Generative Models



Desiderata for generative models

• **Probability evaluation**: given a sample, it is computationally efficient to evaluate the probability of this sample.

• Flexible model family: it is easy to incorporate any neural network models.

• **Easy sampling:** it is computationally efficient to sample a data from the probabilistic model.

Desiderata for generative models



Slide credit to Yang Song

Taxonomy of generative models



- (Nonlinear) ICA
- Normalizing flows

Image credits to Andrej Risteski

Key challenge for building generative models



Slide credit to Yang Song

Slide credit to Yang Song

Key challenge for building generative models

Approximating the normalizing constant

- Variational auto-encoders [Kingma & Welling 2014, Rezende et al. 2014]
- Energy-based models [Ackley et al. 1985, LeCun et al. 2006]

Using restricted neural network models

- Autoregressive models [Bengio & Bengio 2000, van den Oord et al. 2016]
- Normalizing flow models [Dinh et al. 2014, Rezende & Mohamed 2015]

Generative adversarial networks (GANs)

• Model the generation process, not the probability distribution [Goodfellow et al. 2014]







Variational Autoencoder



Architecture

- Auto-encoder: $x \rightarrow z \rightarrow x$
- Encoder: $q(z | x; \phi) : x \to z$
- Decoder: $p(x | z; \theta) : z \to x$

- Isomorphic Gaussian:
- $q(z \mid x; \phi) = N(\mu(x; \phi), \operatorname{diag}(\exp(\sigma(x; \phi))))$
- Gaussian prior: p(z) = N(0,I)
- Gaussian likelihood: $p(x | z; \theta) \sim N(f(z; \theta), I)$
- Probabilistic model interpretation: latent variable model.



VAE Training

- Training via optimizing ELBO
 - $L(\phi, \theta; x) = \mathbb{E}_{z \sim q(z|x;\phi)}[\log p(z|x;\theta)] KL(q(z|x;\phi)||p(z))$
 - Likelihood term + KL penalty



- Likelihood term (reconstruction loss):
 - Monte-Carlo estimation
 - Draw samples from $q(z | x; \phi)$
 - Compute gradient of θ :





VAE Training

• Likelihood term (reconstruction loss):

- Gradient for ϕ . Loss: $L(\phi) = \mathbb{E}_{z \sim q(z;\phi)} \left[\log p(x \mid z) \right]$
- Reparameterization trick:

•
$$z \sim N(\mu, \Sigma) \Leftrightarrow z = \mu + \epsilon, \epsilon \sim N(0, \Sigma)$$

• $L(\phi) \propto \mathbb{E}_{z \sim q(z|\phi)} \left[\|f(z;\theta) - x\|_2^2 \right]$ $\propto \mathbb{E}_{\epsilon \sim N(0,I)} \left[\|f(\mu(x;\phi) + \sigma(x;\phi) \cdot \epsilon;\theta) - x\|_2^2 \right]$

• Monte-Carlo estimate for $\nabla L(\phi)$

• End-to-end training



VAE vs. AE

- AE: classical unsupervised representation learning method.
- VAR: a probabilistic model of AE
 - AE + Gaussian noise on z
 - KL penalty: L_2 constraint on the latent vector z



Conditioned VAE

• Semi-supervised learning: some labels are also available



conditioned generation

Comments on VAE

- Pros:
 - Flexible architecture
 - Stable training
- Cons:
 - Inaccurate probability evaluation (approximate inference)

Energy-Based Models



Energy-based Models

- Goal of generative models:
 - a probability distribution of data: P(x)
- Requirements
 - $P(x) \ge 0$ (non-negative) • $\int_{x} P(x)dx = 1$
- Energy-based model:
 - Energy function: $E(x; \theta)$, parameterized by θ

•
$$P(x) = \frac{1}{z} \exp(-E(x;\theta))$$
 (why exp?)
• $z = \int_{z} \exp(-E(x;\theta)) dx$

