

PML 11. Energy-Based Models. Chapter 24, Probabilistic Machine Learning: Advanced Topics

Probabilistic Machine Learning Reading Group

Víctor Gallego

April 8, 2026

Komorebi AI, ICMAT

- Introduction & Motivation
- Maximum Likelihood Training
- Score Matching
- Adversarial Training
- Other Training Methods
- EBMs in Practice: Two Case Studies
- Summary

Introduction & Motivation

Maximum Likelihood Training

Score Matching

Adversarial Training

Other Training Methods

EBMs in Practice: Two Case Studies

Summary

Why Energy-Based Models?

The generative-model zoo so far:

- VAEs, autoregressive models, normalizing flows: all *directed* graphical models.
- They generate data one step at a time, using **locally normalized** distributions.

Sometimes it is easier to specify *constraints* than a generation process.

- Example: proteins that are *thermally stable* **and** *bind to a target receptor*.
- Each constraint is easy to model; combining them in a directed model is hard.

⇒ **Idea:** specify an *unnormalized* distribution via an **energy function**.

An **Energy-Based Model** defines a Gibbs distribution:

$$p_{\theta}(\mathbf{x}) = \frac{\exp(-\mathcal{E}_{\theta}(\mathbf{x}))}{Z_{\theta}}$$

- $\mathcal{E}_{\theta}(\mathbf{x})$: the **energy function**, any function $\mathbb{R}^d \rightarrow \mathbb{R}$ parameterized by θ .
- $Z_{\theta} = \int \exp(-\mathcal{E}_{\theta}(\mathbf{x})) d\mathbf{x}$: the **partition function** (intractable!).

Advantage: \mathcal{E}_{θ} is *unrestricted*, any neural network that outputs a scalar.

No invertibility, no autoregressive ordering, no tractable Jacobian required.

Challenge: we can never compute Z_{θ} (or even its gradient) exactly.

Example: Product of Experts

Suppose we have two “expert” distributions $p_1(\mathbf{x})$ and $p_2(\mathbf{x})$.

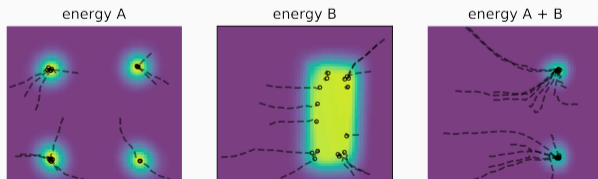
Product of Experts (PoE) [Hinton, 2002]:

$$p_{12}(\mathbf{x}) = \frac{1}{Z_{12}} p_1(\mathbf{x}) p_2(\mathbf{x}).$$

If each expert is an EBM, the combined model is also an EBM with

$$\mathcal{E}_{12}(\mathbf{x}) = \mathcal{E}_1(\mathbf{x}) + \mathcal{E}_2(\mathbf{x}).$$

Each energy acts as a **soft constraint** on the data.

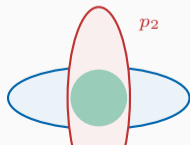


Product of Experts vs. Mixture of Experts

Product of Experts (PoE):

$$p_{\text{PoE}}(\mathbf{x}) = \frac{1}{Z} \prod_m p_m(\mathbf{x})$$

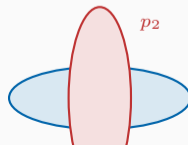
- High probability only where **all experts agree**
- Energies **add**: $\mathcal{E}_{12} = \mathcal{E}_1 + \mathcal{E}_2$
- Result is **sharper** than any single expert
- Natural for EBMs



Mixture of Experts (MoE):

$$p_{\text{MoE}}(\mathbf{x}) = \sum_m \pi_m p_m(\mathbf{x})$$

- High probability where **any expert** assigns mass
- Densities **average**: covers the union
- Result is **broader** than any single expert
- Natural for directed models



*Everything about EBMs is easy...
except for the partition function Z_θ .*

What Z_θ blocks:

- Evaluating $p_\theta(\mathbf{x})$
- Computing the likelihood
- Drawing exact samples

Strategies (this talk):

1. Maximum likelihood + MCMC
2. Score matching
3. Adversarial training
4. Other methods (NCE, Stein, ...)

Introduction & Motivation

Maximum Likelihood Training

Score Matching

Adversarial Training

Other Training Methods

EBMs in Practice: Two Case Studies

Summary

Standard MLE: maximize expected log-likelihood

$$\ell(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x} \sim p_{\mathcal{D}}(\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x})] = -D_{\text{KL}}(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\boldsymbol{\theta}}(\mathbf{x})) + \text{const.}$$

Standard MLE: maximize expected log-likelihood

$$\ell(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x} \sim p_{\mathcal{D}}(\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x})] = -D_{\text{KL}}(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\boldsymbol{\theta}}(\mathbf{x})) + \text{const.}$$

Problem: $\log p_{\boldsymbol{\theta}}(\mathbf{x}) = -\mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x}) - \log Z_{\boldsymbol{\theta}}$. We can't compute $\log Z_{\boldsymbol{\theta}}$.

Standard MLE: maximize expected log-likelihood

$$\ell(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x} \sim p_{\mathcal{D}}(\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x})] = -D_{\text{KL}}(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\boldsymbol{\theta}}(\mathbf{x})) + \text{const.}$$

Problem: $\log p_{\boldsymbol{\theta}}(\mathbf{x}) = -\mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x}) - \log Z_{\boldsymbol{\theta}}$. We can't compute $\log Z_{\boldsymbol{\theta}}$.

