbf{x}'_n) + \log(1 - D(\mathbf{x}'_n, \mathbf{x}_n))] \quad (42)$$

where $\mathbf{x}_1, \dots, \mathbf{x}_N \sim p_D(\mathbf{x})$ and $\mathbf{x}'_1, \dots, \mathbf{x}'_N \sim q(\mathbf{x}'|\mathbf{x})$.

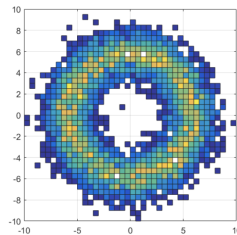
Conditional Noise Contrastive Estimation



(a) Contour plot of the data pdf



(b) NCE noise (histogram)



(c) CNCE noise (histogram)

Figure 14: Data distribution, NCE samples and CNCE samples.

Intuition: Contrastive data closer to real data, so classification is more challenging, training stays informative for longer.

By comparison, NCE classification is much easier. Classification task solved quickly, at which point the EBM stops learning.

Relation between (C)NCE and CD

Yair and Michaeli [2020] view CD as a (C)NCE classification problem.

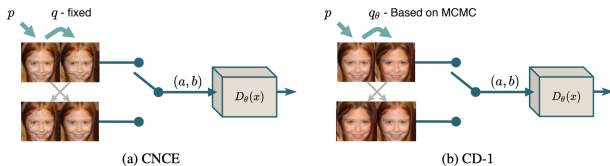


Figure 15: Illustration of CD viewed as a classification problem.

The CNCE parameter update rule can be written as

$$\Delta\theta_{CNCE} = \sum_{n=1}^N (1 - D(\mathbf{x}_n, \mathbf{x}'_n)) [\nabla_\theta \log p_\theta(\mathbf{x}_n) - \nabla_\theta \log p_\theta(\mathbf{x}'_n)]$$

Intuition:

- 1 Term $\nabla_\theta \log p_\theta(\mathbf{x}_n)$ encourages high prob. near the data.
- 2 Term $\nabla_\theta \log p_\theta(\mathbf{x}'_n)$ encourages low prob. at contrastive samples.
- 3 Term $1 - D(\mathbf{x}_n, \mathbf{x}'_n)$ downweights easily-classified pairs.
- 4 CNCE keeps $D(\mathbf{x}_n, \mathbf{x}'_n)$ close to 1/2 for longer, compared to NCE.

Relation between (C)NCE and CD

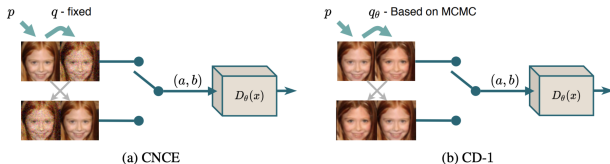


Figure 16: Illustration of CD viewed as a classification problem.

$$\Delta\theta_{CNCE} = \sum_{n=1}^N (1 - D(\mathbf{x}_n, \mathbf{x}'_n)) [\nabla_{\theta} \log p_{\theta}(\mathbf{x}_n) - \nabla_{\theta} \log p_{\theta}(\mathbf{x}'_n)]$$

where $\mathbf{x}_1, \dots, \mathbf{x}_N \sim p_D(\mathbf{x})$ and $\mathbf{x}'_1, \dots, \mathbf{x}'_N \sim p_r(\mathbf{x})$.

Relation to CD: Let $q(\mathbf{x}'|\mathbf{x}) = q_{\theta}(\mathbf{x}'|\mathbf{x})$ be a reversible Markov Chain, with stationary distribution p_{θ} .

The optimal classifier becomes a random guess

$$D(\mathbf{x}_n, \mathbf{x}'_n) = \frac{q(\mathbf{x}'|\mathbf{x})p_{\theta}(\mathbf{x})}{q(\mathbf{x}'|\mathbf{x})p_{\theta}(\mathbf{x}) + q(\mathbf{x}|\mathbf{x}')p_{\theta}(\mathbf{x}')} = \frac{1}{2}. \quad (43)$$

Relation between (C)NCE and CD

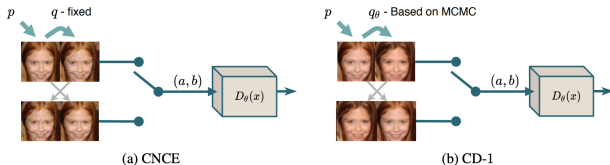


Figure 17: Illustration of CD viewed as a classification problem.