Boltzmann Machine

• Generative model

•
$$E(y) = -\frac{1}{2}y^{\mathsf{T}}Wy$$

• $P(y) = \frac{1}{z}\exp(-\frac{E(y)}{T})$, T: temperature hyper-parameter

- W: parameter to learn
- When y_i is binary, patterns are affecting each other through W



$$z_i = \frac{1}{T} \sum_j w_{ji} s_j$$

$$P(s_i = 1 | s_{j \neq i}) = \frac{1}{1 + e^{-z_i}}$$

Boltzmann Machine: Training

- Objective: maximum likelihood learning (assume T =1):
 - Probability of one sample:

$$P(y) = \frac{\exp(\frac{1}{2}y^{\top}y)}{\sum_{y'} \exp(y'^{\top}Wy')}$$

• Maximum log-likelihood:

$$L(W) = \frac{1}{N} \sum_{y \in D} \frac{1}{2} y^{\mathsf{T}} W y - \log \sum_{y'} \exp(\frac{1}{2} y'^{\mathsf{T}} W y')$$

Boltzmann Machine: Training

Boltzmann Machine: Training

Restricted Bolzmann Machine

- A structured Boltzmann Machine
 - Hidden neurons are only connected to visible neurons
 - No intra-layer connections
 - Invented by Paul Smolensky in '89
 - Became more practical after Hinton invested fast learning algorithms in mid 2000



Restricted Bolzmann Machine

- Computation Rules
 - Iterative sampling

• Hidden neurons
$$h_i: z_i = \sum_j w_{ij} v_j$$
, $P(h_i | v) = \frac{1}{1 + \exp(-z_i)}$
• Visible neurons $v_j: z_j = \sum_i w_{ij} h_i$, $P(v_j | h) = \frac{1}{1 + \exp(-z_j)}$



Restricted Bolzmann Machine

- Sampling:
 - Randomly initialize visible neurons v₀
 - Iterative sampling between hidden neurons and visible neurons
 - Get final sample (v_{∞}, h_{∞})
- Training:
 - MLE
 - Sampling to approximate gradient



Deep Bolzmann Machine

- Can we have a **deep** version of RBM?
 - Deep Belief Net ('06)
 - Deep Boltzmann Machine ('09)
- Sampling?
 - Forward pass: bottom-up
 - Backward pass: top-down
- Deep Bolzmann Machine
 - The very first deep generative model
 - Salakhudinov & Hinton



Deep Bolzmann Machine

Deep Boltzmann Machine



Gaussian visible units (raw pixel data)

Training Samples									
and the second	A	B	JAT.		营				
\$	No.	X	1	E C	B				
		(m)		*	1				
Ŕ	E.	6	-ange	St.	- Mail				
A	0	X	1	E.	(B)				
السفل	ALL STREET	Ť	A.		54 a				

Generated Samples

×	2	80C	U	258	K
南	Z	\$	营	净	1
U	1	3	8	2	t
X	to	٤	N.	10	1
Ģ	1	あい	1	a	父
Ser.	5	乱	X	X	於

Summary

- Pros: powerful and flexible
 - An arbitrarily complex density function $p(x) = \frac{1}{z} \exp(-E(x))$
- Cons: hard to sample / train
 - Hard to sample:
 - MCMC sampling
 - Partition function
 - No closed-form calculation for likelihood
 - Cannot optimize MLE loss exactly
 - MCMC sampling

Normalizing Flows



Intuition about easy to sample

- Goal: design p(x) such that
 - Easy to sample
 - Tractable likelihood (density function)
- Easy to sample
 - Assume a continuous variable z
 - e.g., Gaussian $z \sim N(0,1)$, or uniform $z \sim \text{Unif}[0,1]$
 - x = f(z), x is also easy to sample

Intuition about tractable density

- Goal: design $f(z; \theta)$ such that
 - Assume *z* is from an "easy" distribution
 - $p(x) = p(f(z; \theta))$ has tractable likelihood
- Uniform: $z \sim \text{Unif}[0,1]$
 - Density p(z) = 1
 - x = 2z + 1, then p(x) = ?



