# EBM Learning and Inference: A Retrospection and Beyond

Zijing Ou ouzj@mail2.sysu.edu.cn

# My Story on EBM

First meet EBMs



#### **Binary representation learning on text**

• Generative hashing

$$p_{\theta}(x,z) = p_{\theta}(x|z)p(z)$$

$$x = \{w_1, w_2, \dots, w_{|x|}\} \qquad z = \{0,1\}^m$$
One-hot representation Binary codes

• Softmax decoder

$$p_{\theta}(w_i|z) = \frac{\exp(z^T E w_i + b_i)}{\sum_{j=1}^{|V|} \exp(z^T E w_j + b_j)}$$

 $p_{\theta}(x|z) = \prod_{i=1}^{|x|} p_{\theta}(w_i|z)$  (iid. assumption)

#### Variational inference

• Evidence lower bound

$$\log p(x) = ELBO + KL(q_{\phi}(z|x)||p(z))$$
$$\geq ELBO = E_{q_{\phi}(z|x)} \left[ \log \frac{p_{\theta}(x|z)p(z)}{q_{\phi}(z|x)} \right]$$

• Consistent learning & inference

$$\mathcal{L} = E_{p(x)} E_{q_{\phi}(z|x)} \left[ \log \frac{p_{\theta}(x|z)p(z)}{q_{\phi}(z|x)} \right]$$
 Independent among data

 $q_{\phi}(z|x) \coloneqq \mathcal{N}(z|\mu_{\phi}(x), diag(\sigma_{\phi}^2(x)))$  Independent among bits

#### Our focus: breaking the independence assumption!!!

#### **Boltzmann machine as posterior**



Lower rank perturbation

$$b(z) = \frac{1}{Z} exp(\frac{1}{2}z^T\Sigma z + \mu^T z)$$

 $\Sigma = D + UU^T$  (diagonal + low-rand perturbation matrix)

• Reparameterization  $b(z) = \int p(z|r)p(r)dr$   $p(r) = \frac{1}{z} \prod_{i=1}^{m} (e^{r_i} + 1)\mathcal{N}(r;\mu,\Sigma)$   $p(z|r) = \prod_{i=1}^{m} Bernoulli(\sigma(r_i))$ 

$$q_{\phi}(z|x) = E_{q_{\phi}(r|x)}[Bernoulli(z;r)]$$



#### Boltzmann machine as variational posterior

Zheng and Su. Generative Semantic Hashing Enhanced via Boltzmann Machines. ACL 2020

### **Cancel out partition function**

• Evidence lower bound  $\mathcal{L} = E_{q_{\phi}(Z|X)} \left[ \log \frac{p_{\theta}(x|z)p(z)}{e^{-E_{\phi}(z)}} \right] + \log Z_{\phi}$   $E_{\phi}(z) \coloneqq -\frac{1}{2} z^{T} \Sigma_{\phi}(x) z - \mu_{\sigma}^{T}(x) z$ • An asymptotically-exact lower bound  $h_{k}(z) \coloneqq \frac{1}{k} \sum_{k=1}^{k} p(z|r^{(i)})$ 

$$\begin{split} \widetilde{\mathcal{L}_{k}} &= \mathcal{L} - E_{q_{\phi}\left(r^{(1\dots k)}|x\right)}[KL(h_{k}(z)||q_{\phi}(z|x))] \\ &= E_{q_{\phi}(z|x)}[\log p_{\theta}(x,z)] - E_{q_{\phi}\left(r^{(1\dots k)}|x\right)}E_{h_{k}(z)}[\log h_{k}(z)] \\ &\widetilde{\mathcal{L}_{k}} < \widetilde{\mathcal{L}_{k+1}} \quad \lim_{k \to \infty} \widetilde{\mathcal{L}_{k}} = \mathcal{L} \end{split}$$

Ranganath et al. Hierarchical Variational Models. ICML 2016 Yin and Zhou. Semi-Implicit Variational Inference. ICML 2018 Zheng and Su. Generative Semantic Hashing Enhanced via Boltzmann Machines. ACL 2020



#### Data dependence prior



## **Spanning-tree** approximations



$$p_{\mathrm{T}}(Z) = \prod_{i \in \mathcal{V}} p_G(z_i) \prod_{(i,j) \in \mathcal{E}} \frac{p_G(z_i, z_j)}{p_G(z_i) p_G(z_j)}$$

$$q_{\mathrm{T}}(Z|X) = \prod_{i \in \mathcal{V}} q_{\phi}(z_i|x_i) \prod_{(i,j) \in \mathcal{E}} \frac{q_{\phi}(z_i, z_j|x_i, x_j)}{q_{\phi}(z_i|x_i)q_G(z_j|x_j)}$$