Main insight: we don't need $Z_{\boldsymbol{\theta}}$ itself, only its *gradient* w.r.t. $\boldsymbol{\theta}$:

$$\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x}) = \underbrace{-\nabla_{\boldsymbol{\theta}} \mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x})}_{\text{easy (autodiff)}} - \underbrace{\nabla_{\boldsymbol{\theta}} \log Z_{\boldsymbol{\theta}}}_{\text{hard}}$$

The Log-Partition Gradient

Theorem. The gradient of $\log Z_{\theta}$ equals a model expectation:

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

The Log-Partition Gradient

Theorem. The gradient of $\log Z_{\theta}$ equals a model expectation:

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

Proof sketch:

$$\begin{aligned} \nabla_{\theta} \log Z_{\theta} &= \frac{1}{Z_{\theta}} \int \exp(-\mathcal{E}_{\theta}(\mathbf{x})) (-\nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x})) \, d\mathbf{x} \\ &= \int p_{\theta}(\mathbf{x}) (-\nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x})) \, d\mathbf{x} = \mathbb{E}_{p_{\theta}}[-\nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x})]. \end{aligned}$$

The Log-Partition Gradient

Theorem. The gradient of $\log Z_\theta$ equals a model expectation:

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

Proof sketch:

$$\begin{aligned} \nabla_{\theta} \log Z_{\theta} &= \frac{1}{Z_{\theta}} \int \exp(-\mathcal{E}_{\theta}(\mathbf{x})) (-\nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x})) \, d\mathbf{x} \\ &= \int p_{\theta}(\mathbf{x}) (-\nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x})) \, d\mathbf{x} = \mathbb{E}_{p_{\theta}}[-\nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x})]. \end{aligned}$$

Full gradient of log-likelihood:

$$\nabla_{\theta} \ell = \underbrace{-\mathbb{E}_{p_{\mathcal{D}}}[\nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x})]}_{\text{data term ("clamped")}} + \underbrace{\mathbb{E}_{p_{\theta}}[\nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x})]}_{\text{model term ("unclamped")}}$$

Push data energies **down**, push model-sample energies **up**.

Sampling with Langevin MCMC

To estimate the model expectation, we need samples $\tilde{\mathbf{x}}_s \sim p_{\theta}(\mathbf{x})$.

Observation, the score function:

$$\nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) = -\nabla_{\mathbf{x}} \mathcal{E}_{\theta}(\mathbf{x}) - \underbrace{\nabla_{\mathbf{x}} \log Z_{\theta}}_{=0} = -\nabla_{\mathbf{x}} \mathcal{E}_{\theta}(\mathbf{x}).$$

The score is *free*: no partition function involved!

Sampling with Langevin MCMC

To estimate the model expectation, we need samples $\tilde{\mathbf{x}}_s \sim p_{\theta}(\mathbf{x})$.

Observation, the score function:

$$\nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) = -\nabla_{\mathbf{x}} \mathcal{E}_{\theta}(\mathbf{x}) - \underbrace{\nabla_{\mathbf{x}} \log Z_{\theta}}_{=0} = -\nabla_{\mathbf{x}} \mathcal{E}_{\theta}(\mathbf{x}).$$

The score is *free*: no partition function involved!

Langevin MCMC (overdamped Langevin diffusion, discretized):

$$\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k - \underbrace{\frac{\epsilon^2}{2} \nabla_{\mathbf{x}} \mathcal{E}_{\theta}(\mathbf{x}^k)}_{\text{gradient step}} + \underbrace{\epsilon \mathbf{z}^k}_{\text{noise}}, \quad \mathbf{z}^k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

As $\epsilon \rightarrow 0$ and $K \rightarrow \infty$, the iterates converge to $p_{\theta}(\mathbf{x})$.

Persistent MCMC-SGD (PCD)

Running MCMC to convergence at every step is **expensive**. [Younes, 1989; Tieleman, 2008]: parameters change slowly, so the chain from step $t-1$ is still a good starting point at step t .

Algorithm (Persistent MCMC-SGD):

1. Initialize parameters θ and chains $\tilde{x}_{1:S}$ randomly.
2. For each training step t :
 - Compute data gradients: $\mathbf{g}_b = \nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x}_b)$
 - **Continue** each chain \tilde{x}_s for N Langevin steps (don't restart!)
 - Compute model gradients: $\tilde{\mathbf{g}}_s = \nabla_{\theta} \mathcal{E}_{\theta}(\tilde{\mathbf{x}}_s)$
 - Update: $\theta \ += \ \eta \left(-\frac{1}{B} \sum \mathbf{g}_b - \frac{1}{S} \sum \tilde{\mathbf{g}}_s \right)$

Replay buffer variant [Du & Mordatch, 2019]:

- Store historical MCMC states in a buffer
- Init new chains by sampling from buffer (95%) or noise (5%)
- Encourages diversity, avoids mode collapse

Caveat: truncated MCMC \Rightarrow biased gradients. Corrections exist (coupled MCMC, entropy terms).

Historical Note: CD and RBMs

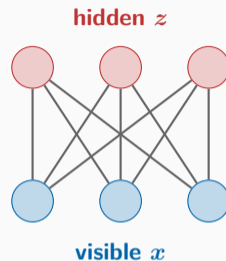
Contrastive Divergence [Hinton, 2002]:

- Initialize chain at the *data point* \mathbf{x}_n , run only T steps (often $T=1$)
- Minimizes
$$CD_T = D_{\text{KL}}(p_0 \| p_\infty) - D_{\text{KL}}(p_T \| p_\infty)$$
- Biased but fast; historically important for training **Restricted Boltzmann Machines**

RBMs: latent-variable EBM with bipartite structure and block Gibbs sampling:

$$\mathcal{E}(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = \sum_{d,k} x_d z_k W_{dk} + \sum_d x_d b_d + \sum_k z_k c_k$$

CD-1 gradient: $\nabla_{\mathbf{W}} \ell \approx \mathbf{x}_n \boldsymbol{\mu}_n^\top - \mathbf{x}'_n (\boldsymbol{\mu}'_n)^\top$
(Hebbian – anti-Hebbian).



PCD superseded CD: persistent chains explore better and don't need the bipartite trick.

