高等机器学习

# 生成式模型

刘 畅
微软亚洲研究院

Microsoft

清华大学
Tsinghua University

# Generative Model: Overview

- Generative Models:
  Models that define $p(\text{data})$: $\textcolor{red}{p(x)\ (\text{unsupervised})}$ or $\textcolor{red}{p(x, y)\ (\text{supervised})}$.
  - By computing the p.d.f/p.m.f of $p(\text{data})$: data generation can be done in principle.
  - By specifying a generating process of data: the distribution $p(\text{data})$ is implicitly defined.

Unsupervised:

$$\{x^{(1)}, \ldots, x^{(N)}\} = \left\{ \boxed{2}, \boxed{7}, \boxed{1}, \boxed{5}, \ldots, \boxed{0} \right\} \sim p(x)$$

Supervised:

$$\{(x^{(1)}, y^{(1)}), \ldots, (x^{(N)}, y^{(N)})\} = \left\{ (\boxed{2}, "2"), \ldots, (\boxed{7}, "7") \right\} \sim p(x, y)$$

# Generative Model: Overview

- Non-Generative Models:

Discriminative models
(e.g., feedforward neural networks):
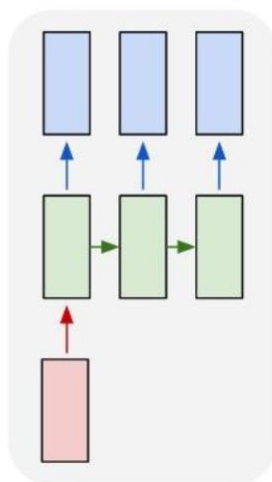only $p(y|x)$ is available.

$x$ 

$f$

$p(y|x)$



"0" "1" "2" "3" "4" "5" "6" "7" "8" "9"

Recurrent neural networks:

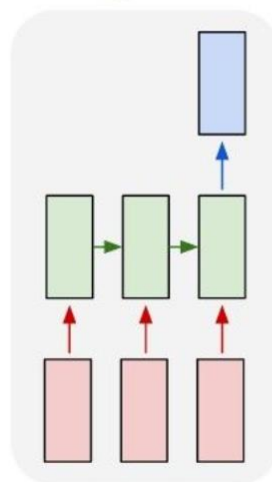only $p(\ \boxed{\ }\ |\ \boxed{\ }\ )$ is

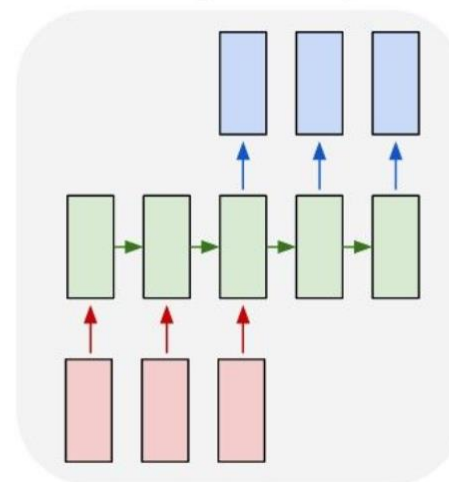available.



one to one    one to many    many to one    many to many    many to many

# Generative Model: Overview

- What can generative models do:

  1. Generate new data.



Generation $p(x)$ [KW14]



Conditional Generation
$p(x|y)$ [LWZZ18]

# Generative Model: Overview

- What can generative models do:

1. Generate new data.

*"the cat sat on the mat"* $\sim p(x)$: Language Model.



$p(x_1 = \text{the}) \quad p(\text{cat}|x_1) \quad p(\text{sat}|x_{1\ldots2}) \quad p(\text{on}|x_{1\ldots3}) \quad p(\text{the}|x_{1\ldots4}) \quad p(\text{mat}|x_{1\ldots5}) \quad p(</s>|x_{1\ldots6})$

# Generative Model: Overview

- What can generative models do:

  2. Infer unobserved variables.



Did it *Rain* if we see *GrassWet*?
-- Query $p(R|G = 1)$ from $p(S, R, G)$.



Missing Value Imputation (Completion) [OKK16].
-- Query $p(x_{\text{hidden}}|x_{\text{observed}})$ from $p(x_{\text{hidden}}, x_{\text{observed}})$.

# Generative Model: Overview

- What can generative models do:

  3. Density estimation $p(x)$.

    - Uncertainty estimate.

    - Anomaly detection.



[Ritchie Ng]

# Generative Model: Overview

- What can generative models do:

4. Representation learning: semantic and concise (via latent variable $z$).



$x$ (documents)

| "ENGINES" | speed | product | introduced | designs |
|-----------|-------|---------|-----------|---------|
| "ROYAL" | britain | queen | sir | earl |
| "ARMY" | commander | forces | war | general |
| "STUDY" | analysis | space | program | user |
| "PARTY" | act | office | judge | justice |
| "DESIGN" | size | glass | device | memory |
| "PUBLIC" | report | health | community | industry |

$z$ (topics) [PT13]



$x$ (image)



$z$ (semantic regions) [DFD+18]

# Generative Model: Overview

- What can generative models do:

4. Representation learning: semantic and concise (via latent variable $z$).

Dimensionality
Reduction:



$x \in \mathbb{R}^{\#\text{vocabulary}}$

Topic proportion
$z \in \mathbb{R}^{\#\text{topic}}$

[PT13]

$x \in \mathbb{R}^{28 \times 28}$

$z \in \mathbb{R}^{20}$ [DFD+18]

# Generative Model: Overview

- What can generative models do:

5. Supervised Learning: query $p(y|x)$ from $p(x,y)$.



$Y$   {bird, mammal}

$X_1$   $X_2$   $X_3$   $X_4$

has beak?   can fly?   has fur?   has four legs?

Naive Bayes

$z$: topics

| "ENGINES" | speed | product | introduced |
| "ROYAL" | britain | queen | sir |
| "ARMY" | commander | forces | war |
| "STUDY" | analysis | space | program |

$x_1$: doc 1     $y_1$: science & tech

$x_2$: doc 2     $y_2$: politics

Supervised LDA [MB08]

# Generative Model: Overview

- What can generative models do:

  5. Supervised Learning: query $p(y|x)$ from $p(x, y)$.

  Semi-Supervised Learning:

  Unlabeled data $\{x^{(n)}\}$ can be utilized to learn a better $p(x, y)$.

# Generative Model: Benefits

*"What I cannot create, I do not understand."*           *—Richard Feynman*

- Natural for generation (*randomness/diversity, high-dimensional*).
- For representation learning: responsible and faithful knowledge of data.
- For supervised learning:
  - Leverage unlabeled data: semi-supervised learning.
  - Data-efficient: for logistic regression (discriminative) and naive Bayes (generative) [NJ01],

$$\epsilon_{\text{Dis},N} \leq \epsilon_{\text{Dis},\infty} + O\left(\sqrt{\frac{d}{N}}\right)$$

$$\epsilon_{\text{Gen},N} \leq \epsilon_{\text{Gen},\infty} + O\left(\sqrt{\frac{\log d}{N}}\right)$$

$d$: data dimension.
$N$: data size.

# Generative Model: Taxonomy

- Plain Generative Models: Directly model $p(x)$; no latent variable. $p_\theta(x)$ (x)

- Latent Variable Models:
  - Deterministic Generative Models: Dependency between $x$ and $z$ is *deterministic*: $x = f_\theta(z)$.

  - Probabilistic Graphical Models: Dependency between $x$ and $z$ is *probabilistic*: $(x, z) \sim p_\theta(x, z)$.

# Generative Model: Taxonomy

- Latent Variable Models
  - Probabilistic Graphical Models (PGM):

  - Directed PGM:
    $p(x, z)$ specified by $p(z)$ and $p(x|z)$.

  - Undirected PGM:
    $p(x, z)$ specified by an Energy function:
    $p_\theta(x, z) \propto \exp(-E_\theta(x, z))$.



$p(z)$ $z$

$x \sim p_\theta(x|z)$

$x$

$p_\theta(x)$



$z$

$p_\theta(x, z) \propto \exp(-E_\theta(x, z))$

$x$

# Generative Model: Taxonomy

- Overview

# Outline

- Generative Models: Overview
- **Plain Generative Models**
  - **Autoregressive Models**
- Latent Variable Models
  - Deterministic Generative Models
    - Generative Adversarial Nets
    - Flow-Based Models
  - Probabilistic Graphical Models
    - Directed PGMs
      - Bayesian Inference (variational inference, MCMC)
      - Topic models (LDA, LightLDA, sLDA)
      - Deep Bayesian Models (VAE)
    - Undirected PGMs (Boltzmann machines, energy-based models)
    - Diffusion-Based Models

# Plain Generative Models

- Directly model $p_\theta(x)$ (parameter $\theta$) without latent variable.

- Easy to learn (no normalization issue of data likelihood) and use (data generation).

- Learning: **Maximum Likelihood Estimation (MLE)**.

$$\theta^* = \arg\max_\theta \mathbb{E}_{\hat{p}(x)}[\log p_\theta(x)] = \arg\min_\theta \mathrm{KL}(\hat{p}, p_\theta)$$

$$\approx \arg\max_\theta \frac{1}{N} \sum_{n=1}^{N} \log p_\theta(x^{(n)}).$$

Kullback-Leibler divergence

$$\mathrm{KL}(\hat{p}, p_\theta) := \mathbb{E}_{\hat{p}(x)}\left[\log\frac{\hat{p}(x)}{p_\theta(x)}\right]$$

- First example: Gaussian Mixture Model

$$p_\theta(x) = \sum_{k=1}^{K} \alpha_k \mathcal{N}(x|\mu_k, \Sigma_k),$$

$$\theta = (\alpha, \mu, \Sigma).$$

# Autoregressive Models



Model $p(x)$ by each conditional $p(x_i|x_{<i})$ ($i$ indices components).

- Full dependency can be restored.
- Conditionals are easier to model.
- Easy data generation:
$$x \sim p(x) \Longleftrightarrow x_1 \sim p(x_1), x_2 \sim p(x_2|x_1), \dots, x_d \sim p(x_d|x_1, \dots, x_{d-1}).$$
  But **non-parallelizable**.

# Autoregressive Models

- Fully Visible Sigmoid Belief Network [Fre98]
$$p(x_i|x_{<i}) = \text{Bern}\left(x_i\middle|\sigma\left(\sum_{j<i} W_{ij}x_j\right)\right)$$

Sigmoid function
$$\sigma(r) = \frac{1}{1+e^{-r}}$$

- Neural Autoregressive Distribution Estimator [LM11]
$$p(x_i|x_{<i}) = \text{Bern}\left(x_i\middle|\sigma\left(V_{i,:}\sigma\left(W_{:,<i}x_{<i} + a\right) + b_i\right)\right)$$

- A typical language model: Use a hidden state to represent the dependency on previous items.

$p(\mathbf{x} =$ "*the cat sat on the mat*")

$= p(x_1 = \text{the})\ p(\text{cat}|x_1)\ p(\text{sat}|x_{1\dots2})\ p(\text{on}|x_{1\dots3})\ p(\text{the}|x_{1\dots4})\ p(\text{mat}|x_{1\dots5})\ p(</s>|x_{1\dots6})$

# Autoregressive Models

- WaveNet [ODZ+16]
  - Construct $p(x_i|x_{<i})$ via Causal Convolution

# Autoregressive Models

- PixelCNN & PixelRNN [OKK16]
  - Autoregressive structure of an image:



  - PixelCNN: model conditional distributions via (masked) convolution:
$$h_i = K * x_{<i},$$
$$p(x_i|x_{<i}) = \text{NN}(h_i).$$
    - Bounded receptive field.
    - Likelihood evaluation: parallel

# Autoregressive Models

- PixelCNN & PixelRNN [OKK16]
  - PixelRNN: model conditional distributions via recurrent connection:

  $$[h_i, c_i] = \text{LSTM}\left(\overbrace{K * h_{(\lfloor i/n \rfloor n - n):\lfloor i/n \rfloor n}}^{\text{1D convolution}}, c_{i-1}, x_{i-1}\right),$$
  $$p(x_i|x_{<i}) = \text{NN}(h_i).$$

  - Unbounded receptive field.
  - Likelihood evaluation (in-row): parallel
    Likelihood evaluation (inter-row): sequential

# Autoregressive Models

- PixelCNN & PixelRNN [OKK16]



Image Generation



Image Completion

# Outline

- Generative Models: Overview
- Plain Generative Models
  - Autoregressive Models
- **Latent Variable Models**
  - Deterministic Generative Models
    - Generative Adversarial Nets
    - Flow-Based Models
  - Probabilistic Graphical Models
    - Directed PGMs
      - Bayesian Inference (variational inference, MCMC)
      - Topic models (LDA, LightLDA, sLDA)
      - Deep Bayesian Models (VAE)
    - Undirected PGMs (Boltzmann machines, energy-based models)
    - Diffusion-Based Models

# Latent Variable Models

- Latent Variable:
  - Abstract knowledge of data; enables various tasks.

Knowledge Discovery

| "ENGINES" | speed | product | introduced |
| "ROYAL" | britain | queen | sir |
| "ARMY" | commander | forces | war |
| "STUDY" | analysis | space | program |
| "PARTY" | act | office | judge |
| "DESIGN" | size | glass | device |
| "PUBLIC" | report | health | community |



Dimensionality Reduction

Manipulated Generation

# Latent Variable Models

- Latent Variable:
  - Compact representation of dependency.

    De Finetti's Theorem (1955): if $(x_1, x_2, \ldots)$ are *infinitely exchangeable*, then $\exists$ r.v. $z$ and $p(\cdot \,|z)$ s.t. $\forall n$,

$$p(x_1, \ldots, x_n) = \int \left( \prod_{i=1}^{n} p(x_i|z) \right) p(z)\, \mathrm{d}z \,.$$



*Infinite exchangeability*:

For all $n$ and permutation $\sigma$, $p(x_1, \ldots, x_n) = p\left(x_{\sigma(1)}, \ldots, x_{\sigma(n)}\right)$.

# Outline

- Generative Models: Overview
- Plain Generative Models
  - Autoregressive Models
- **Latent Variable Models**
  - **Deterministic Generative Models**
    - **Generative Adversarial Nets**
    - Flow-Based Models
  - Probabilistic Graphical Models
    - Directed PGMs
      - Bayesian Inference (variational inference, MCMC)
      - Topic models (LDA, LightLDA, sLDA)
      - Deep Bayesian Models (VAE)
    - Undirected PGMs (Boltzmann machines, energy-based models)
    - Diffusion-Based Models

# Generative Adversarial Nets

- Deterministic $f_\theta : z \mapsto x$, modeled by a neural network.
  + Flexible modeling ability.
  + Good generation performance.
  - Hard to infer $z$ of a data point $x$.
  - Unavailable p.d.f/p.m.f $p_\theta(x)$.
  - Mode-collapse.

- Learning: $\min_\theta \mathrm{discr}\big(\hat{p}(x), p_\theta(x)\big)$.

  - discr. $= \mathrm{KL}(\hat{p}, p_\theta) \implies$ MLE: $\max_\theta \mathbb{E}_{\hat{p}}[\log p_\theta]$, but the p.d.f/p.m.f $p_\theta(x)$ is unavailable!

  - discr. $=$ Jensen-Shannon divergence [GPM+14].
  - discr. $=$ Wasserstein distance [ACB17].

$p(z)$ $z$

$x = f_\theta(z)$
(Neural Nets)

$x$

$p_\theta(x)$

# Generative Adversarial Nets

- Learning: $\min_\theta \text{discr}(\hat{p}(x), p_\theta(x))$.