Under this model, the update rule becomes identical to CD

$$\Delta\theta_{CNCE} = \frac{1}{2} \sum_{n=1}^N [\nabla_{\theta} \log p_{\theta}(\mathbf{x}_n) - \nabla_{\theta} \log p_{\theta}(\mathbf{x}'_n)].$$

Important detail (stopping gradients)

Sample $\mathbf{x}'_n \sim q_{\theta}(\mathbf{x}'|\mathbf{x})$ depends on θ . It is necessary to stop gradients through q_{θ} for equivalence between CD and CNCE. Stopping gradients is equivalent to ignoring the intractable CD term.

Relation between (C)NCE and CD

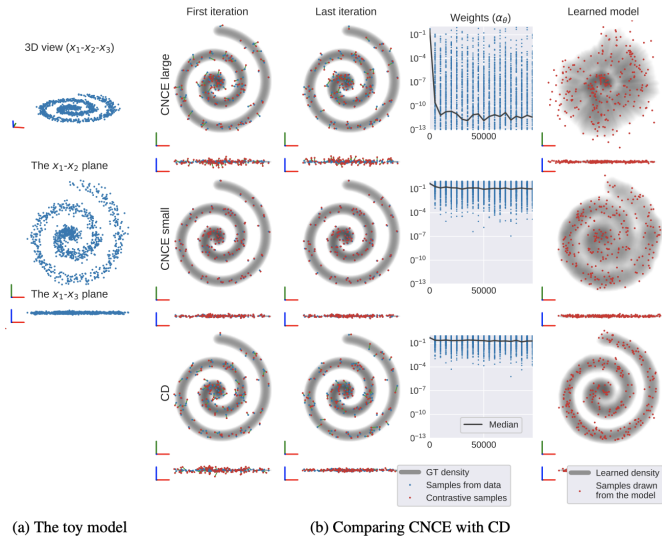


Figure 18: Illustration of CNCE and CD training.

Relation between SM and CD

Hyvarinen [2007] relates SM with CD under a particular Markov Chain.

Suppose we generate samples using Langevin Dynamics

$$\mathbf{x}' = \mathbf{x} + \frac{\eta^2}{2} \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) + \eta \epsilon, \quad (44)$$

where $\eta > 0$ is a step size and ϵ is standard Gaussian noise. Starting from the CD update rule, if we

- 1 Stop gradients flowing through p_{θ} in equation (44),
- 2 Take the limit $\eta \rightarrow 0$,

we recover the SM update rule.

Theorem (informal) Hyvarinen [2007]

Under an appropriate gradient stopping approximation, the update rule of CD is equivalent to the update rule of SM in the limit $\eta \rightarrow 0$.

EBMs in the Wild

Classifiers are EBM's? [Grathwohl et al., 2020b] I

Classifier:

$$p_{\theta}(y | \mathbf{x}) = \frac{\exp(f_{\theta}(\mathbf{x})[y])}{\sum_{y'} \exp(f_{\theta}(\mathbf{x})[y'])}. \quad (45)$$

Joint:

$$p_{\theta}(\mathbf{x}, y) = \frac{\exp(f_{\theta}(\mathbf{x})[y])}{Z_{\theta}}. \quad (46)$$

Marginal:

$$p_{\theta}(\mathbf{x}) = \sum_y p_{\theta}(\mathbf{x}, y) = \frac{\sum_y \exp(f_{\theta}(\mathbf{x})[y])}{Z_{\theta}}.$$

Energy:

$$E_{\theta}(\mathbf{x}) = -\log \sum_y \exp(f_{\theta}(\mathbf{x})[y]). \quad (47)$$

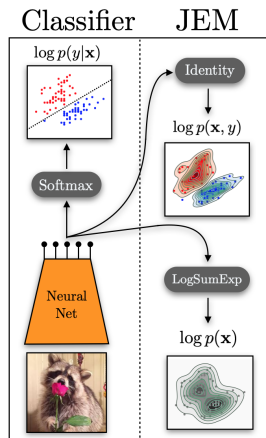


Figure 19: JEM.

Classifiers are EBM's? [Grathwohl et al., 2020b] II

Note:

$$p_{\theta}(y | \mathbf{x}) = \frac{p_{\theta}(\mathbf{x}, y)}{p_{\theta}(\mathbf{x})} = \frac{\exp(f_{\theta}(\mathbf{x})[y])}{Z_{\theta}} \frac{Z_{\theta}}{\sum_{y'} \exp(f_{\theta}(\mathbf{x})[y'])}, \quad (48)$$

as before.