Introduction & Motivation

Maximum Likelihood Training

Score Matching

Adversarial Training

Other Training Methods

EBMs in Practice: Two Case Studies

Summary

Score Matching: Idea

Observation: if two log-pdfs have the *same first derivatives* everywhere, they are equal (up to normalization).

Define the **(Stein) score function**:

$$s_{\theta}(\mathbf{x}) \triangleq \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) = -\nabla_{\mathbf{x}} \mathcal{E}_{\theta}(\mathbf{x}). \quad (\text{no } Z_{\theta}!)$$

Score Matching: Idea

Observation: if two log-pdfs have the *same first derivatives* everywhere, they are equal (up to normalization).

Define the **(Stein) score function**:

$$s_{\theta}(\mathbf{x}) \triangleq \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) = -\nabla_{\mathbf{x}} \mathcal{E}_{\theta}(\mathbf{x}). \quad (\text{no } Z_{\theta}!)$$

Score Matching [Hyvärinen, 2005] minimizes the **Fisher divergence**:

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

Problem: the data score $\nabla_{\mathbf{x}} \log p_{\mathcal{D}}(\mathbf{x})$ is unknown!

Under regularity conditions, the Fisher divergence can be rewritten as:

$$D_F(p_{\mathcal{D}} \parallel p_{\theta}) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} \left[\frac{1}{2} \|\mathbf{s}_{\theta}(\mathbf{x})\|^2 + \text{tr}(\mathbf{J}_{\mathbf{x}} \mathbf{s}_{\theta}(\mathbf{x})) \right] + \text{const}$$

- The unknown $\nabla_{\mathbf{x}} \log p_{\mathcal{D}}$ has **disappeared!**
- Only involves the model score \mathbf{s}_{θ} and its Jacobian.
- The constant doesn't depend on θ , just drop it.

Under regularity conditions, the Fisher divergence can be rewritten as:

$$D_F(p_{\mathcal{D}} \parallel p_{\theta}) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} \left[\frac{1}{2} \|\mathbf{s}_{\theta}(\mathbf{x})\|^2 + \text{tr}(\mathbf{J}_{\mathbf{x}} \mathbf{s}_{\theta}(\mathbf{x})) \right] + \text{const}$$

- The unknown $\nabla_{\mathbf{x}} \log p_{\mathcal{D}}$ has **disappeared!**
- Only involves the model score \mathbf{s}_{θ} and its Jacobian.
- The constant doesn't depend on θ , just drop it.

Downside: computing $\text{tr}(\mathbf{J}_{\mathbf{x}} \mathbf{s}_{\theta})$ costs $O(d^2)$: one backprop per input dimension.

⇒ Only practical for low-dimensional or specially structured energy functions.

Denosing Score Matching (DSM)

[Vincent, 2011] — an elegant workaround to avoid second derivatives entirely.

Idea: corrupt data with noise $\tilde{\mathbf{x}} = \mathbf{x} + \epsilon$, $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$.

The noisy conditional has a *known* score:

$$\nabla_{\tilde{\mathbf{x}}} \log q(\tilde{\mathbf{x}} | \mathbf{x}) = \frac{\mathbf{x} - \tilde{\mathbf{x}}}{\sigma^2}.$$

Denosing Score Matching (DSM)

[Vincent, 2011] — an elegant workaround to avoid second derivatives entirely.

Idea: corrupt data with noise $\tilde{\mathbf{x}} = \mathbf{x} + \epsilon$, $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$.

The noisy conditional has a *known* score:

$$\nabla_{\tilde{\mathbf{x}}} \log q(\tilde{\mathbf{x}} | \mathbf{x}) = \frac{\mathbf{x} - \tilde{\mathbf{x}}}{\sigma^2}.$$

DSM objective:

$$\frac{1}{2} \mathbb{E}_{q(\mathbf{x}, \tilde{\mathbf{x}})} \left[\left\| \mathbf{s}_{\theta}(\tilde{\mathbf{x}}) - \frac{\mathbf{x} - \tilde{\mathbf{x}}}{\sigma^2} \right\|^2 \right]$$

Train the score network to **denoise**: predict the direction back to clean data.

Pro: no Jacobian traces, simple to implement.

Con: learns the score of the *noisy* distribution $q(\tilde{\mathbf{x}})$, not $p_{\mathcal{D}}(\mathbf{x})$.

Sliced Score Matching (SSM)

[Song et al., 2019] : project scores onto random directions.

Sliced Fisher divergence:

$$D_{SF}(p_{\mathcal{D}} \parallel p_{\theta}) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} \mathbb{E}_{p(\mathbf{v})} \left[\frac{1}{2} (\mathbf{v}^{\top} \nabla_{\mathbf{x}} \log p_{\mathcal{D}}(\mathbf{x}) - \mathbf{v}^{\top} \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}))^2 \right]$$

After integration by parts (same trick!), with $p(\mathbf{v}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$:

$$D_{SF} = \mathbb{E}_{p_{\mathcal{D}}} \mathbb{E}_{\mathbf{v} \sim \mathcal{N}} \left[\frac{1}{2} (\mathbf{v}^{\top} \mathbf{s}_{\theta}(\mathbf{x}))^2 + \mathbf{v}^{\top} \mathbf{J} \mathbf{v} \right] + C$$

Autodiff theory: the Hessian-vector product $\mathbf{J} \mathbf{v}$ costs $O(d)$, not $O(d^2)$.