  - GAN [GPM+14]: discr. = Jensen-Shannon divergence.

$$\text{JS}(\hat{p}, p_\theta) := \frac{1}{2}\left(\text{KL}\left(\hat{p}, \frac{p_\theta + \hat{p}}{2}\right) + \text{KL}\left(p_\theta, \frac{p_\theta + \hat{p}}{2}\right)\right)$$

$$= \frac{1}{2}\max_{T(\cdot)} \mathbb{E}_{\hat{p}(x)}\left[\log \sigma(T(x))\right] + \underbrace{\mathbb{E}_{p_\theta(x)}\left[\log\left(1 - \sigma(T(x))\right)\right]}_{= \mathbb{E}_{p(z)}\left[\log\left(1 - \sigma\left(T(f_\theta(z))\right)\right)\right]} + \log 2.$$

- $\sigma(T(x))$ is the discriminator; $T$ implemented as a neural network.

- Expectations can be estimated by samples.

$p(z)$   $z$

$x = f_\theta(z)$

(Neural Nets)

$x$

$p_\theta(x)$

# Generative Adversarial Nets

- Learning: $\min_{\theta} \text{discr}(\hat{p}(x), p_{\theta}(x))$.

  - WGAN [ACB17]: discr. = Wasserstein distance:
  $$d_W(\hat{p}, p_{\theta}) = \inf_{\gamma \in \Gamma(\hat{p}, p_{\theta})} \mathbb{E}_{\gamma(x,y)}[c(x, y)]$$
  $$= \sup_{\phi \in \text{Lip}_1} \mathbb{E}_{\hat{p}}[\phi] - \mathbb{E}_{p_{\theta}}[\phi].$$

    - Choose $\phi$ as a neural network with parameter clipping.
    - Benefit: $d_W$ has more alleviative reaction to distribution difference than JS.
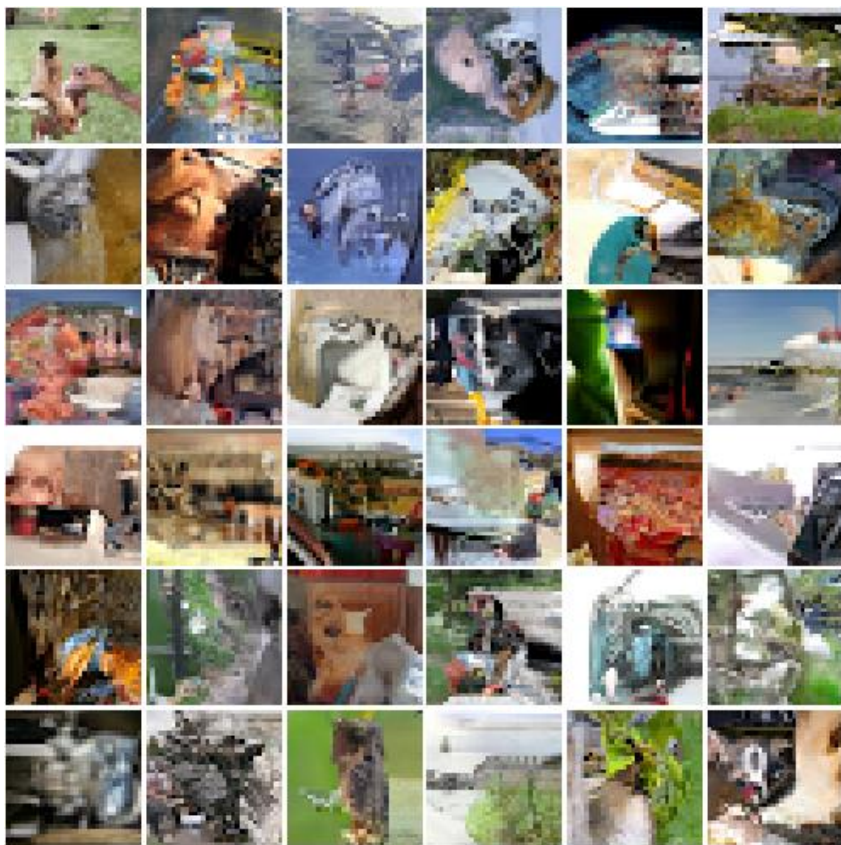
$p(z)$  $z$

$x = f_{\theta}(z)$

(Neural Nets)
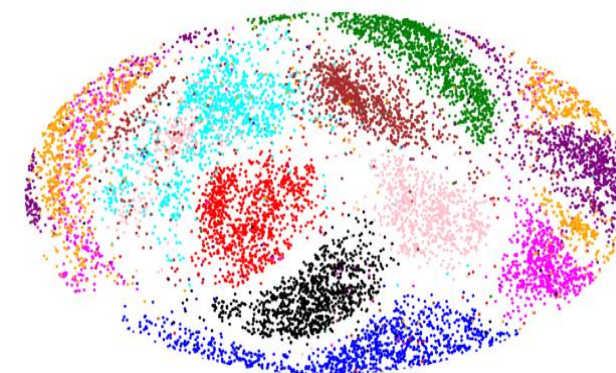
$x$

$p_{\theta}(x)$

# Outline

- Generative Models: Overview
- Plain Generative Models
  - Autoregressive Models
- **Latent Variable Models**
  - **Deterministic Generative Models**
    - Generative Adversarial Nets
    - **Flow-Based Models**
  - Probabilistic Graphical Models
    - Directed PGMs
      - Bayesian Inference (variational inference, MCMC)
      - Topic models (LDA, LightLDA, sLDA)
      - Deep Bayesian Models (VAE)
    - Undirected PGMs (Boltzmann machines, energy-based models)
    - Diffusion-Based Models

# Flow-Based Models



- Deterministic and invertible $f_\theta : z \mapsto x$.
  + Available density function!

$$p_\theta(x) = p\left(z = f_\theta^{-1}(x)\right)\left|\frac{\partial f_\theta^{-1}}{\partial x}\right| \quad \text{(rule of change of variables)}.$$

  + Easy inference: $z = f_\theta^{-1}(x)$.
  - Redundant representation: dim. $z$ = dim. $x$.
  - Restricted $f_\theta$: deliberative design; either $f_\theta$ or $f_\theta^{-1}$ computes costly.

Jacobian determinant, $\left(\frac{\partial f_\theta^{-1}}{\partial x}\right)_{ij} := \frac{\partial (f_\theta^{-1})_i}{\partial x_j}$.

- Learning: $\min_\theta \text{KL}\left(\hat{p}(x), p_\theta(x)\right) \Rightarrow$ MLE: $\max_\theta \mathbb{E}_{\hat{p}(x)}[\log p_\theta(x)]$.

- Examples:
  - NICE [DKB15], RealNVP [DSB17], MAF [PPM17], GLOW [KD18].
  - Also used for variational inference [RM15, KSJ+16].

# Flow-Based Models

- RealNVP [DSB17]
  - Building block: **Coupling**: $y = g(x)$,

$$\begin{cases} y_{1:d} & = x_{1:d} \\ y_{d+1:D} & = x_{d+1:D} \odot \exp\left(s(x_{1:d})\right) + t(x_{1:d}) \end{cases}$$

$$\Leftrightarrow \begin{cases} x_{1:d} & = y_{1:d} \\ x_{d+1:D} & = \left(y_{d+1:D} - t(y_{1:d})\right) \odot \exp\left(-s(y_{1:d})\right), \end{cases}$$



(a) Forward propagation  (b) Inverse propagation

where $s$ and $t: \mathbb{R}^{D-d} \to \mathbb{R}^{D-d}$ are general functions for scale and translation.

- Jacobian Determinant: $\left|\frac{\partial g}{\partial x}\right| = \exp(\sum_{j=1}^{D-d} s_j(x_{1:d}))$.

- Partitioning $x$ using a binary mask $b$:

$$y = b \odot x + (1 - b) \odot \left(x \odot \exp\left(s(b \odot x)\right) + t(b \odot x)\right).$$

# Flow-Based Models

- RealNVP [DSB17]
  - Building block: **Squeezing**: from $s \times s \times c$ to $\frac{s}{2} \times \frac{s}{2} \times 4c$:



  - Combining with a multi-scale architecture:

$$h^{(0)} = x$$

$$(z^{(i+1)}, h^{(i+1)}) = f^{(i+1)}(h^{(i)})$$

$$z^{(L)} = f^{(L)}(h^{(L-1)})$$

$$z = (z^{(1)}, \ldots, z^{(L)}).$$

where each $f$ follows a "coupling-squeezing-coupling" architecture.

# Flow-Based Models

- RealNVP [DSB17]

# Flow-Based Models

- GLOW [KD18]

One step of $f_\theta$



| affine coupling layer |
| invertible 1x1 conv |
| actnorm |

Combination of the steps to form $f_\theta$



| step of flow | × K |
| squeeze | |
| split | → $z_i$ |
| step of flow | × K | × (L−1) |
| squeeze | |

$z_L$ — $x$

Component Details

| Description | Function | Reverse Function | Log-determinant |
|---|---|---|---|
| Actnorm. See Section 3.1. | $\forall i,j : \mathbf{y}_{i,j} = \mathbf{s} \odot \mathbf{x}_{i,j} + \mathbf{b}$ | $\forall i,j : \mathbf{x}_{i,j} = (\mathbf{y}_{i,j} - \mathbf{b})/\mathbf{s}$ | $h \cdot w \cdot \mathtt{sum}(\log|\mathbf{s}|)$ |
| Invertible $1 \times 1$ convolution. $\mathbf{W} : [c \times c]$. See Section 3.2. | $\forall i,j : \mathbf{y}_{i,j} = \mathbf{W}\mathbf{x}_{i,j}$ | $\forall i,j : \mathbf{x}_{i,j} = \mathbf{W}^{-1}\mathbf{y}_{i,j}$ | $h \cdot w \cdot \log|\det(\mathbf{W})|$ or $h \cdot w \cdot \mathtt{sum}(\log|\mathbf{s}|)$ (see eq. (10)) |
| Affine coupling layer. See Section 3.3 and (Dinh et al., 2014) | $\mathbf{x}_a, \mathbf{x}_b = \mathtt{split}(\mathbf{x})$ $(\log \mathbf{s}, \mathbf{t}) = \mathtt{NN}(\mathbf{x}_b)$ $\mathbf{s} = \exp(\log \mathbf{s})$ $\mathbf{y}_a = \mathbf{s} \odot \mathbf{x}_a + \mathbf{t}$ $\mathbf{y}_b = \mathbf{x}_b$ $\mathbf{y} = \mathtt{concat}(\mathbf{y}_a, \mathbf{y}_b)$ | $\mathbf{y}_a, \mathbf{y}_b = \mathtt{split}(\mathbf{y})$ $(\log \mathbf{s}, \mathbf{t}) = \mathtt{NN}(\mathbf{y}_b)$ $\mathbf{s} = \exp(\log \mathbf{s})$ $\mathbf{x}_a = (\mathbf{y}_a - \mathbf{t})/\mathbf{s}$ $\mathbf{x}_b = \mathbf{y}_b$ $\mathbf{x} = \mathtt{concat}(\mathbf{x}_a, \mathbf{x}_b)$ | $\mathtt{sum}(\log(|\mathbf{s}|))$ |

# Flow-Based Models

- GLOW [KD18]

Generation Results (Interpolation)

Generation Results (Manipulation; each semantic direction = $\bar{z}_{\text{pos}} - \bar{z}_{\text{neg}}$)



(a) Smiling

(b) Pale Skin

(c) Blond Hair

(d) Narrow Eyes

(e) Young

(f) Male

# Flow-Based Models

- Autoregressive flows

  [KSJ+16, PPM17].

  - Tractable inverse & JacDet.
  - One direction is non-parallelable.
  - Universal approximator [TIT+20].



Masked Autoregressive Flow (MAF)     Inverse Autoregressive Flow (IAF)

[source]

- Continuous normalizing flow [GCB+18].

$$\partial_t z_t = f_t(z_t) \implies \frac{\mathrm{d}}{\mathrm{d}t} \log p_t(z_t) = -\nabla \cdot f_t(z_t) = -\mathrm{tr}\left(\frac{\partial f_t}{\partial z}\right).$$

  - Use ODE solver for fwd/bwd map and $\log p_{t_1}\big(z(t_1)\big) = \log p_{t_0}\big(z(t_0)\big) - \int_{t_0}^{t_1} \mathrm{tr}\left(\frac{\partial f_t}{\partial z}\right) \mathrm{d}t.$

# Flow-Based Models

- Residual flows.
  - ResNet block $x_{t+1} := F_{\theta_t}(x_t) := x_t + g_{\theta_t}(x_t)$ is invertible if $\text{Lip}(g_{\theta_t}) < 1$.
  - Inverse map: fixed-point iteration.
  - JacDet: $\ln \det J_F = \sum_{k=1}^{\infty} (-1)^{k+1} \frac{\text{tr}(J_g^k)}{k}$.
    - Hutchinson's trace estimator: $\text{tr}(A) = \mathbb{E}[v^\top A v]$, where $\mathbb{E}[v] = 0$, $\text{cov}[v] = I$.
    - Truncated estimate [BGC19].
    - Unbiased "Russian roulette" estimator [CBD19].
- Approximation theory.
  - [KC20, TIT+20].

| 1. Det Identities | 2. Coupling Blocks | 3. Autoregressive | 4. Unbiased Estimation |
|---|---|---|---|
| Planar NF Sylvester NF ... | NICE Real NVP Glow ... | Inverse AF Neural AF Masked AF ... | FFJORD **Residual Flows** |
| (Low rank) | (Lower triangular + structured) | (Lower triangular) | (Arbitrary) |

Jacobian

# Outline

- Generative Models: Overview
- Plain Generative Models
  - Autoregressive Models
- Latent Variable Models
  - Deterministic Generative Models
    - Generative Adversarial Nets
    - Flow-Based Models
  - **Probabilistic Graphical Models**
    - Directed PGMs
      - Bayesian Inference (variational inference, MCMC)
      - Topic models (LDA, LightLDA, sLDA)
      - Deep Bayesian Models (VAE)
    - Undirected PGMs (Boltzmann machines, energy-based models)
    - Diffusion-Based Models

# Classical Probabilistic Graphical Models

- Generally, they may or may not have latent variables.
- Intuitively: represent variable **relations** by a graph.
- Formally: a way to represent a joint distribution by making **conditional independence (CI)** assumptions.



$$p(S, R, G) = p(S)p(R)p(G|S, R)$$



$$p(x) \propto \exp\underbrace{\left(-\sum_{(i,j)\in\mathcal{E}} J(x_i, x_j) - \sum_i H(x_i)\right)}_{\text{Energy function } -E(x)}$$

# Directed Probabilistic Graphical Models

- Represented by a **Directed Acyclic Graph** (DAG).
- Synonyms: Bayesian/belief/causal network.



$$p(S, R, G) = p(S)p(R)p(G|S, R)$$

Markovianess

CI assumptions:
- $S \perp R$ since $p(S, R) = p(S)p(R)$.
- $S \,!\perp R|G$ since $p(S, R|G) \neq p(S|G)p(R|G)$ in general.