Intuition about tractable density



- Assume *z* is from an "easy" distribution
- $p(x) = p(f(z; \theta))$ has tractable likelihood
- Uniform: $z \sim \text{Unif}[0,1]$
 - Density p(z) = 1

•
$$x = 2z + 1$$
, then $p(x) = 1/2$

• x = az + b, then p(x) = 1/|a| (for $a \neq 0$)

•
$$x = f(z)$$
, then $p(x) = p(z) \left| \frac{dz}{dx} \right| = |f'(z)|^{-1} p(z)$

• Assume f(z) is a bijection



Change of variable

• Suppose x = f(z) for some general non-linear $f(\cdot)$

• The linearized change in volume is determined by the Jacobian of $f(\ \cdot\)$:

$$\frac{\partial f(z)}{\partial z} = \begin{vmatrix} \frac{\partial f_1(x)}{\partial z_1} & \cdots & \frac{\partial f_1(z)}{\partial z_d} \\ \vdots & \cdots & \cdots & \vdots \\ \frac{\partial f_d(z)}{\partial z_1} & \cdots & \frac{\partial f_d(z)}{\partial z_d} \end{vmatrix}$$

Given a bijection $f(z) : \mathbb{R}^d \to \mathbb{R}^d$
• $z = f^{-1}(x)$
• $p(x) = p(f^{-1}(x)) \left| \det\left(\frac{\partial f^{-1}(x)}{\partial x}\right) \right| = p(z) \left| \det\left(\frac{\partial f^{-1}(x)}{\partial x}\right) \right|$
• Since $\frac{\partial f^{-1}}{\partial x} = \left(\frac{\partial f}{\partial x}\right)^{-1}$ (Jacobian of invertible function)
• $p(x) = p(z) \left| \det\left(\frac{\partial f^{-1}(x)}{\partial x}\right) \right| = p(z) \left| \det\left(\frac{\partial f(z)}{\partial z}\right) \right|^{-1}$

Normalizing Flow

- Idea
 - Sample z_0 from an "easy" distribution, e.g., standard Gaussian
 - Apply *K* bijections $z_i = f_i(z_{i-1})$
 - The final sample $x = f_K(z_K)$ has tractable desnity
- Normalizing Flow
 - $z_0 \sim N(0,I), z_i = f_i(z_{i-1}), x = Z_K$ where $x, z_i \in \mathbb{R}^d$ and f_i is invertible
 - Every revertible function produces a normalized density function



Normalizing Flow

- Generation is trivial
 - Sample z_0 then apply the transformations
- Log-likelihood

•
$$\log p(x) = \log p(Z_{k-1}) - \log \left| \det \left(\frac{\partial f_K}{\partial z_{K-1}} \right) \right|$$

• $\log p(x) = \log p(z_0) - \sum_i \log \left| \det \left(\frac{\partial f_i}{\partial z_{i-1}} \right) \right|$ **O** (d^3) !!!!



Normalizing Flow

- Naive flow model requires extremely expensive computation
 - Computing determinant of $d \times d$ matrices
- Idea:
 - Design a good bijection $f_i(z)$ such that the determinant is easy to compute

Plannar Flow

- Technical tool: Matrix Determinant Lemma:
 - $\det(A + uv^{\top}) = (1 + v^{\top}A^{-1}u) \det A$
- Model:
 - $f_{\theta}(z) = z + u \odot h(w^{\top}z + b)$
 - $h(\cdot)$ chosen to be $tanh(\cdot)(0 < h'(\cdot) < 1)$

•
$$\theta = [u, w, b], \det\left(\frac{\partial f}{\partial z}\right) = \det(I + h'(w^{\mathsf{T}}z + b)uw^{\mathsf{T}}) = 1 + h'(w^{\mathsf{T}}z + b)u^{\mathsf{T}}w$$

- Computation in O(d) time
- Remarks:
 - $u^{\top}w > -1$ to ensure invertibility
 - Require normalization on u and w

Planar Flow (Rezende & Mohamed, '16)

- $f_{\theta}(z) = z + uh\left(w^{\mathsf{T}}z + b\right)$
- 10 planar transformations can transform simple distributions into a more complex one



Extensions

- Other flow models uses triangular Jacobian (NICE, Dinh et al. '14)
- Invertible 1x1 convolutions (Kingma et al. '18)
- Auto-regressive flow:
 - WaveNet (Deepmind '16)
 - PixelCNN (Deepmind '16)

Summary

- Pros:
 - Easy to sample by transforming from a simple distribution
 - Easy to evaluate the probability
 - Easy training (MLE)
- Con
 - Most restricted neural network structure
 - Trade expressiveness for tractability