For optimization, factorise the joint:

$$\log p_{\theta}(\mathbf{x}, y) = \log p_{\theta}(\mathbf{x}) + \log p_{\theta}(y|\mathbf{x}). \quad (49)$$

$\log p_{\theta}(\mathbf{x})$ is optimised using persistent CD.

$\log p_{\theta}(y|\mathbf{x})$ is optimized using the standard CE loss.

JEM is a good classifier *and* generative model. The classifier is well calibrated, and the generative model knows what it doesn't know.

VAE-EBM Hybrids [Xiao et al., 2020] I

Product of a VAE and an EBM: $h_{\psi,\theta}(\mathbf{x}, \mathbf{z}) = \frac{1}{Z_{\psi,\theta}} p_{\theta}(\mathbf{x}, \mathbf{z}) e^{-E_{\psi}(\mathbf{x})}$.

Marginalizing out \mathbf{z} gives

$$h_{\psi,\theta}(\mathbf{x}) = \frac{1}{Z_{\psi,\theta}} \int p_{\theta}(\mathbf{x}, \mathbf{z}) e^{-E_{\psi}(\mathbf{x})} d\mathbf{z} = \frac{1}{Z_{\psi,\theta}} p_{\theta}(\mathbf{x}) e^{-E_{\psi}(\mathbf{x})}. \quad (50)$$

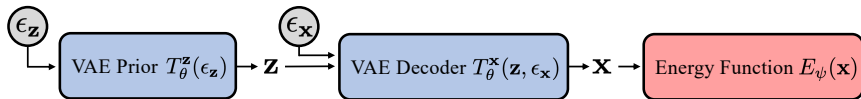


Figure 20: VAE-EBM Computational Model

VAE-EBM Hybrids [Xiao et al., 2020] II

ψ, θ are trained by maximizing the marginal log-likelihood, in 2 steps:

$$\begin{aligned} \log h_{\psi, \theta}(\mathbf{x}) &= \log p_{\theta}(\mathbf{x}) - E_{\psi}(\mathbf{x}) - \log Z_{\psi, \theta} & (51) \\ &\geq \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - KL(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))}_{\mathcal{L}_{\text{vae}}(\mathbf{x}, \theta, \phi)} - \underbrace{E_{\psi}(\mathbf{x}) - \log Z_{\psi, \theta}}_{\mathcal{L}_{\text{EBM}}(\mathbf{x}, \psi, \theta)}. \end{aligned}$$

Step 1: Train the VAE via \mathcal{L}_{vae} .

Step 2: Fix the VAE, train the EBM via \mathcal{L}_{EBM} with CD.



Figure 21: CelebA HQ 256 – Qualitative Results

Energy-based OOD Detection [Liu et al., 2020]

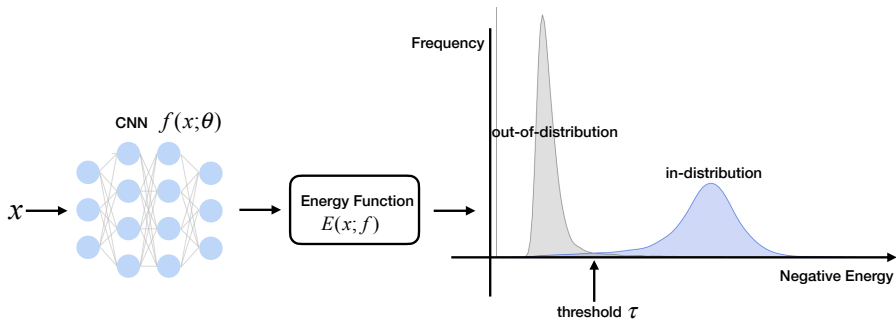


Figure 22: Energy-based OOD Detection.

Either apply the energy score to pre-trained NNs, or use it as an additional loss during training.

Outperforms JEM for OOD detection.

Conclusions

EBMs are flexible class of models, which expand our modelling toolbox.

Several approaches for training EBMs, each with pros and cons:

- 1 Contrastive Divergence (CD).
- 2 Score Matching (SM), denoising SM, sliced SM.
- 3 Noise Contrastive Estimation (NCE) and Conditional NCE (CNCE).

Some of these are equivalent under certain conditions.