Faithfulness

CI assumptions:
- $C \perp N|E$ since $p(C, N|E) = p(C|E)p(N|E)$.
- $C \,!\perp N$ since $p(C, N) \neq p(C)p(N)$ in general.

$$p(E, C, N) = p(E)p(C|E)p(N|E)$$

# Directed Probabilistic Graphical Models

- **d-separation**: read off encoded CI assumptions in general.
  - A path is called **d-separated** by a set of nodes $\mathbf{S}$, if $P$
    either has an *emitter* $X$ ("$\to X \to$" or "$\leftarrow X \to$") in $\mathbf{S}$,
    or has a *collider* $X$ ("$\to X \leftarrow$") that is not in $\mathbf{S}$ nor is any descendant of $X$.

  $p$ is Markovian w.r.t the DAG

  - $A \perp B | \mathbf{S}$    $\Longleftarrow$      All paths between $A$ and $B$ are d-separated by $\mathbf{S}$.

  $p$ is faithful w.r.t the DAG

$YF \perp AD$, but $YF \;!\perp AD | LC$.



[YGL+20]

# Directed Probabilistic Graphical Models

As a language of causality

- Formal definition of causality:
  "*two variables have a causal relation, if **intervening** the cause may change the effect, but not vice versa*" [Pearl09, PJS17].
  - **Intervention**: change the value of a variable by leveraging mechanisms and changing variables out of the considered system.
- Example: for the $A$ltitude and average $T$emperature of a city, $A \rightarrow T$.
  - Running a huge heater (intv. $T$) does not lower $A$.
  - Raising the city by a huge elevator (intv. $A$) lowers $T$.
- Causality contains more information than observation (== static/observational data, joint distribution, CIs).
  - Both $p(A)p(T|A)$ ($A \rightarrow T$) and $p(T)p(A|T)$ ($T \rightarrow A$) can describe $p(A, T)$,
  - but they give different outcomes under intervention.

# Directed Probabilistic Graphical Models

As a language of causality

- Pearl's surgery: describing intervention.



Intervening $T$ with $t$

do-operation

$$p(Y|t) = \int p(Y|X,t)p(X|t)\,\mathrm{d}X$$

$$p(Y|do(t)) = \tilde{p}(Y|t) = \int p(Y|X,t)p(X)\,\mathrm{d}X$$

- Explaining spurious correlation:

(Confounding bias)



Intervening *Choco.* with $c$

$$C \;!\perp N$$

$$p(N|do(c)) = p(N)$$

# Directed Probabilistic Graphical Models

As a language of causality

- **Causal inference**: estimate causal effect $\mathbb{E}[Y|do(t=1)] - \mathbb{E}[Y|do(t=0)]$.



Intervening $T$ with $t$

do-operation

$$p(Y|t) = \int p(Y|X,t)\,p(X|t)\,\mathrm{d}X \qquad\qquad p(Y|do(t)) = \tilde{p}(Y|t) = \int p(Y|X,t)\,p(X)\,\mathrm{d}X$$

  - Under some assumptions, it is identifiable from observation [IW09, Pearl15].

- **Causal discovery**: recover the causal DAG from observation.
  - Constraint-based (e.g., PC alg. [SG91]):
    CIs could recover some structures (e.g., $A \perp B, A \; ! \perp B|C \Longrightarrow A \to C \leftarrow B$).
  - Score-based (i.e., likelihood-based): some DAGs could better fit observation data.
    - Additive noise assumption: a function class restriction makes identifiability.

# Undirected Probabilistic Graphical Models

- For **symmetric** relations (e.g., image pixels), it is **unnatural** to assign a direction.
  - Side effect: there would be undesired or arbitrary CI assertions.

- Represent the relation by an **undirected** graph.
  - **Synonyms**: Markov random field, energy-based model.
  - **d-separation**: every path between $A$ and $B$ contains a node in **S**.
  - **Markovianess** (Hammersley-Clifford theorem):
    $p$ satisfies graph CI properties if it factorizes as one
    term per maximal *clique* (fully connected subgraph).

Ising model

Markovianess

$$p(x) \propto \exp \underbrace{\left( - \sum_{(i,j) \in \mathcal{E}} J(x_i, x_j) - \sum_i H(x_i) \right)}_{\text{Energy function } -E(x)}$$

# Probabilistic Graphical Models

- Directed and Undirected PGMs cover **different** distributions.

- Not all PGMs are generative
(e.g., Bayesian neural networks, conditional random fields).

- Classical PGMs do emphasize the "graph" information.

- Deep PGMs often have simple graphs, and focus on learning the edge relation:

Dependency between $x$ and $z$ is *probabilistic*: $(x, z) \sim p_\theta(x, z)$.

Directed PGM:

$p(z)$ ⬡ $z$

$x \sim p_\theta(x|z)$

$p_\theta(x)$ ⬤ $x$

Undirected PGM:

$z$

$p_\theta(x, z) \propto \exp(-E_\theta(x, z))$

$x$

# Outline

- Generative Models: Overview
- Plain Generative Models
  - Autoregressive Models
- Latent Variable Models
  - Deterministic Generative Models
    - Generative Adversarial Nets
    - Flow-Based Models
  - Probabilistic Graphical Models
    - **Directed PGMs**
      - Bayesian Inference (variational inference, MCMC)
      - Topic models (LDA, LightLDA, sLDA)
      - Deep Bayesian Models (VAE)
    - Undirected PGMs (Boltzmann machines, energy-based models)
    - Diffusion-Based Models

# Directed PGMs

Bayesian models

- Model structure (*Bayesian Modeling*):
  - *Prior $p(z)$*: *initial* belief of $z$.
  - *Likelihood $p(x|z)$*: dependence of $x$ on $z$.

- Learning: MLE.
$$\theta^* = \arg\max_{\theta} \mathbb{E}_{\hat{p}(x)}[\log p_\theta(x)],$$

*Evidence $p(x) = \int p(z,x)\,\mathrm{d}z$.*

- Feature/representation learning (*Bayesian Inference*):
$$Posterior\ p(z|x) = \frac{p(z,x)}{p(x)} = \frac{p(z)p(x|z)}{\int p(z,x)\,\mathrm{d}z} \ \text{(Bayes' rule)}$$

  represents the *updated* information that observation $x$ conveys to latent $z$.

Prior $p(z)$ — Latent Variable $z$ — Posterior $p(z|x)$

*Bayesian Modeling* — *(Bayesian Inference)*

Likelihood $p(x|z)$

Evidence $p(x)$ — Data Variable $x$

# Directed PGMs

Not all Bayesian models are generative:



| | **Generative** | **Non-generative** |
|---|---|---|
| Supervised | Naive Bayes, Supervised LDA | Bayesian Logistic Regression, Bayesian Neural Networks |
| Unsupervised | BayesNets (LDA, VAE), MRFs (BM, RBM, DBM) | (invalid task) |

# Directed PGMs

Benefits of Bayesian models:

- Robust to small data.

- Stable training process.

- Principled and natural inference $p(z|x)$ via Bayes' rule.

- Natural to incorporate prior knowledge:

Problem of knowledge-agnostic conditional generation:

Moustache



v.s.



[KSDV18]

# Outline

- Generative Models: Overview
- Plain Generative Models
  - Autoregressive Models
- Latent Variable Models
  - Deterministic Generative Models
    - Generative Adversarial Nets
    - Flow-Based Models
  - Probabilistic Graphical Models
    - Directed PGMs
      - **Bayesian Inference** (**variational inference**, MCMC)
      - Topic models (LDA, LightLDA, sLDA)
      - Deep Bayesian Models (VAE)
    - Undirected PGMs (Boltzmann machines, energy-based models)

# Bayesian Inference

Estimate the posterior $p(z|x)$.



Bayesian Modeling

Bayesian Inference

Bayes' rule: $\textit{Posterior } p(z|x) = \dfrac{p(x,z)}{p(x)} = \dfrac{p(x,z)}{\int p(x,z)\,\mathrm{d}z} \propto p(x,z) = p(z)p(x|z).$

# Bayesian Inference

Estimate the posterior $p(z|x)$.

- Infer unobserved variables from observation.

Naive Bayes: $z = y$.

$p(y = 0|x)$

$$= \frac{p(x|y = 0)p(y = 0)}{p(x|y = 0)p(y = 0) + p(x|y = 1)p(y = 1)}.$$

$f(x) = \arg\max_{y} p(y|x)$ achieves the lowest error

$\int p\big(y = (1 - f(x)) \,|\, x\big) \, p(x) \, \mathrm{d}x.$

$Y$ $\{\text{bird, mammal}\}$

$X_1$  $X_2$  $X_3$  $X_4$

has beak?  can fly?  has fur?  has four legs?

# Bayesian Inference

Estimate the posterior $p(z|x)$.

- Extract knowledge/representation from data.



| | | | | | |
|---|---|---|---|---|---|
| "ENGINES" | speed | product | introduced | designs | fuel |
| "ROYAL" | britain | queen | sir | earl | died |
| "ARMY" | commander | forces | war | general | military |
| "STUDY" | analysis | space | program | user | research |
| "PARTY" | act | office | judge | justice | legal |
| "DESIGN" | size | glass | device | memory | engine |
| "PUBLIC" | report | health | community | industry | conference |
| "CHURCH" | prayers | communion | religious | faith | historical |

Prior $p(z)$

Posterior $p(z|x)$

Topics

Latent Variable $z$

Bayesian Modeling

Bayesian Inference

Likelihood $p(x|z)$

Data Variable $x$

Documents

# Bayesian Inference

Estimate the posterior $p(z|x)$.

- For prediction:

$$p(y^*|x^*, \{x, y\}_{\text{train}}) = \begin{cases} \int p(y^*|z, x^*)p(z|x^*, \{x, y\}_{\text{train}}) \, \mathrm{d}z\,, \\[2em] \int p(y^*|z, x^*)p(z|\{x, y\}_{\text{train}}) \, \mathrm{d}z\,. \end{cases}$$

(Generative)

(Non-Generative)

# Bayesian Inference

Estimate the posterior $p(z|x)$.

$$p(z|x) = \frac{p(x,z)}{p(x)} = \frac{p(x,z)}{\int p(x,z)\, \mathrm{d}z}$$

Intractable!

# Bayesian Inference

- Variational inference (VI)

  Use a *tractable* variational distribution $q(z)$ to approximate $p(z|x)$:
  $$\min_{q \in \mathcal{Q}} \mathrm{KL}\big(q(z), p(z|x)\big).$$

  Tractability: known density function, or samples are easy to draw.

  - Parametric VI: use a parameter $\phi$ to represent $q_\phi(z)$.

  - Particle-based VI: use a set of particles $\{z^{(i)}\}_{i=1}^N$ to represent $q(z)$.

- Monte Carlo (MC)

  - Draw samples from $p(z|x)$.

  - Typically done by simulating a *Markov chain* (i.e., MCMC) for tractability.

# Bayesian Inference: Variational Inference

*"Feed two birds with one scone."*

- In model learning: $\mathbb{E}_{\hat{p}(x)}[\log p_\theta(x)] = \frac{1}{N}\sum_{n=1}^{N}\log p_\theta(x^{(n)})$.

  - Introduce a *variational distribution* $q(z)$:
$$\log p_\theta(x) = \mathcal{L}_\theta[q(z)] + \text{KL}(q(z), p_\theta(z|x)),$$
$$\mathcal{L}_\theta[q(z)] := \mathbb{E}_{q(z)}[\log p_\theta(z,x)] - \mathbb{E}_{q(z)}[\log q(z)].$$

  - $\mathcal{L}_\theta[q(z)] \leq \log p_\theta(x)$ ➜ **Evidence Lower BOund (ELBO)**!
  - $\mathcal{L}_\theta[q(z)]$ is easier to estimate.

- (Variational) Expectation-Maximization Algorithm:

  (a) E-step: Let $\mathcal{L}_\theta[q(z)] \approx \log p_\theta(x)$, that is $\overbrace{\min_{q\in\mathcal{Q}}\text{KL}(q(z), p_\theta(z|x))}^{\text{Bayesian Inference}}$;

  (b) M-step: $\max_\theta \mathcal{L}_\theta[q(z)]$.

  - Classical EM: take $q(z) = p_\theta(z|x)$ (i.e., with exact inference).

# Bayesian Inference: Variational Inference

*"Feed two birds with one scone."*

- To do Bayesian inference by: $\min\limits_{q \in \mathcal{Q}} \text{KL}\big(q(z), p(z|x)\big),$

  $\text{KL}\big(q(z), p_\theta(z|x)\big)$ is hard to compute...

  Note $\qquad \log p_\theta(x) = \mathcal{L}_\theta[q(z)] + \text{KL}\big(q(z), p_\theta(z|x)\big),$

  so $\qquad \min\limits_{q \in \mathcal{Q}} \text{KL}\big(q(z), p(z|x)\big) \iff \max\limits_{q \in \mathcal{Q}} \mathcal{L}_\theta[q(z)].$

  The ELBO $\mathcal{L}_\theta[q(z)] = \mathbb{E}_{q(z)}[\log p_\theta(z, x)] - \mathbb{E}_{q(z)}[\log q(z)]$ is easier to compute.

# Bayesian Inference: Variational Inference

- **Parametric variational inference**: use a parameter $\phi$ to represent $q_\phi(z)$.

$$\max_\phi \left( \textcolor{red}{\mathcal{L}_\theta[q_\phi(z)]} = \mathbb{E}_{q_\phi(z)}[\log p_\theta(z, x)] - \mathbb{E}_{q_\phi(z)}[\log q_\phi(z)] \right).$$

- For model-specifically designed $q_\phi(z)$, $\textcolor{red}{\mathcal{L}_\theta[q_\phi(z)]}$ has closed form (e.g., [SJJ96] for SBN, [BNJ03] for LDA).

- Main Challenge:
  - $\mathcal{Q}$ should be as large/general/flexible as possible,
  - while enables practical optimization of the ELBO.



$$p(z|x)$$

$$\mathcal{Q}$$

$$q^* = \arg\min_{q \in \mathcal{Q}} \mathrm{KL}\big(q(z), p(z|x)\big)$$

# Bayesian Inference: Variational Inference

- **Parametric variational inference**: use a parameter $\phi$ to represent $q_\phi(z)$.

$$\max_\phi \left( \textcolor{red}{\mathcal{L}_\theta[q_\phi(z)]} = \mathbb{E}_{q_\phi(z)}[\log p_\theta(z, x)] - \mathbb{E}_{q_\phi(z)}[\log q_\phi(z)] \right).$$

  - **Explicit variational inference**: specify the form of the density function $q_\phi(z)$.
    - [GHB12, HBWP13, RGB14]: model-agnostic $q_\phi(z)$ (e.g., mixture of Gaussians).
    - [RM15, KSJ+16]: define $q_\phi(z)$ by a flow-based generative model.
  - **Implicit variational inference**: define $q_\phi(z)$ by a GAN-like generative model.
    - More flexible but more difficult to optimize.
    - Density ratio estimation: [MNG17, SSZ18a].

$$\textcolor{red}{\mathcal{L}_\theta[q_\phi(z)]} = \mathbb{E}_{q_\phi(z)}[\log p_\theta(x|z)] - \mathbb{E}_{q_\phi(z)}\left[\log \frac{q_\phi(z)}{p(z)}\right].$$

    - Gradient Estimation $\nabla \log q_\phi(z)$: [VLBM08, LT18, SSZ18b].

# Bayesian Inference: Variational Inference

$$\min_{q \in \mathcal{Q}} \mathrm{KL}\big(q(z), p(z|x)\big).$$

- **Particle-based variational inference**: use particles $\left\{z^{(i)}\right\}_{i=1}^{N}$ to represent $q(z)$.

  To minimize $\mathrm{KL}\big(q(z), p(z|x)\big)$, simulate its gradient flow on the Wasserstein space.