Method	Consistent?	Cost
Basic SM	✓	$O(d^2)$
DSM	× (noisy target)	$O(d)$
SSM	✓	$O(d)$

Connection: Score Matching \longleftrightarrow Contrastive Divergence

These are more related than they look! [Hyvärinen, 2007a]

Consider 1-step Langevin CD with step size ϵ :

$$\mathbf{x}^1 = \mathbf{x} - \frac{\epsilon}{2} \nabla_{\mathbf{x}} \mathcal{E}_{\theta}(\mathbf{x}) + \sqrt{\epsilon} \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

Connection: Score Matching \longleftrightarrow Contrastive Divergence

These are more related than they look! [Hyvärinen, 2007a]

Consider 1-step Langevin CD with step size ϵ :

$$\mathbf{x}^1 = \mathbf{x} - \frac{\epsilon}{2} \nabla_{\mathbf{x}} \mathcal{E}_{\theta}(\mathbf{x}) + \sqrt{\epsilon} \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

Taylor-expanding the CD gradient in ϵ gives:

$$\mathbb{E}_{p_{\mathcal{D}}}[\nabla_{\theta} \log p_{\theta}(\mathbf{x})] \approx \frac{\epsilon^2}{2} \nabla_{\theta} D_F(p_{\mathcal{D}} \parallel p_{\theta}) + o(\epsilon^2).$$

Takeaway: contrastive divergence (KL-based) and score matching (Fisher-based) are linked by **de Bruijn's identity**:

$$\frac{d}{dt} D_{\text{KL}}(q_t \parallel p_{\theta,t}) = -\frac{1}{2} D_F(q_t \parallel p_{\theta,t}),$$

where $q_t, p_{\theta,t}$ are Gaussian-smoothed versions (variance t) of $p_{\mathcal{D}}, p_{\theta}$.

Introduction & Motivation

Maximum Likelihood Training

Score Matching

Adversarial Training

Other Training Methods

EBMs in Practice: Two Case Studies

Summary

Motivation: Avoiding MCMC with a Learned Sampler

Recall: MLE requires samples $\tilde{\mathbf{x}} \sim p_{\theta}(\mathbf{x})$ at every gradient step.

Problems with MCMC sampling:

- Multiple MCMC steps per training iteration \Rightarrow expensive
- Requires careful tuning of step size, chain length, ...
- Truncated chains \Rightarrow biased gradients

Motivation: Avoiding MCMC with a Learned Sampler

Recall: MLE requires samples $\tilde{\mathbf{x}} \sim p_{\theta}(\mathbf{x})$ at every gradient step.

Problems with MCMC sampling:

- Multiple MCMC steps per training iteration \Rightarrow **expensive**
- Requires careful tuning of step size, chain length, ...
- Truncated chains \Rightarrow biased gradients

Idea: learn an **auxiliary model** $q_{\phi}(\mathbf{x})$ that approximates $p_{\theta}(\mathbf{x})$.

- Replace MCMC sampling with **direct sampling** from q_{ϕ}
- q_{ϕ} acts as a *variational approximation* to the EBM
- Leads naturally to a **minimax** (adversarial) objective

Deriving the Minimax Objective (I)

Goal: maximize the expected log-likelihood $\mathbb{E}_{p_{\mathcal{D}}}[\log p_{\theta}(\mathbf{x})]$.

Introduce a variational distribution $q_{\phi}(\mathbf{x})$ and expand:

$$\begin{aligned}\mathbb{E}_{p_{\mathcal{D}}}[\log p_{\theta}(\mathbf{x})] &= \mathbb{E}_{p_{\mathcal{D}}}[-\mathcal{E}_{\theta}(\mathbf{x})] - \log Z_{\theta} \\ &= \mathbb{E}_{p_{\mathcal{D}}}[-\mathcal{E}_{\theta}(\mathbf{x})] - \log \int e^{-\mathcal{E}_{\theta}(\mathbf{x})} d\mathbf{x} \\ &= \mathbb{E}_{p_{\mathcal{D}}}[-\mathcal{E}_{\theta}(\mathbf{x})] - \log \int q_{\phi}(\mathbf{x}) \frac{e^{-\mathcal{E}_{\theta}(\mathbf{x})}}{q_{\phi}(\mathbf{x})} d\mathbf{x}\end{aligned}$$

Deriving the Minimax Objective (I)

Goal: maximize the expected log-likelihood $\mathbb{E}_{p_{\mathcal{D}}}[\log p_{\theta}(\mathbf{x})]$.

Introduce a variational distribution $q_{\phi}(\mathbf{x})$ and expand:

$$\begin{aligned}\mathbb{E}_{p_{\mathcal{D}}}[\log p_{\theta}(\mathbf{x})] &= \mathbb{E}_{p_{\mathcal{D}}}[-\mathcal{E}_{\theta}(\mathbf{x})] - \log Z_{\theta} \\ &= \mathbb{E}_{p_{\mathcal{D}}}[-\mathcal{E}_{\theta}(\mathbf{x})] - \log \int e^{-\mathcal{E}_{\theta}(\mathbf{x})} d\mathbf{x} \\ &= \mathbb{E}_{p_{\mathcal{D}}}[-\mathcal{E}_{\theta}(\mathbf{x})] - \log \int q_{\phi}(\mathbf{x}) \frac{e^{-\mathcal{E}_{\theta}(\mathbf{x})}}{q_{\phi}(\mathbf{x})} d\mathbf{x}\end{aligned}$$

Apply **Jensen's inequality** ($-\log$ is convex):

$$-\log \int q_{\phi}(\mathbf{x}) \frac{e^{-\mathcal{E}_{\theta}(\mathbf{x})}}{q_{\phi}(\mathbf{x})} d\mathbf{x} \leq - \int q_{\phi}(\mathbf{x}) \log \frac{e^{-\mathcal{E}_{\theta}(\mathbf{x})}}{q_{\phi}(\mathbf{x})} d\mathbf{x}$$

Deriving the Minimax Objective (II)

Expanding the RHS of Jensen's bound:

$$\begin{aligned} - \int q_\phi(\mathbf{x}) \log \frac{e^{-\mathcal{E}_\theta(\mathbf{x})}}{q_\phi(\mathbf{x})} d\mathbf{x} &= - \int q_\phi(\mathbf{x}) \left[-\mathcal{E}_\theta(\mathbf{x}) - \log q_\phi(\mathbf{x}) \right] d\mathbf{x} \\ &= \mathbb{E}_{q_\phi}[\mathcal{E}_\theta(\mathbf{x})] - \underbrace{H(q_\phi)}_{\int q \log q} \end{aligned}$$