Training Method	Fast	Stable training	High dimensions	No aux. model	Unrestricted architecture	Approximates likelihood
Markov chain Monte Carlo	✗	✗	✓	✓	✓	✓
Score Matching Approaches	✓	✗	✓	✓	✗	✗
Noise Contrastive Approaches	✓	✓	✗	✗	✓	✗

(Adapted from [Grathwohl et al., 2020a].)

Conclusions

Energy Based Models are great!

Many things that we didn't get to talk about...

Other training methods:

- Minimizing KL differences – generalized CD and SM.
- Minimizing Stein's discrepancy. E.g. LSD [Grathwohl et al., 2020c].
- Adversarial training

Other cool EBM papers:

- Your GAN is also a secret EBM [Che et al., 2020].
- Generalised Energy Based Models (GAN-EBM hybrids) [Arbel et al., 2020].

And of course, many more details on CD, DM, NCE, their variants, their connections, and more in [Song and Kingma, 2021].

References I

- M. Arbel, L. Zhou, and A. Gretton. Generalized energy based models. *arXiv preprint arXiv:2003.05033*, 2020.
- C. Ceylan and M. U. Gutmann. Conditional noise-contrastive estimation of unnormalised models. In *International Conference on Machine Learning*, pages 726–734. PMLR, 2018.
- T. Che, R. Zhang, J. Sohl-Dickstein, H. Larochelle, L. Paull, Y. Cao, and Y. Bengio. Your gan is secretly an energy-based model and you should use discriminator driven latent sampling. *arXiv preprint arXiv:2003.06060*, 2020.
- Y. Du, S. Li, J. Tenenbaum, and I. Mordatch. Improved contrastive divergence training of energy based models. *arXiv preprint arXiv:2012.01316*, 2020.
- Y. Freund and D. Haussler. Unsupervised learning of distributions of binary vectors using two layer networks. 1994.

References II

- W. Grathwohl, J. Kelly, M. Hashemi, M. Norouzi, K. Swersky, and D. Duvenaud. No mcmc for me: Amortized sampling for fast and stable training of energy-based models. *arXiv preprint arXiv:2010.04230*, 2020a.
- W. Grathwohl, K.-C. Wang, J.-H. Jacobsen, D. Duvenaud, M. Norouzi, and K. Swersky. Your classifier is secretly an energy based model and you should treat it like one, 2020b.
- W. Grathwohl, K.-C. Wang, J.-H. Jacobsen, D. Duvenaud, and R. Zemel. Learning the stein discrepancy for training and evaluating energy-based models without sampling. In *International Conference on Machine Learning*, pages 3732–3747. PMLR, 2020c.

References III

- M. Gutmann and A. Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 297–304. JMLR Workshop and Conference Proceedings, 2010.
- G. E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800, 2002.
- G. E. Hinton, T. J. Sejnowski, et al. Learning and relearning in boltzmann machines. *Parallel distributed processing: Explorations in the microstructure of cognition*, 1(282-317):2, 1986.
- A. Hyvarinen. Connections between score matching, contrastive divergence, and pseudolikelihood for continuous-valued variables. *IEEE Transactions on neural networks*, 18(5):1529–1531, 2007.

References IV

- A. Hyvärinen and P. Dayan. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6 (4), 2005.
- Y. LeCun, S. Chopra, R. Hadsell, M. Ranzato, and F. Huang. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006.
- W. Liu, X. Wang, J. D. Owens, and Y. Li. Energy-based out-of-distribution detection. *arXiv preprint arXiv:2010.03759*, 2020.
- R. Salakhutdinov and G. Hinton. Deep boltzmann machines. In *Artificial intelligence and statistics*, pages 448–455. PMLR, 2009.
- P. Smolensky. Information processing in dynamical systems: Foundations of harmony theory. Technical report, Colorado Univ at Boulder Dept of Computer Science, 1986.
- Y. Song and D. P. Kingma. How to train your energy-based models. *arXiv preprint arXiv:2101.03288*, 2021.

References V

- I. Sutskever and T. Tieleman. On the convergence properties of contrastive divergence. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 789–795. JMLR Workshop and Conference Proceedings, 2010.
- P. Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.
- Z. Xiao, K. Kreis, J. Kautz, and A. Vahdat. Vaebm: A symbiosis between variational autoencoders and energy-based models. In *International Conference on Learning Representations*, 2020.
- O. Yair and T. Michaeli. Contrastive divergence learning is a time reversal adversarial game. *arXiv preprint arXiv:2012.03295*, 2020.