Variance of  $q_{\phi}(z_i, z_j | x_i, x_j)$ :

 $\gamma_{ij} \in (0, 1)$ : positive correlated

$$\begin{bmatrix} diag(\sigma_i^2) & diag(\gamma_{ij} \odot \sigma_i \odot \sigma_j) \\ diag(\gamma_{ij} \odot \sigma_i \odot \sigma_j) & diag(\sigma_j^2) \end{bmatrix}$$

#### **Spanning-tree** approximations



$$\mathcal{L}_{T} = \sum_{i \in \mathcal{V}} E_{q_{\phi}}[\log p_{\theta}(x_{i}|z_{i})] - KL(q_{\phi}(z_{i})||p_{G}(z_{i}))$$

$$= \sum_{(i,j)\in\mathcal{E}} \left( KL\left(q_{\phi}(z_{i},z_{j})||p_{G}(z_{i},z_{j})\right) - KL\left(q_{\phi}(z_{i})||p_{G}(z_{i})\right) - KL\left(q_{\phi}(z_{j})||p_{G}(z_{j})\right) \right)$$

#### **Extending to multiple trees**



 $T_G = \{T_1, T_2, T_3\}$ 

 $T_1$ 

 $T_2$ 

 $T_3$ 

 $\chi_2$ 

## Looking forward



Structured representation learning

- Expressive prior & posterior
- Combine with advances in GNNs
- Graphical models meet Deep Learning (GNN as message passing)

Loss

• Fitting student-t with score matching

$$F(\mathbf{p}, q) = \frac{1}{2} E_{\mathbf{p}}[\|\nabla_x \log q(x) - \nabla_x \log \mathbf{p}(x)\|^2]$$

$$p_{data}(x) = T_{v=5}(0, 0.3)$$
  $q_{\theta}(x) = T_{v=5}(\theta, 0.3)$ 



With diffusion matrix m(x)**Fisher Divergence Diffusion Score Matching**  $F_m(\mathbf{p}, \mathbf{q}) = \frac{1}{2} E_{\mathbf{p}}[\|\mathbf{m}(\mathbf{x})^{\mathrm{T}}(\nabla_x \log \mathbf{q}(x) - \nabla_x \log \mathbf{p}(x))\|^2]$  $F(p,q) = \frac{1}{2} E_p[\|\nabla_x \log q(x) - \nabla_x \log p(x)\|^2]$ SM and DSM Loss for Student-t distribution at different  $\theta$ 20 DSM with manual flow Manual flow  $\theta$ 15 Score Matching Fast region (Manual flow) Ground Truth 10 5 0  $(\boldsymbol{x}- heta)^2$ -5  $\boldsymbol{m}(\boldsymbol{x})$ 

0

Mean  $\theta$ 

2

4

Barp et al. Minimum Stein Discrepancy Estimators. NIPS 2019 Gong and Li. Interpreting diffusion score matching using normalizing flow. ICML 2021 INNF+ workshop

-4

-2

-10

• Interpreting DSM using flows

$$p_{X}, q_{X} \xrightarrow{\text{invertible flow Y} = T(X)} p_{Y}, q_{Y}$$

$$m(x) = (\nabla_{x}T(X))^{-1} \xrightarrow{p_{Y}, q_{Y}} q_{Y}(y) = p_{X}(T^{-1}(y)) \left| \det \frac{\partial T^{-1}(y)}{\partial y} \right|$$

$$q_{Y}(y) = q_{X}(T^{-1}(y)) \left| \det \frac{\partial T^{-1}(y)}{\partial y} \right|$$
Fisher Divergence
$$F(p_{X}, q_{X}) \xrightarrow{m(x) = (\nabla_{x}T(X))^{-1}} \xrightarrow{\text{Diffusion Score Matching}} F(p_{Y}, q_{Y}) = F_{m}(p_{X}, q_{X})$$

• Interpreting DSM using flows

$$p_{X}, q_{X} \xrightarrow{\text{invertible flow Y} = T(X)} p_{Y}, q_{Y}$$

$$p_{Y}(y) = p_{X}(T^{-1}(y)) \left| det \frac{\partial T^{-1}(y)}{\partial y} \right| \qquad q_{Y}(y) = q_{X}(T^{-1}(y)) \left| det \frac{\partial T^{-1}(y)}{\partial y} \right|$$

$$F(p_{Y}, q_{Y}) \coloneqq \frac{1}{2} E_{q_{Y}}[\|\nabla_{Y} \log p_{Y}(y) - \nabla_{y} \log q_{Y}(y)\|^{2}]$$