Deriving the Minimax Objective (II)

Expanding the RHS of Jensen's bound:

$$\begin{aligned} - \int q_\phi(\mathbf{x}) \log \frac{e^{-\mathcal{E}_\theta(\mathbf{x})}}{q_\phi(\mathbf{x})} d\mathbf{x} &= - \int q_\phi(\mathbf{x}) \left[-\mathcal{E}_\theta(\mathbf{x}) - \log q_\phi(\mathbf{x}) \right] d\mathbf{x} \\ &= \mathbb{E}_{q_\phi}[\mathcal{E}_\theta(\mathbf{x})] - \underbrace{H(q_\phi)}_{\int q \log q} \end{aligned}$$

Combining with the data term, we get an **upper bound on expected log-likelihood**:

$$\mathbb{E}_{p_{\mathcal{D}}}[\log p_\theta(\mathbf{x})] \leq \underbrace{\mathbb{E}_{q_\phi}[-\mathcal{E}_\theta(\mathbf{x})] - \mathbb{E}_{p_{\mathcal{D}}}[-\mathcal{E}_\theta(\mathbf{x})]}_{\text{energy gap (sampler vs. data)}} - \underbrace{H(q_\phi)}_{\text{entropy of sampler}}$$

Deriving the Minimax Objective (II)

Expanding the RHS of Jensen's bound:

$$\begin{aligned} - \int q_\phi(\mathbf{x}) \log \frac{e^{-\mathcal{E}_\theta(\mathbf{x})}}{q_\phi(\mathbf{x})} d\mathbf{x} &= - \int q_\phi(\mathbf{x}) \left[-\mathcal{E}_\theta(\mathbf{x}) - \log q_\phi(\mathbf{x}) \right] d\mathbf{x} \\ &= \mathbb{E}_{q_\phi}[\mathcal{E}_\theta(\mathbf{x})] - \underbrace{H(q_\phi)}_{\int q \log q} \end{aligned}$$

Combining with the data term, we get an **upper bound on expected log-likelihood**:

$$\mathbb{E}_{p_{\mathcal{D}}}[\log p_\theta(\mathbf{x})] \leq \underbrace{\mathbb{E}_{q_\phi}[-\mathcal{E}_\theta(\mathbf{x})] - \mathbb{E}_{p_{\mathcal{D}}}[-\mathcal{E}_\theta(\mathbf{x})]}_{\text{energy gap (sampler vs. data)}} - \underbrace{H(q_\phi)}_{\text{entropy of sampler}}$$

Tightness: the bound is tight when $q_\phi(\mathbf{x}) = p_\theta(\mathbf{x})$, i.e. when the sampler perfectly matches the EBM.

Minimize the upper bound w.r.t. ϕ , maximize likelihood w.r.t. θ :

$$\max_{\theta} \min_{\phi} \mathbb{E}_{q_{\phi}}[\mathcal{E}_{\theta}(\mathbf{x})] - \mathbb{E}_{p_{\mathcal{D}}}[\mathcal{E}_{\theta}(\mathbf{x})] - H(q_{\phi})$$

Inner min (over ϕ):

- q_{ϕ} learns to generate **low-energy** samples
- Entropy $H(q_{\phi})$ prevents collapse to a single mode
- $q_{\phi} \rightarrow p_{\theta}$ at optimality

Outer max (over θ):

- Assign low energy to **real data**
- Assign high energy to **sampler outputs**
- Widens the energy gap

Analogy with GANs (Ch. 26): q_{ϕ} plays the role of the *generator*, \mathcal{E}_{θ} the *discriminator*.

But here \mathcal{E}_{θ} is an **explicit energy model**: $p_{\theta} \propto e^{-\mathcal{E}_{\theta}}$ defines a proper density.

Implementing q_ϕ : Normalizing Flows vs. Neural Samplers

Requirements for q_ϕ : (1) **fast sampling** (replace MCMC), and (2) **tractable entropy** $H(q_\phi)$.

Normalizing flows [Dai+19]

- Invertible $x = f_\phi(z)$, $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- Exact density via change-of-variables
- $H(q_\phi)$ tractable via $\log |\det \mathbf{J}_{f_\phi}|$
- ✓ Both requirements satisfied
- ✗ Invertibility constrains architecture

Neural samplers [KB16; Zha+16]

- Deterministic $x = g_\phi(z)$, $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- Like a GAN generator
- ✓ Unrestricted architecture
- ✓ Fast sampling
- ✗ $H(q_\phi)$ **intractable**

Implementing q_ϕ : Normalizing Flows vs. Neural Samplers

Requirements for q_ϕ : (1) **fast sampling** (replace MCMC), and (2) **tractable entropy** $H(q_\phi)$.

Normalizing flows [Dai+19]

- Invertible $\mathbf{x} = f_\phi(\mathbf{z})$, $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- Exact density via change-of-variables
- $H(q_\phi)$ tractable via $\log |\det \mathbf{J}_{f_\phi}|$
- ✓ Both requirements satisfied
- ✗ Invertibility constrains architecture

Neural samplers [KB16; Zha+16]

- Deterministic $\mathbf{x} = g_\phi(\mathbf{z})$, $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- Like a GAN generator
- ✓ Unrestricted architecture
- ✓ Fast sampling
- ✗ $H(q_\phi)$ **intractable**

Entropy estimation heuristics for neural samplers:

- Mutual information: $H(q_\phi(\mathbf{x})) = I(g_\phi(\mathbf{z}), \mathbf{z}) + H(\mathbf{z})$ [Kumar+19]
- Variational lower bounds [Nguyen+10; Nowozin+16]
- With fixed base p_0 : replace $-H(q_\phi)$ by $D_{\text{KL}}(q_\phi \| p_0)$ [Dai+19]

Main idea: add noise to make the entropy gradient tractable.

$$\mathbf{x} = g_{\phi}(\mathbf{z}) + \sigma \boldsymbol{\epsilon}, \quad \mathbf{z}, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

This defines $q_{\phi}(\mathbf{x}) = \int \mathcal{N}(\mathbf{x} \mid g_{\phi}(\mathbf{z}), \sigma^2 \mathbf{I}) p(\mathbf{z}) d\mathbf{z}$.