  - Wasserstein space:
    an abstract space of distributions.
  - Wasserstein tangent vector
    $\Longleftrightarrow$ vector field.

# Bayesian Inference: Variational Inference

$$\min_{q \in \mathcal{Q}} \mathrm{KL}\big(q(z), p(z|x)\big).$$

- **Particle-based variational inference**: use particles $\left\{z^{(i)}\right\}_{i=1}^{N}$ to represent $q(z)$.

$$V := \mathrm{grad}_q \, \mathrm{KL}(q, p) = \nabla \log(q/p).$$

$$z^{(i)} \leftarrow z^{(i)} + \varepsilon V\big(z^{(i)}\big).$$

$$= \textstyle\sum_j \big(z^{(i)} - z^{(j)}\big) K_{ij}$$
for Gaussian Kernel:
Repulsive force!

$V\big(z^{(i)}\big) \approx$

- SVGD [LW16]: $\sum_j K_{ij} \, \nabla_{z^{(j)}} \log p\big(z^{(j)}\big|x\big) + \sum_j \nabla_{z^{(j)}} K_{ij}$.

- Blob [CZW+18]: $\nabla_{z^{(i)}} \log p\big(z^{(i)}\big|x\big) - \dfrac{\sum_j \nabla_{z^{(i)}} K_{ij}}{\sum_k K_{ik}} - \sum_j \dfrac{\nabla_{z^{(i)}} K_{ij}}{\sum_k K_{jk}}$.

- GFSD [LZC+19]: $\nabla_{z^{(i)}} \log p\big(z^{(i)}\big|x\big) - \dfrac{\sum_j \nabla_{z^{(i)}} K_{ij}}{\sum_k K_{ik}}$.

- GFSF [LZC+19]: $\nabla_{z^{(i)}} \log p\big(z^{(i)}\big|x\big) + \sum_{j,k} (K^{-1})_{ik} \nabla_{z^{(j)}} K_{kj}$.

# Bayesian Inference: Variational Inference

- **Particle-based variational inference**: use particles $\left\{z^{(i)}\right\}_{i=1}^{N}$ to represent $q(z)$.
  - Unified view as Wasserstein gradient flow: [LZC+19].
  - Asymptotic analysis: SVGD [Liu17] ($N \to \infty, \varepsilon \to 0$).
  - Non-asymptotic analysis
    - w.r.t $\varepsilon$: e.g., [RT96] (as WGF).
    - w.r.t $N$: [CMG+18, FCSS18, ZZC18].
  - Accelerating ParVIs: [LZC+19, LZZ19].
  - Add particles dynamically: [CMG+18, FCSS18].
  - Solve the Wasserstein gradient by optimal transport: [CZ17, CZW+18].
  - Manifold support space: [LZ18].

# Outline

- Generative Models: Overview
- Plain Generative Models
  - Autoregressive Models
- Latent Variable Models
  - Deterministic Generative Models
    - Generative Adversarial Nets
    - Flow-Based Models
  - Probabilistic Graphical Models
    - Directed PGMs
      - **Bayesian Inference** (variational inference, **MCMC**)
      - Topic models (LDA, LightLDA, sLDA)
      - Deep Bayesian Models (VAE)
    - Undirected PGMs (Boltzmann machines, energy-based models)

# Bayesian Inference: MCMC

- Monte Carlo
  - Directly draw (i.i.d.) samples from $p(z|x)$.
  - Almost always impossible to directly do so (esp. w/ unnormalized $p(z|x)$).
- Markov Chain Monte Carlo (MCMC):
  Simulate a Markov chain whose stationary distribution is $p(z|x)$.
  - Easier to implement: only requires unnormalized $p(z|x)$ (*e.g.*, $p(z, x)$).
  - Asymptotically accurate.
  - Drawback/Challenge: sample auto-correlation.
    Less effective than i.i.d. samples.

[GC11]

# Bayesian Inference: MCMC

A fantastic MCMC animation site: https://chi-feng.github.io/mcmc-demo/

## The Markov-chain Monte Carlo Interactive Gallery

Click on an algorithm below to view interactive demo:

- Random Walk Metropolis Hastings
- Adaptive Metropolis Hastings [1]
- Hamiltonian Monte Carlo [2]
- No-U-Turn Sampler [2]
- Metropolis-adjusted Langevin Algorithm (MALA) [3]
- Hessian-Hamiltonian Monte Carlo (H2MC) [4]
- Stein Variational Gradient Descent (SVGD) [5]
- Nested Sampling with RadFriends (RadFriends-NS) [6]

View the source code on github: https://github.com/chi-feng/mcmc-demo.

# Bayesian Inference: MCMC

Classical MCMC

- Metropolis-Hastings framework [MRR+53, Has70]:

    Draw $z^* \sim q\left(z^*|z^{(k)}\right)$ and take $z^{(k+1)}$ as $z^*$ with probability

    $$\min\left\{1, \frac{q\left(z^{(k)}|z^*\right)p(z^*|x)}{q\left(z^*|z^{(k)}\right)p\left(z^{(k)}|x\right)}\right\},$$

    else take $z^{(k+1)}$ as $z^{(k)}$.

    - Note that $\dfrac{p(z^*|x)}{p\left(z^{(k)}|x\right)} = \dfrac{p(z^*,x)}{p(z^{(k)},x)}$ can be evaluated.

    - Proposal distribution $q(z^*|z)$: e.g., taken as $\mathcal{N}(z^*|z, \sigma^2)$.

# Bayesian Inference: MCMC

Classical MCMC

- Gibbs sampling [GG87]:

  Iteratively sample from conditional distributions, which are easier to draw:

  $$z_1^{(1)} \sim p\left(z_1 \mid z_2^{(0)}, z_3^{(0)}, \ldots, z_d^{(0)}, x\right),$$

  $$z_2^{(1)} \sim p\left(z_2 \mid z_1^{(1)}, z_3^{(0)}, \ldots, z_d^{(0)}, x\right),$$

  $$z_3^{(1)} \sim p\left(z_3 \mid z_1^{(1)}, z_2^{(1)}, \ldots, z_d^{(0)}, x\right),$$

  $$\ldots,$$

  $$z_i^{(k+1)} \sim p\left(z_i \mid z_1^{(k+1)}, \ldots, z_{i-1}^{(k+1)}, z_{i+1}^{(k)}, \ldots, z_d^{(k)}, x\right).$$

# Bayesian Inference: MCMC

Dynamics-based MCMC

- Simulates a jump-free continuous-time Markov process (dynamics):

$$\mathrm{d}z = \underbrace{b(z)\,\mathrm{d}t}_{\text{drift}} + \underbrace{\sqrt{2D(z)}\,\mathrm{d}B_t(z)}_{\text{diffusion}},$$

Pos. semi-def. matrix

Brownian motion

$$\Delta z = b(z)\varepsilon + \mathcal{N}(0, 2D(z)\varepsilon) + o(\varepsilon),$$

with appropriate $b(z)$ and $D(z)$ so that $p(z|x)$ is kept stationary/invariant.

- Informative transition using gradient $\nabla_z \log p(z|x)$.

- Some are compatible with *stochastic gradient* (SG): more efficient.

$$\nabla_z \log p(z|x) = \nabla_z \log p(z) + \sum_{n \in \mathcal{D}} \nabla_z \log p\big(x^{(n)}\big|z\big),$$

$$\widetilde{\nabla}_z \log p(z|x) = \nabla_z \log p(z) + \frac{|\mathcal{D}|}{|\mathcal{S}|} \sum_{n \in \mathcal{S}} \nabla_z \log p\big(x^{(n)}\big|z\big), \mathcal{S} \subset \mathcal{D}.$$

# Bayesian Inference: MCMC

Dynamics-based MCMC

- Langevin Dynamics [RS02] (compatible with SG [WT11, CDC15, TTV16]):
$$z^{(k+1)} = z^{(k)} + \varepsilon \nabla \log p\big(z^{(k)}\big|x\big) + \mathcal{N}(0, 2\varepsilon).$$

- Hamiltonian Monte Carlo [DKPR87, Nea11, Bet17]

  (*incompatible* with SG [CFG14, Bet15]; leap-frog integrator [CDC15]):
$$r^{(0)} \sim \mathcal{N}(0, \Sigma), \quad \begin{cases} r^{(k+1/2)} = r^{(k)} + (\varepsilon/2)\nabla \log p\big(z^{(k)}\big|x\big), \\ z^{(k+1)} = z^{(k)} + \varepsilon \Sigma^{-1} r^{(k+1/2)}, \\ r^{(k+1)} = r^{(k+1/2)} + (\varepsilon/2)\nabla \log p\big(z^{(k+1)}\big|x\big). \end{cases}$$

- Stochastic Gradient Hamiltonian Monte Carlo [CFG14] (compatible with SG):
$$\begin{cases} z^{(k+1)} = z^{(k)} + \varepsilon \Sigma^{-1} r^{(k)}, \\ r^{(k+1)} = r^{(k)} + \varepsilon \nabla \log p\big(z^{(k)}\big|x\big) - \varepsilon C \Sigma^{-1} r^{(k)} + \mathcal{N}(0, 2C\varepsilon). \end{cases}$$

- ...

# Bayesian Inference: MCMC

Dynamics-based MCMC

- Complete framework for MCMC dynamics: [MCF15].
- Interpretation on the Wasserstein space: [JKO98, LZZ19].
- Integrators and their non-asymptotic analysis (with SG): [CDC15].
- For manifold support space:
  - LD: [GC11]; HMC: [GC11, BSU12, BG13, LSSG15]; SGLD: [PT13]; SGHMC: [MCF15, LZS16]; SGNHT: [LZS16]
- Different kinetic energy (other than Gaussian):
  - Monomial Gamma [ZWC+16, ZCG+17].
- Fancy Dynamics:
  - Relativistic: [LPH+16]
  - Magnetic: [TRGT17]

# Bayesian Inference: Comparison

| | **Parametric VI** | **Particle-Based VI** | **MCMC** |
|---|---|---|---|
| Asymptotic Accuracy | No | Yes | Yes |
| Approximation Flexibility | Limited | Unlimited | Unlimited |
| Empirical Convergence Speed | High | High | Low |
| Particle Efficiency | (Do not apply) | High | Low |
| High-Dimensional Efficiency | High | Low | High |

# Outline

- Generative Models: Overview
- Plain Generative Models
  - Autoregressive Models
- Latent Variable Models
  - Deterministic Generative Models
    - Generative Adversarial Nets
    - Flow-Based Models
  - Probabilistic Graphical Models
    - Directed PGMs
      - Bayesian Inference (variational inference, MCMC)
      - **Topic models (LDA, LightLDA, sLDA)**
      - Deep Bayesian Models (VAE)
    - Undirected PGMs (Boltzmann machines, energy-based models)
    - Diffusion-Based Models

# Topic Models

Separate *global* (dataset abstraction) and *local* (datum representation) latent variables.



[SG07]

# Latent Dirichlet Allocation

Model Structure [BNJ03]:



- Data variable: Words/Documents $w = \{w_{dn}\}_{n=1:N_d, d=1:D}, w_{dn} \in \{1 \ldots W\}$.

- Latent variables:
  - *Global*:      topics $\beta = \{\beta_k\}_{k=1:K}, \beta_k \in \Delta^W$.
  - *Local*:      topic proportions $\theta = \{\theta_d\}, \theta_d \in \Delta^K$,
    
               topic assignments $z = \{z_{dn}\}, z_{dn} \in \{1 \ldots K\}$.

- Prior: $p(\beta_k|b) = \text{Dir}(b), p(\theta_d|a) = \text{Dir}(a), p(z_{dn}|\theta_d) = \text{Mult}(\theta_d)$.

- Likelihood: $p(w_{dn}|z_{dn}, \beta) = \text{Mult}\left(\beta_{z_{dn}}\right)$.

# Latent Dirichlet Allocation

Variational inference [BNJ03]:

- Take variational distribution (mean-field approximation):

$$q_{\lambda,\gamma,\phi}(\beta,\theta,z) := \prod_{k=1}^{K} \mathrm{Dir}(\beta_k|\lambda_k) \prod_{d=1}^{D} \mathrm{Dir}(\theta_d|\gamma_d) \prod_{n=1}^{N_d} \mathrm{Mult}(z_{dn}|\phi_{dn}).$$

- ELBO$(\lambda,\gamma,\phi;a,b)$ is available in <span style="color:red">closed form</span>.

- E-step: update $\lambda,\gamma,\phi$ by maximizing ELBO;

- M-step: update $a,b$ by maximizing ELBO.

# Latent Dirichlet Allocation

MCMC: Gibbs sampling [GS04]

Model structure $\implies p(\beta, \theta, z, w) = AB\left(\prod_{k,w} \beta_{kw}^{N_{kw}+b_w-1}\right)\left(\prod_{d,k} \theta_{dk}^{N_{kd}+a_k-1}\right)$

$\implies p(z,w) = AB\left(\prod_k \frac{\prod_w \Gamma(N_{kw}+b_w)}{\Gamma(N_k+W\bar{b})}\right)\left(\prod_d \frac{\prod_k \Gamma(N_{kd}+a_k)}{\Gamma(N_d+K\bar{a})}\right).$

($N_{kw}$: #times word $w$ is assigned to topic $k$; $N_{kd}$: #times topic $k$ appears in document $d$.)

- Unacceptable cost to directly compute $p(z|w) = p(z,w)/p(w)$.

- Use Gibbs sampling to draw from $p(z|w)$!

$$p(z_{dn} = k | z^{-dn}, w) \propto \frac{N_{kw}^{-dn} + b_w}{N_k^{-dn} + W\bar{b}}\left(N_{kd}^{-dn} + a_k\right).$$

- For $\beta$ and $\theta$, use MAP estimate:

$$\hat{\beta} := \arg\max_\beta \log p(\beta|w) \approx \frac{N_{kw} + b_w}{N_k + W\bar{b}},$$

$$\hat{\theta}_{dk} := \arg\max_\theta \log p(\theta|w) \approx \frac{N_{kd} + a_k}{N_d + K\bar{a}}.$$

Estimated by samples of $z$

# Latent Dirichlet Allocation

MCMC: LightLDA [YGH+15]

$$p(z_{dn} = k | z^{-dn}, w) \propto (N_{kd}^{-dn} + a_k) \frac{N_{kw}^{-dn} + b_w}{N_k^{-dn} + W\overline{b}}.$$

- Direct implementation: $O(K)$ time.
- Amortized $O(1)$ multinomial sampling: alias table.



$$\left[\frac{3}{8}, \frac{1}{16}, \frac{1}{8}, \frac{7}{16}\right] \Rightarrow \text{Alias Table:} \left[\left(4, \frac{3}{16}\right), \left(1, \frac{1}{16}\right), \left(4, \frac{1}{8}\right), \left(4, \frac{1}{4}\right)\right] = [(h_i, v_i)]$$

- $O(1)$ sampling: $i \sim \text{Unif}\{1, \dots, K\}, v \sim \text{Unif}[0,1], z = i$ if $v < v_i$ else $h_i$.
- $O(K)$ time to build the Alias Table $\Rightarrow$ Amortized $O(1)$ time for $K$ samples.
- What if the target changes (slightly): use Metropolis Hastings (MH) to correct.

