A Tale of Two Flows: Cooperative Learning of Langevin Flow and Normalizing Flow Toward Energy-Based Model

Jianwen Xie, Yaxuan Zhu, Jun Li, Ping Li Cognitive Computing Lab, Baidu Research



(1) An EBM parameterized by a ConvNet can be trained by MCMC-based maximum likelihood estimation (Xie et al. $2016)^1$.

(2) However, the Langevin MCMC sampling on a highly multi-modal energy function is generally not mixing, which leads to a short-run Langevin flow model (Nijkamp et al. 2019)².

(3) Motivated by reducing the number of steps in the Langevin flow, this paper proposes the CoopFlow model that trains a Langevin flow jointly with a normalizing flow (Jonathan et al. 2019) ³ in a cooperative learning scheme (Xie et al. 2018)⁴.

¹Jianwen Xie, et al. "A Theory of Generative ConvNet." ICML, 2016. ²Erik Nijkamp, et al. "Learning Non-Convergent Non-Persistent Short-Run MCMC Toward Energy-Based Model." NeurIPS, 2019.

³Jonathan Ho, et al. "Flow++: Improving Flow-Based Generative Models with Variational Dequantization and Architecture Design." ICML, 2019

⁴ Jianwen Xie, et al. "Cooperative Training of Descriptor and Generator Networks." TPAMI, 2018

Energy-based model

Let $x \in \mathbb{R}^D$ be the observed signal, e.g., an image. An energy-based model defines an unnormalized probability distribution of x:

$$p_{\theta}(x) = \frac{1}{Z(\theta)} \exp[f_{\theta}(x)],$$

where $f_{\theta} : \mathbb{R}^{D} \to \mathbb{R}$ is the negative energy function and defined by a bottom-up neural network whose parameters are θ . The normalizing constant $Z(\theta) = \int \exp[f_{\theta}(x)] dx$ is analytically intractable and difficult to compute due to high dimensionality of x.

Maximum likelihood learning

Suppose we observe training examples $\{x_i, i = 1, ..., n\}$ from unknown data distribution $p_{data}(x)$, the energy-based model can be trained by maximum likelihood estimation. The learning gradient is given by

$$L'(\theta) = \mathbb{E}_{p_{data}}[\nabla_{\theta}f_{\theta}(x)] - \mathbb{E}_{p_{\theta}}[\nabla_{\theta}f_{\theta}(x)] \approx \frac{1}{n} \sum_{i=1}^{n} \nabla_{\theta}f_{\theta}(x_{i}) - \frac{1}{n} \sum_{i=1}^{n} \nabla_{\theta}f_{\theta}(\tilde{x}_{i}),$$

where the expectations are approximated by averaging over the observed examples $\{x_i\}$ and the synthesized examples $\{\tilde{x}_i\}$ generated from the current model $p_{\theta}(x)$ via MCMC, respectively.

Langevin dynamics as MCMC

Generating synthesized examples from $p_{\theta}(x)$ can be accomplished with Langevin dynamics, which is applied as follows

$$x_{t+1} = x_t + \frac{\delta^2}{2} \nabla_x f_\theta(x_t) + \delta \epsilon_t; \quad x_0 \sim p_0(x), \quad \epsilon_t \sim \mathcal{N}(0, I_D), \quad t = 1, \cdots, T,$$

where t indexes the Langevin time step, δ denotes the Langevin step size, ϵ_t is a Brownian motion that explores different modes, and $p_0(x)$ is a uniform distribution that initializes the MCMC chains.

Langevin Flow

Let $\tilde{p}_{\theta}(x)$ be the distribution of x_{T} , which is the resulting distribution of x after T steps of Langevin update starting from $x_0 \sim p_0(x)$. Due to fixed initial distribution $p_0(x)$ and fixed K and δ , the distribution $\tilde{p}_{\theta}(x)$ is

$$ilde{p}_{ heta}(x) = (\mathcal{K}_{ heta} p_0)(x) = \int p_0(z) \mathcal{K}_{ heta}(x|z) dz,$$

where \mathcal{K}_{θ} denotes the transition kernel of \mathcal{T} steps of Langevin dynamics that samples p_{θ} .