$$= \frac{1}{2} E_{q_{Y}}\left[ \|\nabla_{y} \log p_{X}(T^{-1}(y)) - \nabla_{y} \log q_{X}(T^{-1}(y)) \|^{2} \right]$$

$$= \frac{1}{2} E_{q_{Y}}\left[ \|\nabla_{y} T^{-1}(y)^{T}(\nabla_{T^{-1}(y)} \log p_{X}(T^{-1}(y)) - \nabla_{T^{-1}(y)} \log q_{X}(T^{-1}(y)) \|^{2} \right]$$

$$= \frac{1}{2} E_{q_{X}}[\|\nabla_{x} T(x)^{-T}(\nabla_{x} \log p_{X}(x) - \nabla_{x} \log q_{X}(x))\|^{2}]$$

• Interpreting DSM using flows

$$p_{X}, q_{X} \xrightarrow{\text{invertible flow Y} = T(X)} p_{Y}, q_{Y}$$

$$p_{Y}(y) = p_{X}(T^{-1}(y)) \left| det \frac{\partial T^{-1}(y)}{\partial y} \right| \qquad q_{Y}(y) = q_{X}(T^{-1}(y)) \left| det \frac{\partial T^{-1}(y)}{\partial y} \right|$$

$$F(p_{Y}, q_{Y}) = \frac{1}{2} \mathbb{E}_{q_{Y}} [ \|\nabla_{Y} \log p_{Y}(y) - \nabla_{y} \log q_{Y}(y)\|^{2} ]$$

$$= \frac{1}{2} E_{q_{Y}} \left[ \|\nabla_{y} \log p_{X}(T^{-1}(y)) - \nabla_{y} \log q_{X}(T^{-1}(y)) \|^{2} \right]$$

$$= \frac{1}{2} E_{q_{Y}} \left[ \|\nabla_{y} T^{-1}(y)^{T} (\nabla_{T^{-1}(y)} \log p_{X}(T^{-1}(y)) - \nabla_{T^{-1}(y)} \log q_{X}(T^{-1}(y)) \|^{2} \right]$$

$$= \frac{1}{2} E_{q_{X}} \left[ \|\nabla_{x} T(x)^{-T} (\nabla_{x} \log p_{X}(x) - \nabla_{x} \log q_{X}(x)) \|^{2} \right]$$

• Interpreting DSM using flows

$$p_{X}, q_{X} \xrightarrow{\text{invertible flow Y} = T(X)} p_{Y}, q_{Y}$$

$$p_{Y}(y) = p_{X}(T^{-1}(y)) \left| det \frac{\partial T^{-1}(y)}{\partial y} \right| \qquad q_{Y}(y) = q_{X}(T^{-1}(y)) \left| det \frac{\partial T^{-1}(y)}{\partial y} \right|$$

$$F(p_{Y}, q_{Y}) \coloneqq \frac{1}{2} \mathbb{E}_{q_{Y}} [ \|\nabla_{Y} \log p_{Y}(y) - \nabla_{Y} \log q_{Y}(y)\|^{2} ]$$

$$= \frac{1}{2} E_{q_{Y}} \left[ \|\nabla_{Y} \log p_{X}(T^{-1}(y)) - \nabla_{Y} \log q_{X}(T^{-1}(y)) \|^{2} \right]$$

$$= \frac{1}{2} E_{q_{Y}} \left[ \|\nabla_{Y} T^{-1}(y)^{T} (\nabla_{T^{-1}(y)} \log p_{X}(T^{-1}(y)) - \nabla_{T^{-1}(y)} \log q_{X}(T^{-1}(y)) \|^{2} \right]$$

$$= \frac{1}{2} E_{q_{X}} [ \|\nabla_{x} T(x)^{-T} (\nabla_{x} \log p_{X}(x) - \nabla_{x} \log q_{X}(x)) \|^{2} ]$$

SM and DSM Loss for Student-t distribution at different 
$$\theta$$
  

$$p_{x}, q_{x} \xrightarrow{\text{invertible flow Y = T(X)}{m(x) = (\nabla_{x}T(X))^{-1}} p_{Y}, q_{Y}$$

$$p_{data}(x) = T_{v=5}(0, 0.3) \quad q_{\theta}(x) = T_{v=5}(\theta, 0.3)$$

$$p_{m}(p, q) = \frac{1}{2}E_{p} \left[ \left\| m(x)^{T}(\nabla_{x} \log q(x) - \nabla_{x} \log p(x)) \right\|^{2} \right]_{5}$$