# Latent Dirichlet Allocation

- Dynamics-Based MCMC and Particle-Based VI: target $p(\beta|w)$.

$$\nabla_\beta \log p(\beta|w) = \mathbb{E}_{p(z|\beta,w)}\big[\nabla_\beta \log p(\beta,z,w)\big].$$

| Gibbs Sampling | Closed-form known |

- Stochastic Gradient Riemannian Langevin Dynamics [PT13], Stochastic Gradient Nose-Hoover Thermostats [DFB+14], Stochastic Gradient Riemannian Hamiltonian Monte Carlo [MCF15].

- Accelerated particle-based VI [LZC+19, LZZ19].

# Supervised Latent Dirichlet Allocation

Model structure [MB08]:



$z$: topics

| "ENGINES" | speed | product | introduced |
|---|---|---|---|
| "ROYAL" | britain | queen | sir |
| "ARMY" | commander | forces | war |
| "STUDY" | analysis | space | program |

$x_1$: doc 1

$y_1$: science & tech

$x_2$: doc 2

$y_2$: politics

- Variational inference: similar to LDA.
- Prediction: for test document $w_d$,

$$\hat{y}_d := \mathbb{E}_{p(y_d|w_d)}[y_d] = \eta^\top \mathbb{E}_{p(z_d|w_d)}[\bar{z}_d]$$
$$\approx \eta^\top \mathbb{E}_{q(z_d|w_d)}[\bar{z}_d].$$

First do inference (find $q(z_d|w_d)$), then estimate $\hat{y}_d$.

# Supervised Latent Dirichlet Allocation

Variational inference with posterior regularization [ZAX12]

- Regularized Bayes (RegBayes) [ZCX14]:
  - Recall: $p\left(z\middle|\left\{x^{(n)}, y^{(n)}\right\}\right)$
    $= \arg\min_{q(z)}\left\{-\mathcal{L}[q] = \mathrm{KL}\left(q(z), p(z)\right) - \sum_n \mathbb{E}_q\left[\log p\left(x^{(n)}, y^{(n)}|z\right)\right]\right\}.$
  - Regularize posterior towards better prediction:
    $\min_{q(z)} \mathrm{KL}\left(q(z), p(z)\right) - \sum_n \mathbb{E}_q\left[\log p\left(x^{(n)}, y^{(n)}|z\right)\right] + \lambda\ell\left(q(z); \left\{x^{(n)}, y^{(n)}\right\}\right).$
- Maximum entropy discrimination LDA (MedLDA) [ZAX12]:
  - $\ell\left(q; \left\{w^{(n)}, y^{(n)}\right\}\right) = \sum_n \ell_\varepsilon\left(y^{(n)} - \hat{y}^{(n)}\left(q, w^{(n)}\right)\right)$
    $= \sum_n \ell_\varepsilon\left(y^{(n)} - \eta^\top \mathbb{E}_{q\left(z^{(n)}|w^{(n)}\right)}\left[\bar{z}^{(n)}\right]\right),$
    where $\ell_\varepsilon(r) = \max\{0, |r| - \varepsilon\}$ is the hinge (max-margin) loss.
  - Facilitates both prediction and topic representation.

# Outline

- Generative Models: Overview
- Plain Generative Models
  - Autoregressive Models
- Latent Variable Models
  - Deterministic Generative Models
    - Generative Adversarial Nets
    - Flow-Based Models
  - Probabilistic Graphical Models
    - Directed PGMs
      - Bayesian Inference (variational inference, MCMC)
      - Topic models (LDA, LightLDA, sLDA)
      - **Deep Bayesian Models (VAE)**
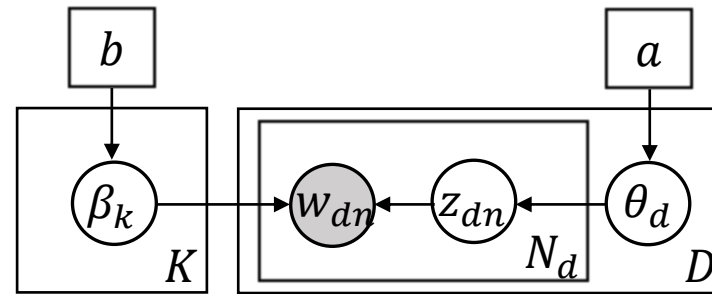    - Undirected PGMs (Boltzmann machines, energy-based models)
    - Diffusion-Based Models

# Variational Auto-Encoder

More *flexible* Bayesian model using *deep learning* tools.

- Model structure (decoder) [KW14]:

$$z_d \sim p(z_d) = \mathcal{N}(z_d|0, I),$$

$$x_d \sim p_\theta(x_d|z_d) = \mathcal{N}\big(x_d|\mu_\theta(z_d), \Sigma_\theta(z_d)\big),$$

where $\mu_\theta(z_d)$ and $\Sigma_\theta(z_d)$ are modeled by neural networks.

# Variational Auto-Encoder

- Variational inference (encoder) [KW14]:

$$q_\phi(z|x) := \prod_{d=1}^{D} q_\phi(z_d|x_d) = \prod_{d=1}^{D} \mathcal{N}\big(z_d\big|\nu_\phi(x_d), \Gamma_\phi(x_d)\big),$$

where $\nu_\phi(x_d), \Gamma_\phi(x_d)$ are also NNs.

- Amortized inference: to approximate local posteriors $\{p(z_d|x_d)\}_{d=1}^{D}$,
  - instead of using $q_{\phi_d}(z_d)$ for each $p(z_d|x_d)$ and learning *local* parameters $\{\phi_d\}$ (like LDA),
  - use $q_\phi(z_d|x_d)$ and learn the *global* parameter $\boldsymbol{\phi}$ (**fast inference for unseen $\boldsymbol{x_d}$**).

- Objective: $\mathbb{E}_{\hat{p}(x_d)}[\log p_\theta(x_d)] \geq \mathbb{E}_{\hat{p}(x_d)}[\mathrm{ELBO}(x_d)]$,

$$\mathrm{ELBO}(x_d) = \mathbb{E}_{q_\phi(z_d|x_d)}\big[\log p_\theta(z_d)p_\theta(x_d|z_d) - \log q_\phi(z_d|x_d)\big].$$



$p(z_d)$

$p_\theta(x_d|z_d)$

$q_\phi(z_d|x_d)$

# Variational Auto-Encoder

- Variational inference (encoder) [KW14]:

$$q_\phi(z|x) := \prod_{d=1}^D q_\phi(z_d|x_d) = \prod_{d=1}^D \mathcal{N}\big(z_d\big|\nu_\phi(x_d), \Gamma_\phi(x_d)\big),$$

where $\nu_\phi(x_d), \Gamma_\phi(x_d)$ are also NNs.

$$\text{ELBO}(x_d) = \mathbb{E}_{q_\phi(z_d|x_d)}\big[\log p_\theta(z_d)p_\theta(x_d|z_d) - \log q_\phi(z_d|x_d)\big].$$

- Gradient estimation with the *reparameterization trick*:

$$z_d \sim q_\phi(z_d|x_d) \iff z_d = g_\phi(x_d, \epsilon) := \nu_\phi(x_d) + \epsilon\sqrt{\Gamma_\phi(x_d)}, \epsilon \sim q(\epsilon) := \mathcal{N}(\epsilon|0, I).$$

  - Gradient estimation: $\nabla_{\phi,\theta}\text{ELBO}(x_d) =$
  
  $$\mathbb{E}_{q(\epsilon)}\Big[\nabla_{\phi,\theta}\big(\log p_\theta\big(g_\phi(x_d,\epsilon)\big)p_\theta\big(x_d|g_\phi(x_d,\epsilon)\big) - \log q_\phi\big(g_\phi(x_d,\epsilon)|x_d\big)\big)\Big].$$

  - Smaller variance than REINFORCE-like estimator [Wil92]:
  
  $$\nabla_\phi \mathbb{E}_{q_\phi}[f_\phi] = \mathbb{E}_{q_\phi}\big[\nabla_\phi f_\phi + f_\phi \nabla_\phi \log q_\phi\big].$$



$\phi \rightarrow z_d$   $p(z_d)$

$q_\phi(z_d|x_d)$    $p_\theta(x_d|z_d)$

$x_d \leftarrow \theta$

$D$

# Variational Auto-Encoder

- Inference with importance-weighted ELBO [BGS15]
  - Conventional ELBO (subscript $d$ omitted):

  $$\mathcal{L}_\theta[q_\phi](x) := \mathbb{E}_{q_\phi(z|x)}\left[\log\frac{p_\theta(z,x)}{q_\phi(z|x)}\right].$$

  - A tighter lower bound:

  $$\mathcal{L}_\theta^{(K)}[q_\phi](x) := \mathbb{E}_{z^{(1)},\ldots,z^{(K)}\sim\text{i.i.d.}\,q_\phi}\left[\log\frac{1}{K}\sum_{i=1}^K\frac{p_\theta(z^{(i)},x)}{q_\phi(z^{(i)}|x)}\right].$$

  Ordering relation:

  $$\mathcal{L}_\theta[q_\phi](x) = \mathcal{L}_\theta^{(1)}[q_\phi](x) \leq \mathcal{L}_\theta^{(2)}[q_\phi](x) \leq \cdots \leq \mathcal{L}_\theta^{(\infty)}[q_\phi](x) = \log p_\theta(x).$$

  - SUMO [LBN+19]: unbiased estimate of $\mathcal{L}_\theta^{(\infty)}[q_\phi](x)$.

  If $\frac{p(z,x)}{q(z|x)}$ is bounded.

# Variational Auto-Encoder



- Semi-supervised VAE [KMRW14, M2]
  - For labeled data:
    - Required encoder: $q_\phi(z_d|x_d, y_d)$.
    - Objective: $\mathbb{E}_{\hat{p}(x_d, y_d)}[\log p_\theta(x_d, y_d)] \geq \mathbb{E}_{\hat{p}(x_d, y_d)}[\mathrm{ELBO}(x_d, y_d)]$,

$$\mathrm{ELBO}(x_d, y_d) = \mathbb{E}_{q_\phi(z_d|x_d, y_d)}[\log p_\theta(z_d) p_\theta(y_d) p_\theta(x_d|z_d, y_d) - \log q_\phi(z_d|x_d, y_d)].$$

  - For unlabeled data:
    - Required encoder: $q_\phi(y_d, z_d|x_d) = q_\phi(y_d|x_d) q_\phi(z_d|x_d, y_d)$.
    - Objective: $\mathbb{E}_{\hat{p}(x_d)}[\log p_\theta(x_d)] \geq \mathbb{E}_{\hat{p}(x_d)}[\mathrm{ELBO}(x_d)]$,

$$\mathrm{ELBO}(x_d) = \mathbb{E}_{q_\phi(y_d, z_d|x_d)}[\log p_\theta(z_d) p_\theta(y_d) p_\theta(x_d|z_d, y_d) - \log q_\phi(y_d, z_d|x_d)]$$
$$= \mathbb{E}_{q_\phi(y_d|x_d)}[\mathrm{ELBO}(x_d, y_d) - \log q_\phi(y_d|x_d)].$$

  - For prediction: use $q_\phi(y_d|x_d)$.

# Variational Auto-Encoder



- Conditional VAE [SYL15]
  - Let the generation of $(z_d, y_d)$ conditioned on $x_d$ (so it is not generative).
  - Model: $p_\theta(z_d, y_d | x_d) = p_\theta(z_d | x_d) p(y_d | x_d, z_d)$.
  - Required encoder: $q_\phi(z_d | x_d, y_d)$.
  - Objective: $\mathrm{ELBO}(y_d | x_d) = \mathbb{E}_{q_\phi(z_d | x_d, y_d)} \big[ \log p_\theta(z_d | x_d) p(y_d | x_d, z_d) - \log q_\phi(z_d | x_d, y_d) \big]$.
  - Prediction: ancestral sampling: $z_d \sim p_\theta(z_d | x_d), y_d \sim p(y_d | x_d, z_d)$.
- VAE with structured prior
  - [LWZZ18] mixture of Gaussian, state-space model.
  - [KSDV18] Causal network.
  - [PHN+20] energy-based prior.

# Variational Auto-Encoder

- Learning disentangled representation
  - InfoGAN [CDH+16]: max mutual_info(part_of_$z$, generated_$x$).
  - $\beta$-VAE [HLP+17]: upscale the KL term ($q(z|x)$ to factorized prior $p(z)$) in ELBO.
  - Total Correlation VAE [CLG+18]: upscale the total-correlation term in a finer decomposition of ELBO.



(a) Varying $c_1$ on InfoGAN (Digit type)

(c) Varying $c_2$ from $-2$ to $2$ on InfoGAN (Rotation)

(d) Varying $c_3$ from $-2$ to $2$ on InfoGAN (Width)

# Variational Auto-Encoder

- Learning disentangled representation
  - Formal definition [HAP+18] (roughly): a class of transformations on $x$ (holding some semantics) changes only one dimension of the representation.
  - Impossibility theorem [LBL+19]:

**Theorem 1.** *For $d > 1$, let $\mathbf{z} \sim P$ denote any distribution which admits a density $p(\mathbf{z}) = \prod_{i=1}^{d} p(\mathbf{z}_i)$. Then, there exists an infinite family of bijective functions $f : \mathrm{supp}(\mathbf{z}) \to \mathrm{supp}(\mathbf{z})$ such that $\frac{\partial f_i(\boldsymbol{u})}{\partial u_j} \neq 0$ almost everywhere for all $i$ and $j$ (i.e., $\mathbf{z}$ and $f(\mathbf{z})$ are completely entangled) and $P(\mathbf{z} \leq \boldsymbol{u}) = P(f(\mathbf{z}) \leq \boldsymbol{u})$ for all $\boldsymbol{u} \in \mathrm{supp}(\mathbf{z})$ (i.e., they have the same marginal distribution).*

  - Works afterwards:
    - Weak supervision: a few labels [LTB+19], pairwise similarity [CB20], paired unsupervised data [LPR+20], rank pairing [SCK+20].
    - If the cause of $z$ is observed, $z$'s suff. stat. can be **identified** up to a permutation [KKM+20].

# Variational Auto-Encoder

Train:     Test: 

- Learning causal representation.
  - Why: Causal relations tend to hold across domains [SJP+12, PJS17, Sch19].
  - Invariance risk min. [ABGL19]: Optimal representation-based classifier is invariant.
  - Causal generative model [LSW+21] (single training domain; [SWZ+21] for multiple tr. dom.):



$p(s,v)$    $q(s,v|x)$    $p(y|s)$    $p(x|s,v)$

**Model**:
- Generative process is more likely causal/invariant than inference process.
- Domain shift comes from the change of **prior** (repr. distr.).
- Not all representation *causes* $y$ ➡ the $s$emantic-$v$ariation split.

  - **Prediction**: use an *independent prior* (if no test data) or a *newly learned prior* (unlabeled test data).
  - **Learning**: using the test-domain inf. model $q^{\perp}(s,v|x)$ or $\tilde{q}(s,v|x)$ suffices.
  - **Theory**: under certain conditions,
    a well-learned model **identifies** the semantics $s$,
    and the test-domain/out-of-distr. prediction error
    is bounded (no test data) or vanishes (unlabeled test data).

