A Unified View on Learning Abhin Shah **Unnormalized Distributions** via Noise-Contrastive Estimation MIT EECS

TL;DR: We unify various methods for learning EBMs via a family of NCE principles and establish finite-sample rates for exponential family distributions

Unnormalized Distributions (a.k.a. Energy-Based Models (EBMs)) **Are Flexible Probabilistic Models**



Application scenarios: graphical models (Markov random fields), physical modeling (Ising models), causal modeling, image modeling, ...

Variant 1. α -Centered NCE

• Consider $f_{\alpha}(\rho) = \frac{\rho^{\alpha} - 1}{\alpha(\alpha - 1)} (\alpha \notin \{0, 1\})$ ("asymmetric power")

• Define a **normalized** model (which we call " α -centered model")

$$\tilde{\phi}_{\theta;\alpha}(\mathbf{x}) \stackrel{\text{\tiny def}}{=} \frac{\phi_{\theta}(\mathbf{x})}{Z_{\alpha}(\theta)}, \quad \text{where } Z_{\alpha}(\theta) \stackrel{\text{\tiny def}}{=} \left(\mathbb{E}_{q_{n}(\mathbf{x})} \left[\left(\frac{\phi_{\theta}(\mathbf{x})}{q_{n}(\mathbf{x})} \right)^{\alpha} \right] \right)^{1/\alpha}$$

such that it's
$$\alpha$$
-centered: $\mathbb{E}_{q_n(\mathbf{x})} \left[\left(\frac{\tilde{\phi}_{\theta;\alpha}(\mathbf{x})}{q_n(\mathbf{x})} \right)^{\alpha} \right] = 1$
(α -CentNCE obj.) $\mathcal{L}_{\alpha}^{\text{cent}}(\phi_{\theta}) \stackrel{\text{def}}{=} \mathcal{L}_{f_{\alpha}}^{\text{nce}}(\tilde{\phi}_{\theta;\alpha}, 1)$

$$\begin{array}{c|c} \hline \text{Objectives} & \alpha = 0 & \alpha = \frac{1}{2} & \alpha = 1 \\ \hline f_{\alpha} \text{-NCE} & \begin{bmatrix} \mathbb{E}_{q_d} \left[\frac{q_n}{\phi_{\theta}} \right] + \mathbb{E}_{q_n} \left[\log \frac{\phi_{\theta}}{q_n} \right] & 2(\mathbb{E}_{q_d} \left[\sqrt{\frac{q_n}{\phi_{\theta}}} \right] + \mathbb{E}_{q_n} \left[\sqrt{\frac{\phi_{\theta}}{q_n}} \right]) & \mathbb{E}_{q_d} \left[\log \frac{q_n}{\phi_{\theta}} \right] + \mathbb{E}_{q_n} \left[\frac{\phi_{\theta}}{q_n} \right] \\ & (\text{InvIS} & (\text{eNCE} & (\text{Importance Sampling (IS)}) \\ & (\text{Pihlaja et al., 2010}) & (\text{Liu et al., 2021})) & (\text{Pihlaja et al., 2010; Riou-Durand & Chopin, 2018)}) \\ \hline \\ \alpha \text{-CentNCE} & \begin{bmatrix} \mathbb{E}_{q_d} \left[\frac{q_n}{\phi_{\theta}} \right] e^{\mathbb{E}_{q_n} \left[\log \frac{\phi_{\theta}}{q_n} \right]} \\ & (\mathbf{GlobalGISO} \\ & (\text{Shah et al., 2023)}) & 2\mathbb{E}_{q_d} \left[\sqrt{\frac{q_n}{\phi_{\theta}}} \right] \mathbb{E}_{q_n} \left[\sqrt{\frac{\phi_{\theta}}{q_n}} \right] \\ & MC-MLE (\text{Geyer, 1994; Jiang et al., 2023)}) \end{array}$$

Challenge: Unknown Normalization

Flexibility comes at the cost of not knowing the normalization constant (a.k.a. partition function)

$$Z_{\theta} \stackrel{\text{\tiny def}}{=} \int \exp(-E_{\theta})$$

Some known approaches

- Noise-contrastive estimation (Gutmann & Hyvarinen, 2010)
- Score matching (Hyvarinen, 2012)
- Contrastive divergence (Hinton, 2002)
- (pseudo likelihood (Besag, 1975), interaction screening (Vuffray et al., 2016), ...)

Q. Is there a unifying principle?

Takeaway 1

 α -CentNCE unifies MLE, MC-MLE (Geyer, 1994), and GlobalGISO $_{\rm (Shah\ et\ al.,\ 2023)}$ as limiting instances

Jongha (Jon) Ryu **Gregory W. Wornell**

correspondence to: jongha@mit.edu

 $d(\mathbf{x}) d\mathbf{x}$

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

 \Rightarrow Can't apply maximum likelihood estimation (MLE)!

Variants and other proposals in some specific settings

Variant 2. *f*-Conditional NCE

• $q_n(\mathbf{x})$, so "conditional NCE" (Ceylan and Gutmann, 2018) was proposed • Conditional distribution $\pi(\mathbf{y}|\mathbf{x})$ (ex: $\mathbf{y} = \mathbf{x} + \epsilon \mathbf{v}, \mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$) • Idea: match "joint" density ratios $\frac{q_{d}(\mathbf{x})\pi(\mathbf{y}|\mathbf{x})}{q_{d}(\mathbf{y})\pi(\mathbf{x}|\mathbf{y})} \approx \frac{\phi_{\theta}(\mathbf{x})\pi(\mathbf{y}|\mathbf{x})}{\phi_{\theta}(\mathbf{y})\pi(\mathbf{x}|\mathbf{y})}$ • We extend CondNCE to *f*-CondNCE in the same way as from NCE to *f*-NCE • For Gaussian $\pi(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{x}, \epsilon^2 \mathbf{I}))$, Ceylan and Gutmann (2018) argued $\lim_{\epsilon \to 0} (CondNCE) = (Score Matching)$ • **Takeaway 2**: However, we reveal that, for any convex *f*: (1) This is only true in the population limit; (2) With finite samples, *f*-CondNCE is dominated

by statistical noise as $\epsilon \to 0$ (so NOT equivalent to SM)



- Unnormalized model $\phi_{\theta}(\mathbf{x}) \stackrel{\text{\tiny def}}{=} \exp(-E_{\theta}(\mathbf{x}))$
- Data distribution $q_{d}(\mathbf{x})$
- Noise distribution $q_n(\mathbf{x})$ (typically Gaussian)

Note: Can be unbiasedly estimated by samples from $q_d(\mathbf{x})$ and $q_n(\mathbf{x})$ **Caveat**: Need to choose $q_n(\mathbf{x})$ carefully

We propose and study **TWO variants**

Takeaway 3. Finite-Sample Analysis

- considered estimators
- This is thanks to our unifying framework
- GlobalGISO (Shah et al., 2023)





• For exponential family $\phi_{\theta}(\mathbf{x}) = \exp(\theta^{\mathsf{T}}\psi(\mathbf{x}))$, we analyze the finite-sample performance of the estimators f-NCE, α -CentNCE, f-CondNCE • Most are **first** convergence rate guarantees for the • Idea: Adapt the convergence rate analysis of **Remark**: We can also extend this unification to

local models (e.g., sparse Markov random fields), which unifies pseudo likelihood (Besag, 1975), and ISO (Vuffray et al., 2016; 2021; Shah et al., 2021a; Ren et al., 2021)