Generally, $\tilde{p}_{\theta}(x)$ is not necessarily equal to $p_{\theta}(x)$. $\tilde{p}_{\theta}(x)$ is dependent on T and s. According to the law of thermodynamics, $\mathbb{D}_{\mathsf{KL}}(\tilde{p}_{\theta}(x)||p_{\theta}(x))$ decreases to zero monotonically as $T \to \infty$.

Normalizing Flow

Let $z \in \mathbb{R}^D$ be the latent vector of the same dimensionality as x. A normalizing flow is of the form

$$x = g_{\alpha}(z); z \sim q_0(z),$$

where $q_0(z)$ is a Gaussian distribution, and g_α consists of a sequence of L invertible transformations.

The successive transformations between x and z can be expressed as a flow $z \xleftarrow{g_1} h_1 \xleftarrow{g_2} h_2 \cdots \xleftarrow{g_L} x$, where we define $z := h_0$ and $x := h_L$ for succinctness. The log-likelihood of x can be easily computed by

$$\log q_{lpha}(x) = \log q_0(z) + \sum_{l=1}^{L} \log \left| \det \left(\frac{\partial h_{l-1}}{\partial h_l} \right) \right|.$$

The flow-based model can be trained by maximizing the data log-likelihood $L(\alpha) = \sum_{i=1}^{n} \log q_{\alpha}(x_{i})$.

(1) We accept the fact that the short-run non-convergent MCMC is inevitable and more affordable in practice.

(2) We treat a non-convergent short-run Langevin flow as a generator

(3) We propose to jointly train the Langevin flow with a normalizing flow as a rapid initializer for more efficient generation.

(4) The resulting generator is called CoopFlow, which consists of a Langevin flow and a normalizing flow.

The distribution of the CoopFlow

The density of the CoopFlow $\pi_{(\theta,\alpha)}(x)$ can be implicitly expressed by

$$\pi_{(heta,lpha)}(x) = (\mathcal{K}_ heta q_lpha)(x) = \int q_lpha(x') \mathcal{K}_ heta(x|x') dx',$$

where θ is the parameters of the EBM, and α is the parameters of the normalizing flow. $\mathcal{K}_{\theta}(x|x')$ is the transition kernel of \mathcal{T} steps of Langevin dynamics that samples p_{θ}

The CoopFlow Algorithm

At each iteration, we perform

(Step 1) For i = 1, ..., m, we first generate $z_i \sim \mathcal{N}(0, I_D)$, and then transform z_i by a normalizing flow to obtain $\hat{x}_i = g_\alpha(z_i)$.

(Step 2) Starting from each \hat{x}_i , we run a Langevin flow (i.e., a finite number of Langevin steps toward an EBM $p_{\theta}(x)$) to obtain \tilde{x}_i .

(Step 3) We update α of the normalizing flow by treating \tilde{x}_i as training data.

(Step 4) We update θ of the Langevin flow according to the learning gradient of the EBM, which is computed with the synthesized examples \tilde{x}_i and the observed examples.

4. Understanding the Learned Two Flows

Moment Matching Estimator

Consider a simple EBM with $f_{\theta}(x) = \langle \theta, h(x) \rangle$, where h(x) is the feature statistics. The CoopFlow π^* is a moment matching estimator, i.e., $\mathbb{E}_{p_{data}}[h(x)] = \mathbb{E}_{\pi^*}[h(x)].$



Figure 1: Illustration of convergence of the CoopFlow.

5. Experiments

Exp 1: Image Generation



Figure 2: Generated examples (32×32 pixels) by CoopFlow models trained from CIFAR-10, SVHN and Celeba datasets respectively. Samples are obtained from the setting of CoopFlow(pre).