Main idea: add noise to make the entropy gradient tractable.

$$\mathbf{x} = g_{\phi}(\mathbf{z}) + \sigma \boldsymbol{\epsilon}, \quad \mathbf{z}, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

This defines $q_{\phi}(\mathbf{x}) = \int \mathcal{N}(\mathbf{x} \mid g_{\phi}(\mathbf{z}), \sigma^2 \mathbf{I}) p(\mathbf{z}) d\mathbf{z}$.

Why this helps:

- The gradient $\nabla_{\phi} H(q_{\phi})$ admits unbiased Monte Carlo estimates
- The noisy sampler is a latent variable model with Gaussian conditionals
- Enables **gradient-based optimization** of the full minimax objective

Main idea: add noise to make the entropy gradient tractable.

$$\mathbf{x} = g_{\phi}(\mathbf{z}) + \sigma \boldsymbol{\epsilon}, \quad \mathbf{z}, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

This defines $q_{\phi}(\mathbf{x}) = \int \mathcal{N}(\mathbf{x} \mid g_{\phi}(\mathbf{z}), \sigma^2 \mathbf{I}) p(\mathbf{z}) d\mathbf{z}$.

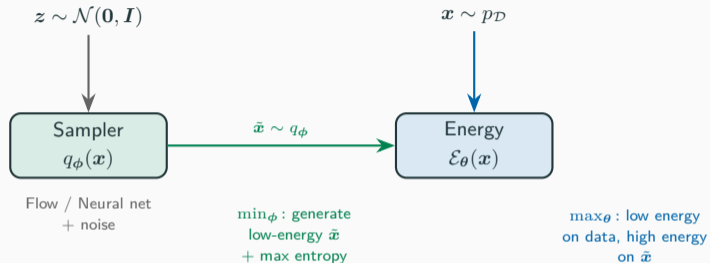
Why this helps:

- The gradient $\nabla_{\phi} H(q_{\phi})$ admits unbiased Monte Carlo estimates
- The noisy sampler is a latent variable model with Gaussian conditionals
- Enables **gradient-based optimization** of the full minimax objective

Connection to normalizing flows:

- A normalizing flow is a *deterministic* invertible sampler ($\sigma = 0$)
- A noisy neural sampler is a *stochastic, non-invertible* sampler ($\sigma > 0$)
- The noise trades off fidelity for architectural freedom and entropy tractability

Adversarial Training: The Full Picture



- **No MCMC** at training time: q_ϕ provides fast approximate samples
- **Explicit density model**: unlike GANs, \mathcal{E}_θ defines $p_\theta \propto e^{-\mathcal{E}_\theta}$
- Can combine with short-run MCMC for refinement [Xie+18]
- Extends to arbitrary f -divergences beyond KL [Yu+20]

Introduction & Motivation

Maximum Likelihood Training

Score Matching

Adversarial Training

Other Training Methods

EBMs in Practice: Two Case Studies

Summary

Noise Contrastive Estimation [GH10]

Reduce density estimation to binary classification (data vs. noise p_n). At the optimum $p_{\theta^*} \equiv p_{\mathcal{D}}$. Bonus: Z_{θ} is a learnable scalar.

With $p_n = p_{\mathcal{D}}(\cdot - \mathbf{v})$, $\|\mathbf{v}\| \rightarrow 0$: recovers SSM.

KL Divergence Tricks [Mov08; SDBD11; Lyu11]

Z_{θ} cancels in *differences/derivatives* of KL:

- Min. velocity / probability flow:
minimize $\frac{d}{dt} D_{\text{KL}}(p_t \| p_{\theta}) \Big|_{t=0}$
- Min. KL contraction:
 $D_{\text{KL}}(p_{\mathcal{D}} \| p_{\theta}) - D_{\text{KL}}(\Phi\{p_{\mathcal{D}}\} \| \Phi\{p_{\theta}\})$

Stein Discrepancy [GM15; LLJ16]

$$D_{\text{Stein}} := \sup_{\mathbf{f} \in \mathcal{F}} \mathbb{E}_{p_{\mathcal{D}}} [\nabla_{\mathbf{x}} \log p_{\theta}^{\top} \mathbf{f} + \text{tr}(\nabla_{\mathbf{x}} \mathbf{f})]$$

- Only needs the **score**: no Z_{θ}
- $\mathcal{F} = \text{RKHS ball} \Rightarrow$ **KSD**: trace is free, scalable

Common theme: every method finds a clever way to *avoid computing* the intractable partition function Z_{θ} .

Introduction & Motivation

Maximum Likelihood Training

Score Matching

Adversarial Training

Other Training Methods

EBMs in Practice: Two Case Studies

Summary

IGEBM: Scaling EBMs to High Dimensions

[Du & Mordatch, NeurIPS 2019] — the theory we've seen, made to work on real images.

Recipe: MLE + Langevin dynamics + three practical tricks:

1. **Replay buffer:** store past MCMC samples; reinitialize 95% of chains from the buffer, 5% from uniform noise. Improves mixing and diversity.
2. **Spectral normalization:** constrain the Lipschitz constant of \mathcal{E}_θ to stabilize Langevin dynamics (smooth energy landscape \Rightarrow stable gradients).
3. **Energy regularization:** weakly penalize $\mathcal{E}_\theta(\mathbf{x}^+)^2 + \mathcal{E}_\theta(\mathbf{x}^-)^2$ to keep energy magnitudes bounded.