$p^{\perp}(s,v) \coloneqq p(s)p(v)$    $q^{\perp}(s,v|x)$     $\tilde{p}(s,v)$    $\tilde{q}(s,v|x)$

# Variational Auto-Encoder

- Bidirectional/Prior-free generative modeling [LTQ+21]
  - Modeling $p(x, z)$ by specifying a prior $p(z)$:
  
  (1) Hard inference. (2) Manifold mismatch.　　(3) Posterior collapse.

  true
  data
  distr.

  learned
  data
  distr.

  class-wise posterior samples

  - Thm (informal): Conditional densities $p(x|z)$, $q(z|x)$ come from a common joint $p(x, z)$ (*compatible*), **iff.** $\frac{p(x|z)}{q(z|x)}$ factorizes as $a(x)b(z)$ on a certain region they determine. Such $p(x, z)$ is unique on the region (*determinacy*).
    - For $p(x|z) = \delta_{f(z)}(x)$, insufficient determinacy (compatible $\Leftrightarrow \exists x_0$ s.t. $q(f^{-1}(\{x_0\})|x_0) = 1$).
- Algorithms are possible!
  - Enforcing compatibility: $\min \mathbb{E}_{p^*(x)q_\phi(z|x)} \left\| \nabla_x \nabla_z^\top \log \left( p_\theta(x|z)/q_\phi(z|x) \right) \right\|_F^2$.
  - Data-fitting: MLE: $\mathbb{E}_{p^*(x)}\left[ \log p_{\theta,\phi}(x) \right] = \mathbb{E}_{p^*(x)}\left[ -\log \mathbb{E}_{q_\phi(z'|x)}[1/p_\theta(x|z')] \right]$.
  - Data gen.: MCMC: $\Delta x^{(t)} = \varepsilon \nabla_{x^{(t)}} \log \frac{p_\theta(x^{(t)}|z^{(t)})}{q_\phi(z^{(t)}|x^{(t)})} + \sqrt{2\varepsilon}\, \eta^{(t)}$, where $z^{(t)} \sim q_\phi(z|x^{(t)})$, $\eta^{(t)} \sim \mathcal{N}(0, I)$.

$p(z)$ ? $z$

$q_\phi(z|x)$ !

$p_\theta(x|z)$

$x$

# Variational Auto-Encoder

- Parametric Variational Inference: towards more flexible approximations.
  - Explicit VI:

    Normalizing flows [RM15, KSJ+16].
  - Implicit VI:

    Adversarial Auto-Encoder [MSJ+15], Adversarial Variational Bayes [MNG17], Wasserstein Auto-Encoder [TBGS17], [SSZ18a], [LT18], [SSZ18b].


- MCMC [LTL17] and Particle-Based VI [FWL17, PGH+17]:
  - Train the encoder as a sample generator.
  - Amortize the update on samples to $\phi$.

# Outline

- Generative Models: Overview
- Plain Generative Models
  - Autoregressive Models
- Latent Variable Models
  - Deterministic Generative Models
    - Generative Adversarial Nets
    - Flow-Based Models
  - Probabilistic Graphical Models
    - Directed PGMs
      - Bayesian Inference (variational inference, MCMC)
      - Topic models (LDA, LightLDA, sLDA)
      - Deep Bayesian Models (VAE)
    - **Undirected PGMs (Boltzmann machines, energy-based models)**
    - Diffusion-Based Models

# Undirected PGMs

Specify $p_\theta(x,z)$ by an energy function $E_\theta(x,z)$:

$$p_\theta(x,z) = \frac{1}{Z_\theta}\exp\big(-E_\theta(x,z)\big), Z_\theta = \int \exp\big(-E_\theta(x',z')\big)\ \mathrm{d}x'\mathrm{d}z'.$$

$p_\theta(x,z) \propto \exp\big(-E_\theta(x,z)\big)$

- Only correlation and no causality: $p(x,z)$ is either $p(z)p(x|z)$ or $p(x)p(z|x)$.

+ Flexible and simple in modeling dependency.

- Harder to learn and generate than directed PGMs.

  =0 if $E = \log p$.

  - Learning: even $p_\theta(x,z)$ is unavailable.

    $$\nabla_\theta \mathbb{E}_{\hat{p}(x)}[\log p_\theta(x)] = -\mathbb{E}_{\hat{p}(x)p_\theta(z|x)}[\nabla_\theta E_\theta(x,z)] + \mathbb{E}_{p_\theta(x,z)}[\nabla_\theta E_\theta(x,z)].$$

    (augmented) data distribution
    (Bayesian inference)
    model distribution
    (generation)

  - Bayesian inference: generally same as directed PGMs.

  - Generation: rely on MCMC or training a generator.

# Undirected PGMs

- Learning: $\nabla_\theta \mathbb{E}_{\hat{p}(x)}[\log p_\theta(x)] = -\mathbb{E}_{\hat{p}(x)\,\color{green}{p_\theta(z|x)}}[\nabla_\theta E_\theta(x,z)] + \mathbb{E}_{\color{blue}{p_\theta(x,z)}}[\nabla_\theta E_\theta(x,z)].$

<span style="color:green">↑<br>Bayesian Inference</span>    <span style="color:blue">↑<br>Generation</span>

- Boltzmann Machine: Gibbs sampling for both inference and generation [HS83].



$$E_\theta(x,z) = -x^\top W z - \tfrac{1}{2} x^\top L x - \tfrac{1}{2} z^\top J z.$$
$$\Longrightarrow$$
$$p_\theta(z_j|x, z_{-j}) = \mathrm{Bern}\left(\sigma\left(\textstyle\sum_{i=1}^{D} W_{ij} x_i + \sum_{m\neq j}^{P} J_{jm} z_j\right)\right),$$
$$p_\theta(x_i|z, x_{-i}) = \mathrm{Bern}\left(\sigma\left(\textstyle\sum_{j=1}^{P} W_{ij} z_j + \sum_{k\neq i}^{D} L_{ik} x_k\right)\right).$$

# Undirected PGMs

- Learning: $\nabla_\theta \mathbb{E}_{\hat{p}(x)}[\log p_\theta(x)] = -\mathbb{E}_{\hat{p}(x)\, \textcolor{green}{p_\theta(z|x)}}[\nabla_\theta E_\theta(x,z)] + \mathbb{E}_{\textcolor{blue}{p_\theta(x,z)}}[\nabla_\theta E_\theta(x,z)].$

<div align="center">

↑ <span style="color:green">Bayesian Inference</span>　　　　↑ <span style="color:blue">Generation</span>

</div>

- Restricted Boltzmann Machine [Smo86]:



$$E_\theta(x,z) = -x^\top W z + b^{(x)^\top} x + b^{(z)^\top} z.$$

- Bayesian Inference is exact:

$$p_\theta(z_k|x) = \mathrm{Bern}\left(\sigma\left(x^\top W_{:k} + b_k^{(z)}\right)\right).$$

- Generation: Gibbs sampling. Iterate:

$$p_\theta(z_k|x) = \mathrm{Bern}\left(\sigma\left(x^\top W_{:k} + b_k^{(z)}\right)\right),$$

$$p_\theta(x_k|z) = \mathrm{Bern}\left(\sigma\left(W_{k:}z + b_k^{(x)}\right)\right).$$

# Undirected PGMs

- Learning: $\nabla_\theta \mathbb{E}_{\hat{p}(x)}[\log p_\theta(x)] = -\mathbb{E}_{\hat{p}(x) p_\theta(z|x)}[\nabla_\theta E_\theta(x,z)] + \mathbb{E}_{p_\theta(x,z)}[\nabla_\theta E_\theta(x,z)].$

<span style="color:green">Bayesian Inference</span>      <span style="color:blue">Generation</span>

- Deep Belief Network [HOT06]
  (hybrid of directed and undirected)

- Deep Boltzmann Machine [SH09]



$p(v, h^{(1)}, \dots, h^{(L)})$
$= p(v|h^{(1)})p(h^{(1)}|h^{(2)}) \dots p(h^{(L-2)}|h^{(L-1)})p(h^{(L-1)}, h^{(L)}).$

$E_\theta(v, h^{(1)}, \dots, h^{(L)})$
$= E_{W^{(1)}}(v, h^{(1)}) + \sum_{l=2}^{L} E_{W^{(l)}}(h^{(l-1)}, h^{(l)}).$

# Undirected PGMs

- Learning: $\nabla_\theta \mathbb{E}_{\hat{p}(x)}[\log p_\theta(x)] = -\mathbb{E}_{\hat{p}(x)\color{green}{p_\theta(z|x)}}[\nabla_\theta E_\theta(x,z)] + \mathbb{E}_{\color{blue}{p_\theta(x,z)}}[\nabla_\theta E_\theta(x,z)].$

<span style="color:green">↑<br>Bayesian Inference</span>           <span style="color:blue">↑<br>Generation</span>

- [Hin02]: estimation with $k$-step MCMC approximates the gradient of $k$-step *Contrastive Divergence* (CD-$k$):

$$\mathrm{CD}_k := \mathrm{KL}(P^0 || P_\theta^\infty) - \mathrm{KL}(P_\theta^k || P_\theta^\infty),$$
$$P^0(x) = \hat{p}(x), P_\theta^k(x) := P^0(x){\color{red}P_\theta(x^{(k)}|x)}.$$

<span style="color:red">$k$-step transition of MCMC from data to model.</span>

# Undirected PGMs

Deep Energy-Based Models:

No latent variable; $E_\theta(x)$ is modeled by a neural network.

$$\nabla_\theta \mathbb{E}_{\hat{p}(x)}[\log p_\theta(x)] = -\mathbb{E}_{\hat{p}(x)}[\nabla_\theta E_\theta(x)] + \mathbb{E}_{p_\theta(x')}[\nabla_\theta E_\theta(x')].$$

- [KB16]: learn a generator

$$x \sim q_\phi(x) \Longleftrightarrow z \sim q(z), x = g_\phi(z),$$

to mimic the generation from $p_\theta(x)$:

$$\arg\min_\phi \mathrm{KL}(q_\phi, p_\theta) = \arg\min_\phi \mathbb{E}_{q(z)}\left[E_\theta\left(g_\phi(z)\right)\right] - \underbrace{\mathbb{H}[q_\phi]}_{\text{approx. by batch normalization Gaussian}}$$

# Undirected PGMs

Deep Energy-Based Models:

No latent variable; $E_\theta(x)$ is modeled by a neural network.

$$\nabla_\theta \mathbb{E}_{\hat{p}(x)}[\log p_\theta(x)] = -\mathbb{E}_{\hat{p}(x)}[\nabla_\theta E_\theta(x)] + \mathbb{E}_{p_\theta(x')}[\nabla_\theta E_\theta(x')].$$
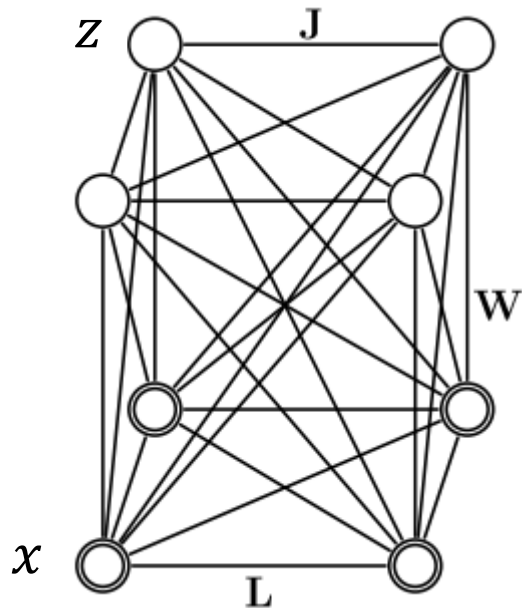
- [DM19]: estimate $\mathbb{E}_{p_\theta(x')}[\cdot]$ by samples drawn by the Langevin Dynamics

$$x^{(k+1)} = x^{(k)} - \varepsilon \nabla_x E_\theta(x^{(k)}) + \mathcal{N}(0, 2\varepsilon).$$

- Replay buffer for initializing the LD chain.
- $L_2$-regularization on the energy function.

# Undirected PGMs

**Deep Energy-Based Models:**

- [DM19]



ImageNet32x32 Generation

| Model | Inception | FID |
|---|---|---|
| **CIFAR-10 Unconditional** | | |
| PixelCNN (Van Oord et al., 2016) | 4.60 | 65.93 |
| PixelIQN (Ostrovski et al., 2018) | 5.29 | 49.46 |
| EBM (single) | 6.02 | 40.58 |
| DCGAN (Radford et al., 2016) | 6.40 | 37.11 |
| WGAN + GP (Gulrajani et al., 2017) | 6.50 | 36.4 |
| EBM (10 historical ensemble) | 6.78 | 38.2 |
| SNGAN (Miyato et al., 2018) | **8.22** | 21.7 |
| **CIFAR-10 Conditional** | | |
| Improved GAN | 8.09 | - |
| EBM (single) | 8.30 | 37.9 |
| Spectral Normalization GAN | **8.59** | 25.5 |
| **ImageNet 32x32 Conditional** | | |
| PixelCNN | 8.33 | 33.27 |
| PixelIQN | 10.18 | 22.99 |
| EBM (single) | **18.22** | **14.31** |
| **ImageNet 128x128 Conditional** | | |
| ACGAN (Odena et al., 2017) | 28.5 | - |
| EBM* (single) | 28.6 | 43.7 |
| SNGAN | **36.8** | **27.62** |

# Undirected PGMs

Deep Energy-Based Models:

- Score-based methods [Hyv05]:
  Learn $\mathbf{s}_\theta(\mathbf{x})$ (represents $\nabla_\mathbf{x} \log p_\theta(\mathbf{x}) = -\nabla_\mathbf{x} E_\theta(\mathbf{x})$) to approx $\nabla_\mathbf{x} \log p_{\text{data}}(\mathbf{x})$, by min:

$$\underbrace{\mathbb{E}_{p_{\text{data}}(\mathbf{x})} \|\mathbf{s}_\theta(\mathbf{x}) - \nabla_\mathbf{x} \log p_{\text{data}}(\mathbf{x})\|_2^2}_{\text{Fisher divergence } (p_\theta, p_{\text{data}})} = \underbrace{\mathbb{E}_{p_{\text{data}}(\mathbf{x})} [\|\mathbf{s}_\theta(\mathbf{x})\|_2^2 + 2\nabla \cdot \mathbf{s}_\theta(\mathbf{x})]}_{\text{score}-\text{matching objective}} + \text{const.},$$

  - The density $p_{\text{data}}(\mathbf{x})$ is not required! Estimate the expectation by sample average.
  - Data generation: run Langevin dynamics with $\mathbf{s}_\theta(\mathbf{x})$.

- Noise Annealing Score Matching [SE19]:
  - $p_{\text{data}}(\mathbf{x})$ may concentrate on a low-dim. manifold $\Longrightarrow \nabla_\mathbf{x} \log p_{\text{data}}(\mathbf{x})$ is ill-posed!
  - Perturb the data by noise with shrinking variance: avoid concentration on manifold.
    Consider $p_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) := \mathcal{N}(\tilde{\mathbf{x}}|\mathbf{x}, \sigma^2\mathbf{I})$, $p_\sigma(\tilde{\mathbf{x}}) := \int p_{\text{data}}(\mathbf{x}) p_\sigma(\tilde{\mathbf{x}}|\mathbf{x})\, d\mathbf{x}$,
    and $\sigma_{\max} = \sigma_1 > \cdots > \sigma_T = \sigma_{\min}$ s.t.: $p_{\sigma_{\min}}(\mathbf{x}) \approx p_{\text{data}}(\mathbf{x})$.