$$DSM: m(x) = \left( 1 + \frac{(x-\theta)^{2}}{0.6} \right)$$

$$Gaussian flow: T_{G}(x) = F_{G}^{-1} \circ F_{\theta}(x)$$

$$-10$$

$$SM \text{ and DSM Loss for Student-t distribution at different  $\theta$ 

$$Gaussian flow \theta$$

$$Gaussian flow \theta$$

$$Score Matching Fast region (Gaussian flow)$$

$$Gound Tuth$$

$$m(x) = (\nabla_{x}T_{G}(x))^{-1}$$

$$m(x) = (\nabla_{x}T_{G}(x))^{-1}$$

$$-10$$

$$-4$$

$$-2$$

$$0$$

$$2$$

$$4$$$$



 $p_{\theta}(x) = \frac{exp(f(x))}{Z}, x \in \{0,1\}^{D}, Z = \int f(x)dx$ 

• Metropolis-Hastings sampling

Proposal distribution  $q(x_{-i}|x)$ :

sample index  $i \sim q(i|x)$ , set  $x_{-i} = flip_{-}dim(x, i)$ 

Acceptance rate:

$$\min\left\{\exp(f(x_{-i}) - f(x))\frac{q(i|x_{-i})}{q(i|x)}, 1\right\}$$

How to design q(i|x)?

$$p_{\theta}(x) = \frac{exp(f(x))}{Z}, x \in \{0,1\}^{D}, Z = \int f(x) dx$$

• Metropolis-Hastings sampling

Acceptance rate:

$$\min\left\{\exp(f(x_{-i}) - f(x))\frac{q(i|x_{-i})}{q(i|x)}, 1\right\}$$
  
nigh to proposals have high likelihood

computational complexity



$$q(i|x) \propto \exp(f(x_{-i}) - f(x))$$

$$p_{\theta}(x) = \frac{exp(f(x))}{Z}, x \in \{0,1\}^{D}, \overline{Z} = \int f(x)dx$$

• Metropolis-Hastings sampling

Proposal distribution

computational complexity

**U**(1

 $q(i|x) \propto exp(f(x_{-i}) - f(x))$  **O(D)** 

First order Taylor approximation  

$$1 if x_i = 0; -1 if x_i = 1$$

$$f(x_{-i}) - f(x) \approx (x_{-i} - x)^T \nabla_x f(x)$$

$$= (-(2x - 1) \odot \nabla_x f(x))_i$$
computational complexity

$$q(i|x) \propto \exp\left(\left(-(2x-1)\odot\nabla_x f(x)\right)_i\right)$$

Grathwohl et al. Oops I Took A Gradient: Scalable Sampling for Discrete Distributions. ICML 2021



Grathwohl et al. Oops I Took A Gradient: Scalable Sampling for Discrete Distributions. ICML 2021

### Applying gradients to ratio matching

• Ratio matching

computational complexity

$$F_R(p;q) = E_p \left[ \sum_i \left( \log \frac{q(x_{-i})}{q(x)} - \log \frac{p(x_{-i})}{p(x)} \right)^2 \right] \quad \boldsymbol{O}(\boldsymbol{D})$$

• Continuous gradients relaxation

$$\log p(x_{-i}) - \log p(x) \propto (x_{-i} - x)^T \nabla_x \log p(x)$$
$$= (-(2x - 1) \odot \nabla_x \log p(x))_i$$

computational complexity

$$F_{RL}(p;q) = E_p \left[ \left\| -(2x-1)^T (\nabla_x \log q(x) - \nabla_x \log p(x)) \right\|^2 \right] \quad O(1)$$

Diffusion/flow score matching:  $m(x) = \nabla_x T(x)^{-1} = -(2x - 1)$ 

cancel out partition function

## Applying gradients to ratio matching

$$p(x) \propto \exp(\frac{1}{2}x^T \Sigma x + \mu^T x) \qquad x = \{0, 1\}^m$$
$$q_{\theta}(x) \propto \exp\left(\frac{1}{2}x^T \Sigma_{\theta} x + \mu_{\theta}^T x\right)$$

• Ratio matching

$$\mathcal{L} = E_p\left[\sum_{i=1}^m f^2\left(\frac{q(x)}{q(x_{-i})}\right)\right], f(x) = \frac{1}{1+x^2}$$

• General score matching



• Gradient ratio matching

 $\mathcal{L} = E_p[\| - (2x - 1)^T (\nabla_x \log q(x) - \nabla_x \log p(x)) \|^2]$ 



## Looking forward



- EBM learning
  - Efficient and accurate gradient estimation (slice, diffusion, stein method)
  - Discrete variables (relaxed by gradients, enhanced by flow)
- EBM inference
  - Sampling from EBMs (graph generation)
  - Latent variables (representation learning)