Training:

$$\Delta\theta \propto \underbrace{\nabla_{\theta}\mathcal{E}_{\theta}(\mathbf{x}^{+})}_{\text{push data energy down}} - \underbrace{\nabla_{\theta}\mathcal{E}_{\theta}(\mathbf{x}^{-})}_{\text{push sample energy up}} + \alpha \nabla_{\theta}(\mathcal{E}_{\theta}(\mathbf{x}^{+})^2 + \mathcal{E}_{\theta}(\mathbf{x}^{-})^2)$$

where $\mathbf{x}^+ \sim p_{\mathcal{D}}$ and \mathbf{x}^- is obtained via K steps of Langevin dynamics.

IGEBM: Generation Results

Competitive image generation on CIFAR-10 and ImageNet:

Model	IS \uparrow	FID \downarrow
<i>CIFAR-10 Unconditional</i>		
PixelCNN	4.60	65.93
WGAN+GP	6.50	36.4
EBM (single)	6.02	40.58
SNGAN	8.22	21.7
<i>ImageNet 32\times32 Conditional</i>		
PixelIQN	10.18	22.99
EBM (single)	18.22	14.31

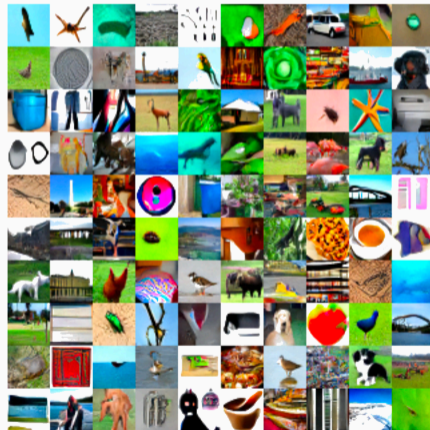


Figure 2: Conditional ImageNet32x32 EBM samples

EBMs as a versatile building block:

Adversarial robustness

- Conditional EBMs on CIFAR-10 achieve 49.6% accuracy under PGD attack, *without* explicit adversarial training
- Outperform L_∞ SOTA
- Implicit generation “refines” inputs toward high-density regions

Trajectory modeling

- Robot hand manipulation: EBMs capture multimodal transitions

Continual learning

- Split MNIST: **64.99%** (EBM) vs. 40.04% (VAE)
- Replay buffer acts as a natural episodic memory
- No task-specific regularization needed

Some advantages over GANs/VAEs:

- ✓ Single network (no generator/discriminator)
- ✓ Adaptive computation (more MCMC = better samples)

JEM: Your Classifier Is Secretly an EBM

[Grathwohl et al., ICLR 2020] — a **nice reinterpretation** of standard classifiers.

A classifier with logits $f_{\theta}(\mathbf{x}) \in \mathbb{R}^K$ defines $p_{\theta}(y | \mathbf{x})$ via softmax.
But the *same* logits also define:

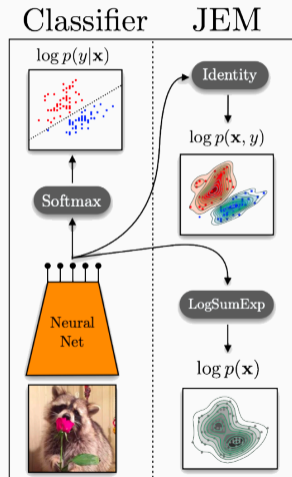
Joint: $p_{\theta}(\mathbf{x}, y) = \exp(f_{\theta}(\mathbf{x})[y]) / Z(\theta)$

Marginal: $p_{\theta}(\mathbf{x}) = \sum_y \exp(f_{\theta}(\mathbf{x})[y]) / Z(\theta)$

$\Rightarrow \mathcal{E}_{\theta}(\mathbf{x}) = -\text{LogSumExp}_y f_{\theta}(\mathbf{x})[y]$

Punchline: dividing $p_{\theta}(\mathbf{x}, y)$ by $p_{\theta}(\mathbf{x})$, $Z(\theta)$ **cancels** \Rightarrow we recover the standard softmax $p_{\theta}(y | \mathbf{x})$!

\Rightarrow Every classifier has a **generative model hidden inside**.
Grathwohl et al. propose to *actually train it* — the **Joint Energy-based Model (JEM)**.



Factor the joint log-likelihood:

$$\log p_{\theta}(\mathbf{x}, y) = \underbrace{\log p_{\theta}(y | \mathbf{x})}_{\text{cross-entropy (standard classification)}} + \underbrace{\log p_{\theta}(\mathbf{x})}_{\text{EBM training (SGLD / PCD)}}$$

- The discriminative term is **exact**: just standard cross-entropy.
- The generative term uses SGLD with persistent chains and replay buffer (following Du & Mordatch).
- Architecture: Wide ResNet.

JEM: Training and the Hybrid Objective

Factor the joint log-likelihood:

$$\log p_{\theta}(\mathbf{x}, y) = \underbrace{\log p_{\theta}(y | \mathbf{x})}_{\text{cross-entropy (standard classification)}} + \underbrace{\log p_{\theta}(\mathbf{x})}_{\text{EBM training (SGLD / PCD)}}$$

- The discriminative term is **exact**: just standard cross-entropy.
- The generative term uses SGLD with persistent chains and replay buffer (following Du & Mordatch).
- Architecture: Wide ResNet.

CIFAR-10 results: best of both worlds:

Model	Accuracy \uparrow	IS \uparrow	FID \downarrow
Wide-ResNet (disc. only)	95.8%	N/A	N/A
SNGAN (gen. only)	N/A	8.59	25.5
JEM (ours)	92.9%	8.76	38.4

Calibration (CIFAR-100):

- Baseline WRN: ECE = 22.32%
- **JEM: ECE = 4.87%**
- Nearly perfectly calibrated, *without* post-hoc methods (Platt scaling, etc.)