# Outline

- Generative Models: Overview
- Plain Generative Models
    - Autoregressive Models
- Latent Variable Models
    - Deterministic Generative Models
        - Generative Adversarial Nets
        - Flow-Based Models
    - Probabilistic Graphical Models
        - Directed PGMs
            - Bayesian Inference (variational inference, MCMC)
            - Topic models (LDA, LightLDA, sLDA)
            - Deep Bayesian Models (VAE)
        - Undirected PGMs (Boltzmann machines, energy-based models)
        - **Diffusion-Based Models**

# Diffusion-Based Models

[SWMG15, HJA20]



$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$
$$q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

- Gradually corrupting images into random noise is easy:

  Let $q(\mathbf{x}_t|\mathbf{x}_{t-1}) \coloneqq \mathcal{N}\left(\mathbf{x}_t; \sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t \boldsymbol{I}\right)$, $\mathbf{x}_0$ be the data variable.
  Then $q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1-\bar{\alpha}_t)\boldsymbol{I}\right)$, $\bar{\alpha}_t \coloneqq \prod_{s=1}^{t}(1-\beta_t)$.
  $q(\mathbf{x}_T|\mathbf{x}_0) \approx \mathcal{N}(\mathbf{x}_T; \boldsymbol{0}, \boldsymbol{I}) \eqqcolon p(\mathbf{x}_T)$ for large $T$!

> Mimics Langevin dynamics targeting std Gaussian:
> $x^{(t)} = x^{(t-1)} + \varepsilon\nabla\log p_{\mathcal{N}}\left(x^{(t-1)}\right) + \mathcal{N}(0,2\varepsilon).$

- The reverse process is data generation.
  - The forward path serves as a guide for recovering data from noise.
  - Training enormous layers is possible.

- Learning the reverse process:
  - Treat $(\mathbf{x}_1, \dots, \mathbf{x}_T)$ as the latent variable $\mathbf{z}$.
  - The forward process defines $q(\mathbf{x}_1, \dots, \mathbf{x}_T|\mathbf{x}_0) \coloneqq q(\mathbf{x}_1|\mathbf{x}_0)\dots q(\mathbf{x}_T|\mathbf{x}_{T-1})$.
  - The reverse process defines $p_\theta(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_T) \coloneqq p(\mathbf{x}_T)p_\theta(\mathbf{x}_{T-1}|\mathbf{x}_T)\dots p_\theta(\mathbf{x}_0|\mathbf{x}_1)$.
  - Learn $\theta$ to make the posterior $p_\theta(\mathbf{x}_1, \dots, \mathbf{x}_T|\mathbf{x}_0)$ match $q(\mathbf{x}_1, \dots, \mathbf{x}_T|\mathbf{x}_0)$ by optimizing ELBO.

# Diffusion-Based Models

As a diffusion process (described by Stochastic Differential Equation) [SSK+21]:



Langevin dynamics targeting std. Gaussian (w/ time dilation $\beta(t)$).

- The forward process: discretizes Variance Preserving (VP) SDE: $d\mathbf{x} = -\frac{1}{2}\beta(t)\mathbf{x}\,dt + \sqrt{\beta(t)}\,d\mathbf{w}$.

- SDE theory gives the reverse process: $d\mathbf{x} = [\mathbf{f}(\mathbf{x},t) - g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x})]dt + g(t)d\bar{\mathbf{w}}$
  - Only the score function needs to be learned:
    $$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{\mathbf{x}(0)} \mathbb{E}_{\mathbf{x}(t)} \left[ \left\| \mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x}(t), t) - \nabla_{\mathbf{x}(t)} \log p_{0t}(\mathbf{x}(t) \mid \mathbf{x}(0)) \right\|_2^2 \right] \right\}.$$
    - When $\mathbf{f}(\mathbf{x}, t)$ is affine (e.g., VP SDE), $p_{0t}(\mathbf{x}_t|\mathbf{x}_0)$ is a Gaussian in closed-form.
    - Otherwise, $p_{0t}(\mathbf{x}_t|\mathbf{x}_0)$ is intractable. Use (sliced) score-matching objective.

# Diffusion-Based Models

As a diffusion process (described by Stochastic Differential Equation) [SSK+21]:

- Relation to noise annealing score-matching:

  Annealing corruption process: discretizes Variance Exploding (VE) SDE:

$$\mathbf{x}_i = \mathbf{x}_{i-1} + \sqrt{\sigma_i^2 - \sigma_{i-1}^2}\,\mathbf{z}_{i-1} \qquad \Longrightarrow \qquad \mathrm{d}\mathbf{x} = \sqrt{\frac{\mathrm{d}\left[\sigma^2(t)\right]}{\mathrm{d}t}}\,\mathrm{d}\mathbf{w}.$$

- Unified algorithm:

**Algorithm 1** PC sampling (VE SDE)

1: $\mathbf{x}_N \sim \mathcal{N}(\mathbf{0}, \sigma_{\max}^2 \mathbf{I})$
2: **for** $i = N - 1$ **to** $0$ **do**
3: $\quad \mathbf{x}_i' \leftarrow \mathbf{x}_{i+1} + (\sigma_{i+1}^2 - \sigma_i^2)\mathbf{s}_{\boldsymbol{\theta}*}(\mathbf{x}_{i+1}, \sigma_{i+1})$
4: $\quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
5: $\quad \mathbf{x}_i \leftarrow \mathbf{x}_i' + \sqrt{\sigma_{i+1}^2 - \sigma_i^2}\,\mathbf{z}$
6: $\quad$ **for** $j = 1$ **to** $M$ **do**
7: $\quad\quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
8: $\quad\quad \mathbf{x}_i \leftarrow \mathbf{x}_i + \epsilon_i \mathbf{s}_{\boldsymbol{\theta}*}(\mathbf{x}_i, \sigma_i) + \sqrt{2\epsilon_i}\,\mathbf{z}$
9: **return** $\mathbf{x}_0$

**Algorithm 2** PC sampling (VP SDE)

1: $\mathbf{x}_N \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
2: **for** $i = N - 1$ **to** $0$ **do**
3: $\quad \mathbf{x}_i' \leftarrow (2 - \sqrt{1 - \beta_{i+1}})\mathbf{x}_{i+1} + \beta_{i+1}\mathbf{s}_{\boldsymbol{\theta}*}(\mathbf{x}_{i+1}, i+1)$
4: $\quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
5: $\quad \mathbf{x}_i \leftarrow \mathbf{x}_i' + \sqrt{\beta_{i+1}}\,\mathbf{z}$      Predictor
6: $\quad$ **for** $j = 1$ **to** $M$ **do**      Corrector
7: $\quad\quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
8: $\quad\quad \mathbf{x}_i \leftarrow \mathbf{x}_i + \epsilon_i \mathbf{s}_{\boldsymbol{\theta}*}(\mathbf{x}_i, i) + \sqrt{2\epsilon_i}\,\mathbf{z}$
9: **return** $\mathbf{x}_0$

# Generative Model: Summary

| Plain Gen. | Latent Variable Models | | | | | |
|---|---|---|---|---|---|---|
| | Deterministic Generative | | Probabilistic Graphical Models | | | |
| Autoregres-sive Models | GANs | Flow-Based | Directed | Dir.: Diffusion | Undirected | Bidirectional |
| + Easy generation<br>+ Explicit llh ( easy learn-ing)<br>- No natural repr.<br>- Slow/seq. generation | + Easy generation | | | | - Hard generation (use MCMC) | |
| | - No llh (hard learning)<br>- Hard repr.<br>+ Flexible model | + Explicit llh ( easy learning)<br>+ Easy repr.<br>- High-dim. repr.<br>- Hard model design | Unnormalized llh: + stable learning,  - need expectation est. | | | |
| | | | + Moderate repr.<br>+ Prior knowledge<br>+ Small-data robust<br>+ Describe causality | + Easy repr.<br>+ Allow big model<br>- High-dim. repr. | - Hard repr.<br>- MCMC in learning<br>+ Simple depen-dency modeling | + Easy & flexible repr.<br>+ Flexible distribution |

**Colors represent:**
**Model component**
**Derived quantity**
**Auxiliary part**



$p(z)$ $z$
$x = f_\theta(z)$ (neural nets)
$p_\theta(x)$ $x$

$p(z)$ $z$
$x = f_\theta(z)$ (invertible)
$p_\theta(x)$ $x$

$p(z)$ $z$ $q_\phi(z|x)$
$p_\theta(x|z)$
$p_\theta(x)$ $x$

$p(z)$ $z$ $q(z|x)$ (fixed)
$p_\theta(x|z)$
$p_\theta(x)$ $x$

$z$
$p_\theta(x,z) \propto \exp(-E_\theta(x,z))$
$x$

$z$ $q_\phi(z|x)$
$p_\theta(x|z)$
$p_{\theta,\phi}(x)$ $x$

# Questions?

# References

# References

- Plain Generative Models
  - Autoregressive Models
    - [Fre98] Frey, Brendan J. (1998). *Graphical models for machine learning and digital communication*. MIT press.
    - [LM11] Larochelle, H. and Murray, I. The neural autoregressive distribution estimator. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2011.
    - [UML14] Uria, B., Murray, I., & Larochelle, H. (2014). A deep and tractable density estimator. In *International Conference on Machine Learning* (pp. 467-475).
    - [GGML15] Germain, M., Gregor, K., Murray, I., & Larochelle, H. (2015). MADE: Masked autoencoder for distribution estimation. In *International Conference on Machine Learning* (pp. 881-889).
    - [OKK16] Oord, A. V. D., Kalchbrenner, N., & Kavukcuoglu, K. (2016). Pixel recurrent neural networks. *arXiv preprint arXiv:1601.06759*.
    - [ODZ+16] Oord, A. V. D., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A., & Kavukcuoglu, K. (2016). WaveNet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*.

# References

- Deterministic Generative Models
  - Generative Adversarial Networks
    - [GPM+14] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative adversarial nets. In *Advances in neural information processing systems* (pp. 2672-2680).
    - [ACB17] Arjovsky, M., Chintala, S., & Bottou, L. (2017). Wasserstein generative adversarial networks. In *International Conference on Machine Learning* (pp. 214-223).
  - Flow-Based Models
    - [DKB15] Dinh, L., Krueger, D., & Bengio, Y. (2015). NICE: Non-linear independent components estimation. *ICLR workshop*.
    - [DSB17] Dinh, L., Sohl-Dickstein, J., & Bengio, S. (2017). Density estimation using real NVP. In *Proceedings of the International Conference on Learning Representations*.
    - [PPM17] Papamakarios, G., Pavlakou, T., & Murray, I. (2017). Masked autoregressive flow for density estimation. In *Advances in Neural Information Processing Systems*.
    - [KD18] Kingma, D. P., & Dhariwal, P. (2018). Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems*.

# References

- Deterministic Generative Models
  - Flow-Based Models
    - [KSJ+16] Kingma, D. P., Salimans, T., Jozefowicz, R., Chen, X., Sutskever, I., & Welling, M. (2016). Improved Variational Inference with Inverse Autoregressive Flow. In *Advances in Neural Information Processing Systems*.
    - [GCB+18] Grathwohl, W., Chen, R. T., Bettencourt, J., Sutskever, I., & Duvenaud, D. (2018). FFJORD: Free-form continuous dynamics for scalable reversible generative models. In *Proceedings of International Conference on Learning Representations*.
    - [BGC19] Behrmann, J., Grathwohl, W., Chen, R. T., Duvenaud, D., & Jacobsen, J. H. (2019). Invertible residual networks. In *International Conference on Machine Learning*.
    - [CBD19] Chen, R. T., Behrmann, J., Duvenaud, D., & Jacobsen, J. H. (2019). Residual Flows for Invertible Generative Modeling. In *Advances in Neural Information Processing Systems*.
    - [KC20] Kong, Z., & Chaudhuri, K. (2020). The expressive power of a class of normalizing flow models. In *International Conference on Artificial Intelligence and Statistics*.
    - [TIT+20] Teshima, T., Ishikawa, I., Tojo, K., Oono, K., Ikeda, M., & Sugiyama, M. (2020). Coupling-based invertible neural networks are universal diffeomorphism approximators. *arXiv preprint arXiv:2006.11469*.

# References

- Bayesian Inference: Variational Inference
  - Explicit Parametric VI:

- [SJJ96] Saul, L. K., Jaakkola, T., & Jordan, M. I. (1996). Mean field theory for sigmoid belief networks. *Journal of artificial intelligence research*, 4, 61-76.
- [BNJ03] Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent Dirichlet Allocation. *Journal of Machine Learning Research*, 3(Jan), pp.993-1022.
- [GHB12] Gershman, S., Hoffman, M., & Blei, D. (2012). Nonparametric variational inference. arXiv preprint arXiv:1206.4665.
- [HBWP13] Hoffman, M. D., Blei, D. M., Wang, C., & Paisley, J. (2013). Stochastic variational inference. *The Journal of Machine Learning Research*, *14*(1), 1303-1347.
- [RGB14] Ranganath, R., Gerrish, S., & Blei, D. (2014). Black box variational inference. In *Artificial Intelligence and Statistics* (pp. 814-822).
- [RM15] Rezende, D.J., & Mohamed, S. (2015). Variational inference with normalizing flows. In *Proceedings of the International Conference on Machine Learning* (pp. 1530-1538).
- [KSJ+16] Kingma, D.P., Salimans, T., Jozefowicz, R., Chen, X., Sutskever, I., & Welling, M. (2016). Improved variational inference with inverse autoregressive flow. In *Advances in neural information processing systems* (pp. 4743-4751).

# References

- Bayesian Inference: Variational Inference
  - Implicit Parametric VI: density ratio estimation
    - [MSJ+15] Makhzani, A., Shlens, J., Jaitly, N., Goodfellow, I., & Frey, B. (2016). Adversarial Autoencoders. In *Proceedings of the International Conference on Learning Representations*.
    - [MNG17] Mescheder, L., Nowozin, S., & Geiger, A. (2017). Adversarial variational Bayes: Unifying variational autoencoders and generative adversarial networks. In *Proceedings of the International Conference on Machine Learning* (pp. 2391-2400).
    - [Hus17] Huszár, F. (2017). Variational inference using implicit distributions. *arXiv preprint* arXiv:1702.08235.
    - [TRB17] Tran, D., Ranganath, R., & Blei, D. (2017). Hierarchical implicit models and likelihood-free variational inference. In *Advances in Neural Information Processing Systems* (pp. 5523-5533).
    - [SSZ18a] Shi, J., Sun, S., & Zhu, J. (2018). Kernel Implicit Variational Inference. In *Proceedings of the International Conference on Learning Representations*.
  - Implicit Parametric VI: gradient estimation
    - [VLBM08] Vincent, P., Larochelle, H., Bengio, Y., & Manzagol, P. A. (2008). Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning* (pp. 1096-1103). ACM.
    - [LT18] Li, Y., & Turner, R. E. (2018). Gradient estimators for implicit models. In *Proceedings of the International Conference on Learning Representations*.
    - [SSZ18b] Shi, J., Sun, S., & Zhu, J. (2018). A spectral approach to gradient estimation for implicit distributions. In *Proceedings of the 35th International Conference on Machine Learning* (pp. 4651-4660).

# References

- Bayesian Inference: Variational Inference
  - Particle-Based VI
    - [LW16] Liu, Q., & Wang, D. (2016). Stein variational gradient descent: A general purpose Bayesian inference algorithm. In *Advances In Neural Information Processing Systems* (pp. 2378-2386).
    - [Liu17] Liu, Q. (2017). Stein variational gradient descent as gradient flow. In *Advances in neural information processing systems* (pp. 3115-3123).
    - [CZ17] Chen, C., & Zhang, R. (2017). Particle optimization in stochastic gradient MCMC. *arXiv preprint arXiv:1711.10927*.
    - [FWL17] Feng, Y., Wang, D., & Liu, Q. (2017). Learning to Draw Samples with Amortized Stein Variational Gradient Descent. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*.
    - [PGH+17] Pu, Y., Gan, Z., Henao, R., Li, C., Han, S., & Carin, L. (2017). VAE learning via Stein variational gradient descent. In *Advances in Neural Information Processing Systems* (pp. 4236-4245).
    - [LZ18] Liu, C., & Zhu, J. (2018). Riemannian Stein Variational Gradient Descent for Bayesian Inference. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence* (pp. 3627-3634).