Exp 1: Image Generation

Approach	Models	FID ↓	Models	FID
VAE	VAE (Kingma & Welling, 2014)	78.41	ABP (Han et al., 2017)	49.71
Autoregressive	PixelCNN (Salimans et al., 2017) PixelIQN (Ostrovski et al., 2018)	65.93 49.46	ABP-SRI (Nijkamp et al., 2020b) ABP-OT (An et al., 2021)	35.23 19.48
GAN	WGAN-GP(Gulrajani et al., 2017) SN-GAN (Miyato et al., 2018) StyleGAN2-ADA (Karras et al., 2020)	36.40 21.70 2.92	VAE (Kingma & Welling, 2014) 2sVAE (Dai & Wipf, 2019) RAE (Ghosh et al., 2020) Glow (Kingma & Dhariwal, 2018) DCGAN (Radford et al., 2016) NT-EBM (Nijkamp et al., 2020a) LP-EBM (Pang et al., 2020)	46.78 42.81 40.02 41.70 21.40 48.01 29.44
Score Based	NCSN (Song & Ermon, 2019) NCSN-v2 (Song & Ermon, 2020) NCSN++ (Song et al., 2021)	25.32 31.75 2.20		
Flow	Glow (Kingma & Dhariwal, 2018) Residual Flow (Chen et al., 2019)	45.99 46.37	EBM-FCE (Gao et al., 2020)	20.19
EBM	LP-EBM (Pang et al., 2020) EBM-SR (Nijkamp et al., 2019) EBM-IG (Du & Mordatch, 2019) CoopVAEBM (Xie et al., 2021b)	70.15 44.50 38.20 36.20	CoopFlow(Long) CoopFlow(Pre) (b) SVHN Models	18.11 16.97 15.32
	CoopNets (Xie et al., 2020a)	33.61		
	Divergence Triangle (Han et al., 2020)	30.10		FID .
	BBM-CD (Du et al., 2021) GEBM (Arbel et al., 2021) GF-EBM Zhao et al. (2021) VAEBM (Xiao et al., 2021) EBM-Diffusion (Gao et al., 2021)	27.30 ABP (Han et al., 2017) 19.31 ABP-SRI (Nijkamp et al., 2 16.71 VAE (Kingma & Welling, 2 12.16 Glow (Kingma & Dhariwal, 9.60 DCGAN (Radford et al., 20	ABP (Han et al., 2017) ABP-SRI (Nijkamp et al., 2020b) VAE (Kingma & Welling, 2014) Glow (Kingma & Dhariwal, 2018) DCGAN (Radford et al., 2016)	51.50 36.84 38.76 23.32 12.50
Flow+EBM	NT-EBM (Nijkamp et al., 2020a) EBM-FCE (Gao et al., 2020)	78.12 37.30	EBM-FCE (Gao et al., 2020) GEBM (Arbel et al., 2021)	12.21 5.21
Ours	CoopFlow CoopFlow(Long) CoopFlow(Pre)	21.16 18.89 15.80	CoopFlow CoopFlow(Long) CoopFlow(Pre)	6.44 4.90 4.15
(a) CIFAR 10			(c) Celeba	

Exp 2: Image Reconstruction

The CoopFlow model is a latent variable generative model:

$$z \sim p_0(z); \hat{x} = g_\alpha(z); x = F_\theta(\hat{x}, e),$$

where z is the latent variables, e are all injected noises in the Langevin flow, and F_{θ} is the mapping realized by a *T*-step Langevin flow.

We can reconstruct any x by inferring the corresponding latent variables z using gradient descent on $L(z) = ||x - F_{\theta}(g_{\alpha}(z), e)||^2$, with z being initialized by p_0 .

5. Experiments

Exp 2: Image Reconstruction



Figure 3: Reconstruction of images.



Figure 4: Image inpainting. Each row represents one different initialization. The last two columns display the masked and original images respectively.

5. Experiments

Exp 3: Interpolation in the Latent Space

The CoopFlow is capable of doing interpolation in the latent space z. Given an image x, we first find its corresponding \hat{x}^* using the reconstruction method. We then infer z by the inversion of the normalizing flow $z^* = g_{\alpha}^{-1}(\hat{x}^*)$.



Figure 5: Image interpolation results on Celeba dataset (32×32) . The leftmost and rightmost columns represent the images we observed. The columns in the middle represent the interploation results between the inferred latent vectors of the observed images.

(1) We study amortized sampling for training a short-run non-mixing Langevin sampler toward an EBM.

(2) We propose a novel framework, the CoopFlow, which cooperatively trains a short-run Langevin flow as a valid generator and a normalizing flow as an amortized sampler for image representation.

(3) We provide both theoretical and empirical justifications for the proposed CoopFlow algorithm.