Adversarial robustness:

- Competitive with explicit adversarial training under L_∞ and L_2 PGD attacks
- MCMC refinement at test time “undoes” adversarial perturbations
- Doesn’t confidently classify noise images (unlike standard CNNs)

OOD detection (AUROC):

Score	Glow	JEM
$s = \log p_\theta(x)$		
SVHN	.05	.67
CelebA	.57	.75
$s = -\ \nabla_x \log p_\theta\ $ (approx. mass)		
SVHN	.27	.95
CelebA	.29	.79

New OOD score: norm of the score function distinguishes typical-set points from high-density anomalies.

JEM consistently outperforms Glow and IGEBM across all OOD metrics.

What works:

- Langevin dynamics + replay buffer is a practical MCMC recipe for images
- Spectral norm + energy regularization stabilize training
- Classifier logits provide a free energy parameterization
- Joint generative + discriminative training is feasible and beneficial

Open challenges:

- Training stability remains fragile (JEM required restarts with lower LR)
- Cannot compute normalized likelihoods: hard to verify learning
- MCMC mixing in high dimensions is slow (biased short-run chains)
- Scaling beyond CIFAR/small ImageNet still difficult

*These papers demonstrated that EBMs are not just theoretically elegant, they are **practically competitive** and bring unique capabilities (compositionality, robustness, calibration) that other generative models lack.*

Introduction & Motivation

Maximum Likelihood Training

Score Matching

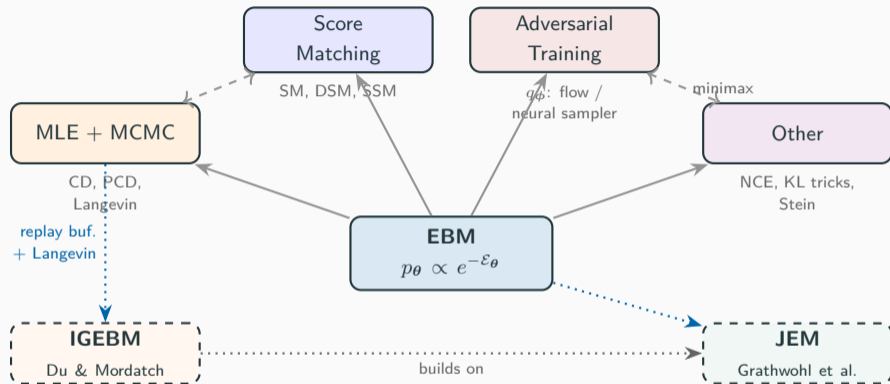
Adversarial Training

Other Training Methods

EBMs in Practice: Two Case Studies

Summary

The Big Picture



Theory: all methods sidestep intractable Z_θ .

Connections:
CD \leftrightarrow SM (de Bruijn),
adversarial \approx variational MLE.

Practice: IGEBM scales to images; JEM unifies gen. + disc. learning.

Takeaways

1. **EBMs are maximally flexible** generative models: \mathcal{E}_θ can be *any* scalar-valued network. The price is an intractable Z_θ .
2. **MLE + MCMC** gives unbiased gradients in principle, but requires sampling from the model: expensive and tricky (CD, PCD, replay buffers).
3. **Score matching** avoids Z_θ by matching *score functions*. DSM adds noise and learns to denoise; SSM projects onto random directions.
4. **Adversarial training** sidesteps MCMC entirely via a learned sampler q_ϕ (normalizing flow or noisy neural sampler), yielding a minimax game analogous to GANs but with an explicit energy model.
5. These methods are **deeply connected**: CD \leftrightarrow SM (de Bruijn), NCE \leftrightarrow SSM (local limit), adversarial training \leftrightarrow variational MLE.
6. **EBMs scale to real data** (IGEBM) and can **unify generative & discriminative learning** (IFM)

Thank you!

Questions?

Diffusion Models (Ch. 25)

Apr 15, 2026

Alberto Suárez

- Hinton, G.E. (2002). Training products of experts by minimizing contrastive divergence. *Neural Computation*.
- Hyvärinen, A. (2005). Estimation of non-normalized statistical models by score matching. *JMLR*.
- Hyvärinen, A. (2007). Some extensions of score matching. *Computational Statistics & Data Analysis*.
- Vincent, P. (2011). A connection between score matching and denoising autoencoders. *Neural Computation*.
- Gutmann, M.U., Hyvärinen, A. (2010). Noise-contrastive estimation. *AISTATS*.
- Song, Y., Garg, S., Shi, J., Ermon, S. (2019). Sliced score matching. *UAI*.
- Tieleman, T. (2008). Training restricted Boltzmann machines using approximations to the likelihood gradient. *ICML*.
- Du, Y., Mordatch, I. (2019). Implicit generation and modeling with energy-based models. *NeurIPS*.

- Grathwohl, W., Wang, K.-C., Jacobsen, J.-H., Duvenaud, D., Swersky, K., Norouzi, M. (2020). Your classifier is secretly an energy based model and you should treat it like one. *ICLR*.
- Grathwohl, W., Wang, K.-C., Jacobsen, J.-H., Duvenaud, D., Zemel, R. (2020). Learning the Stein discrepancy for training and evaluating energy-based models without sampling. *ICML*.
- Kim, T., Bengio, Y. (2016). Deep directed generative models with energy-based probability estimation. *arXiv:1606.03439*.
- Dai, Z., et al. (2019). Exponential family estimation via adversarial dynamics embedding. *NeurIPS*.
- Sohl-Dickstein, J., Battaglino, P., DeWeese, M.R. (2011). Minimum probability flow learning. *ICML*.
- Lyu, S. (2011). Unifying non-maximum-likelihood learning objectives with minimum KL contraction. *NeurIPS*.
- Murphy, K.P. (2023). *Probabilistic Machine Learning: Advanced Topics*. MIT Press.