# References

- Bayesian Inference: Variational Inference
  - Particle-Based VI
    - [CMG+18] Chen, W. Y., Mackey, L., Gorham, J., Briol, F. X., & Oates, C. J. (2018). Stein points. *arXiv preprint arXiv:1803.10161*.
    - [FCSS18] Futami, F., Cui, Z., Sato, I., & Sugiyama, M. (2018). Frank-Wolfe Stein sampling. *arXiv preprint arXiv:1805.07912*.
    - [CZW+18] Chen, C., Zhang, R., Wang, W., Li, B., & Chen, L. (2018). A unified particle-optimization framework for scalable Bayesian sampling. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*.
    - [ZZC18] Zhang, J., Zhang, R., & Chen, C. (2018). Stochastic particle-optimization sampling and the non-asymptotic convergence theory. *arXiv preprint arXiv:1809.01293*.
    - [LZC+19] Liu, C., Zhuo, J., Cheng, P., Zhang, R., Zhu, J., & Carin, L. (2019). Understanding and Accelerating Particle-Based Variational Inference. In *Proceedings of the 36th International Conference on Machine Learning* (pp. 4082-4092).

# References

- Bayesian Inference: MCMC
  - Classical MCMC
    - [MRR+53] Metropolis, N., Rosenbluth, A. W., Rosenbluth, M.N., Teller, A.H., & Teller, E. (1953). Equation of state calculations by fast computing machines. *The Journal of Chemical Physics*, 21(6), pp.1087-1092.
    - [Has70] Hastings, W. (1970). Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, 57(1), pp.97-109.
    - [GG87] Geman, S., & Geman, D. (1987). Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. In *Readings in computer vision* (pp. 564-584).
    - [ADDJ03] Andrieu, C., De Freitas, N., Doucet, A., & Jordan, M. I. (2003). An introduction to MCMC for machine learning. *Machine learning*, *50*(1-2), 5-43.

# References

- Bayesian Inference: MCMC
  - Dynamics-Based MCMC: full-batch
    - [Lan08] Langevin, P. (1908). Sur la théorie du mouvement Brownien. *Compt. Rendus*, *146*, 530-533.
    - [DKPR87] Duane, S., Kennedy, A.D., Pendleton, B.J., Roweth, D. (1987). Hybrid Monte Carlo. *Physics letters B*, 195(2), pp.216-222.
    - [RT96] Roberts, G. O., & Tweedie, R. L. (1996). Exponential convergence of Langevin distributions and their discrete approximations. *Bernoulli*, *2*(4), 341-363.
    - [RS02] Roberts, G.O., & Stramer, O. (2002). Langevin diffusions and Metropolis-Hastings algorithms. *Methodology and computing in applied probability*, 4(4), pp.337-357.
    - [Nea11] Neal, R.M. (2011). MCMC using Hamiltonian dynamics. Handbook of Markov chain Monte Carlo, 2(11), p.2.
    - [ZWC+16] Zhang, Y., Wang, X., Chen, C., Henao, R., Fan, K., & Carin, L. (2016). Towards unifying Hamiltonian Monte Carlo and slice sampling. In *Advances in Neural Information Processing Systems* (pp. 1741-1749).
    - [TRGT17] Tripuraneni, N., Rowland, M., Ghahramani, Z., & Turner, R. (2017, August). Magnetic Hamiltonian Monte Carlo. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70* (pp. 3453-3461).
    - [Bet17] Betancourt, M. (2017). A conceptual introduction to Hamiltonian Monte Carlo. *arXiv preprint arXiv:1701.02434*.

# References

- Bayesian Inference: MCMC
  - Dynamics-Based MCMC: full-batch (manifold support)
    - [GC11] Girolami, M., & Calderhead, B. (2011). Riemann manifold Langevin and Hamiltonian Monte Carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, *73*(2), 123-214.
    - [BSU12] Brubaker, M., Salzmann, M., & Urtasun, R. (2012, March). A family of MCMC methods on implicitly defined manifolds. In *Artificial intelligence and statistics* (pp. 161-172).
    - [BG13] Byrne, S., & Girolami, M. (2013). Geodesic Monte Carlo on embedded manifolds. *Scandinavian Journal of Statistics*, *40*(4), 825-845.
    - [LSSG15] Lan, S., Stathopoulos, V., Shahbaba, B., & Girolami, M. (2015). Markov chain Monte Carlo from Lagrangian dynamics. *Journal of Computational and Graphical Statistics*, *24*(2), 357-378.

# References

- Bayesian Inference: MCMC
  - Dynamics-Based MCMC: stochastic gradient
    - [WT11] Welling, M., & Teh, Y. W. (2011). Bayesian learning via stochastic gradient Langevin dynamics. In *Proceedings of the International Conference on Machine Learning* (pp. 681-688).
    - [CFG14] Chen, T., Fox, E., & Guestrin, C. (2014). Stochastic gradient Hamiltonian Monte Carlo. In *Proceedings of the International conference on machine learning* (pp. 1683-1691).
    - [DFB+14] Ding, N., Fang, Y., Babbush, R., Chen, C., Skeel, R. D., & Neven, H. (2014). Bayesian sampling using stochastic gradient thermostats. In *Advances in neural information processing systems* (pp. 3203-3211).
    - [Bet15] Betancourt, M. (2015). The fundamental incompatibility of scalable Hamiltonian Monte Carlo and naive data subsampling. In *International Conference on Machine Learning* (pp. 533-540).
    - [TTV16] Teh, Y. W., Thiery, A. H., & Vollmer, S. J. (2016). Consistency and fluctuations for stochastic gradient Langevin dynamics. *The Journal of Machine Learning Research*, *17*(1), 193-225.
    - [LPH+16] Lu, X., Perrone, V., Hasenclever, L., Teh, Y. W., & Vollmer, S. J. (2016). Relativistic Monte Carlo. *arXiv preprint arXiv:1609.04388*.
    - [ZCG+17] Zhang, Y., Chen, C., Gan, Z., Henao, R., & Carin, L. (2017, August). Stochastic gradient monomial Gamma sampler. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70* (pp. 3996-4005).
    - [LTL17] Li, Y., Turner, R.E., & Liu, Q. (2017). Approximate inference with amortised MCMC. *arXiv preprint arXiv:1702.08343*.

# References

- Bayesian Inference: MCMC
  - Dynamics-Based MCMC: stochastic gradient (manifold support)
    - [PT13] Patterson, S., & Teh, Y.W. (2013). Stochastic gradient Riemannian Langevin dynamics on the probability simplex. In *Advances in neural information processing systems* (pp. 3102-3110).
    - [MCF15] Ma, Y. A., Chen, T., & Fox, E. (2015). A complete recipe for stochastic gradient MCMC. In *Advances in Neural Information Processing Systems* (pp. 2917-2925).
    - [LZS16] Liu, C., Zhu, J., & Song, Y. (2016). Stochastic Gradient Geodesic MCMC Methods. In *Advances in Neural Information Processing Systems* (pp. 3009-3017).
  - Dynamics-Based MCMC: general theory
    - [JKO98] Jordan, R., Kinderlehrer, D., & Otto, F. (1998). The variational formulation of the Fokker-Planck equation. *SIAM journal on mathematical analysis*, *29*(1), 1-17.
    - [MCF15] Ma, Y. A., Chen, T., & Fox, E. (2015). A complete recipe for stochastic gradient MCMC. In *Advances in Neural Information Processing Systems* (pp. 2917-2925).
    - [CDC15] Chen, C., Ding, N., & Carin, L. (2015). On the convergence of stochastic gradient MCMC algorithms with high-order integrators. In *Advances in Neural Information Processing Systems* (pp. 2278-2286).
    - [LZZ19] Liu, C., Zhuo, J., & Zhu, J. (2019). Understanding MCMC Dynamics as Flows on the Wasserstein Space. In *Proceedings of the 36th International Conference on Machine Learning* (pp. 4093-4103).

# References

- Probabilistic Graphical Models
  - Directed PGM: Causality
    - [SG91] Spirtes, P., Glymour, C. (1991). An algorithm for fast recovery of sparse causal graphs. *Social Science Computer Review*. 9 (1): 62–72.
    - [Pearl09] Pearl, J. (2009). *Causality*. Cambridge university press.
    - [PJS17] Peters, J., Janzing, D., & Schölkopf, B. (2017). *Elements of causal inference: foundations and learning algorithms*. The MIT Press.
    - [IW09] Imbens, G. W. & Wooldridge, J. M. Recent developments in the econometrics of program evaluation. *Journal of economic literature*, 47(1):5–86, 2009.
    - [Pearl15] Pearl, J. Detecting latent heterogeneity. *Sociological Methods & Research*, pp. 0049124115600597, 2015.
    - [SJP+12] Schölkopf, B., Janzing, D., Peters, J., Sgouritsa, E., Zhang, K., & Mooij, J. (2012). On causal and anticausal learning. In *International Conference on Machine Learning*.

# References

- Probabilistic Graphical Models
  - Directed PGM: Causality
    - [Sch19] Schölkopf, B. (2019). Causality for machine learning. *arXiv preprint arXiv:1911.10500*.
    - [YGL+20] Yu, K., Guo, X., Liu, L., Li, J., Wang, H., Ling, Z., & Wu, X. (2020). Causality-based feature selection: Methods and evaluations. *ACM Computing Surveys (CSUR)*, 53(5), 1-36.
    - [LSW+21] Liu, C., Sun, X., Wang, J., Tang, H., Li, T., Qin, T., Chen, W., & Liu, T. Y. (2021). Learning Causal Semantic Representation for Out-of-Distribution Prediction. In *Advances in Neural Information Processing Systems*.
    - [SWZ+21] Sun, X., Wu, B., Zheng, X., Liu, C., Chen, W., Qin, T., & Liu, T. Y. (2021). Recovering Latent Causal Factor for Generalization to Distributional Shifts. In *Advances in Neural Information Processing Systems*.
  - Related:
    - [ABGL19] Arjovsky, M., Bottou, L., Gulrajani, I., & Lopez-Paz, D. (2019). Invariant risk minimization. *arXiv preprint arXiv:1907.02893*.

# References

- Probabilistic Graphical Models
  - Directed PGM: Topic Models
    - [BNJ03] Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent Dirichlet Allocation. *Journal of Machine Learning Research*, 3(Jan), pp.993-1022.
    - [GS04] Griffiths, T.L., & Steyvers, M. (2004). Finding scientific topics. *Proceedings of the National academy of Sciences*, 101 (suppl 1), pp.5228-5235.
    - [SG07] Steyvers, M., & Griffiths, T. (2007). Probabilistic topic models. *Handbook of latent semantic analysis*, *427*(7), 424-440.
    - [MB08] Mcauliffe, J.D., & Blei, D.M. (2008). Supervised topic models. In *Advances in neural information processing systems* (pp. 121-128).
    - [ZAX12] Zhu, J., Ahmed, A., & Xing, E. P. (2012). MedLDA: maximum margin supervised topic models. *Journal of Machine Learning Research*, *13*(Aug), 2237-2278.
    - [PT13] Patterson, S., & Teh, Y.W. (2013). Stochastic gradient Riemannian Langevin dynamics on the probability simplex. In *Advances in neural information processing system* (pp. 3102-3110).
    - [ZCX14] Zhu, J., Chen, N., & Xing, E. P. (2014). Bayesian inference with posterior regularization and applications to infinite latent SVMs. *The Journal of Machine Learning Research*, *15*(1), 1799-1847.

# References

- Probabilistic Graphical Models
  - Directed PGM: Topic Models
    - [LARS14] Li, A.Q., Ahmed, A., Ravi, S., & Smola, A.J. (2014). Reducing the sampling complexity of topic models. In *Proceedings of the ACM SIGKDD international conference on Knowledge discovery and data mining* (pp. 891-900).
    - [YGH+15] Yuan, J., Gao, F., Ho, Q., Dai, W., Wei, J., Zheng, X., Xing, E.P., Liu, T.Y., & Ma, W.Y. (2015). LightLDA: Big topic models on modest computer clusters. In *Proceedings of the 24th International Conference on World Wide Web* (pp. 1351-1361).
    - [CLZC16] Chen, J., Li, K., Zhu, J., & Chen, W. (2016). WarpLDA: a cache efficient o(1) algorithm for latent Dirichlet allocation. *Proceedings of the VLDB Endowment*, 9(10), pp.744-755.

# References

- Probabilistic Graphical Models
  - Directed PGM: Variational Auto-Encoders
    - [KW14] Kingma, D.P., & Welling, M. (2014). Auto-encoding variational Bayes. In *Proceedings of the International Conference on Learning Representations*.
    - [KMRW14] Kingma, D. P., Mohamed, S., Rezende, D. J., & Welling, M. (2014). Semi-supervised learning with deep generative models. In *Advances in neural information processing systems* (pp. 3581-3589).
    - [SLY15] Sohn, K., Lee, H., & Yan, X. (2015). Learning structured output representation using deep conditional generative models. In *Advances in Neural Information Processing Systems*, 28, 3483-3491.
    - [GDG+15] Gregor, K., Danihelka, I., Graves, A., Rezende, D. J., & Wierstra, D. (2015). DRAW: A recurrent neural network for image generation. In *Proceedings of the 32nd International Conference on Machine Learning*.
    - [BGS15] Burda, Y., Grosse, R., & Salakhutdinov, R. (2015). Importance weighted autoencoders. *arXiv preprint arXiv:1509.00519*.

# References

- Probabilistic Graphical Models
  - Directed PGM: Variational Auto-Encoders

- [DFD+18] Davidson, T.R., Falorsi, L., De Cao, N., Kipf, T., & Tomczak, J.M. (2018). Hyperspherical variational auto-encoders. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*.
- [MSJ+15] Makhzani, A., Shlens, J., Jaitly, N., Goodfellow, I., & Frey, B. (2016). Adversarial Autoencoders. In *Proceedings of the International Conference on Learning Representations*.
- [CDH+16] Chen, X., Duan, Y., Houthooft, R., Schulman, J., Sutskever, I., & Abbeel, P. (2016). InfoGAN: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in neural information processing systems* (pp. 2172-2180).
- [MNG17] Mescheder, L., Nowozin, S., & Geiger, A. (2017). Adversarial variational Bayes: Unifying variational autoencoders and generative adversarial networks. In *Proceedings of the International Conference on Machine Learning* (pp. 2391-2400).
- [TBGS17] Tolstikhin, I., Bousquet, O., Gelly, S., & Schoelkopf, B. (2017). Wasserstein Auto-Encoders. *arXiv preprint arXiv:1711.01558*.

# References

- Probabilistic Graphical Models
  - Directed PGM: Variational Auto-Encoders

  - [FWL17] Feng, Y., Wang, D., & Liu, Q. (2017). Learning to Draw Samples with Amortized Stein Variational Gradient Descent. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*.
  - [PGH+17] Pu, Y., Gan, Z., Henao, R., Li, C., Han, S., & Carin, L. (2017). VAE learning via Stein variational gradient descent. In *Advances in Neural Information Processing Systems* (pp. 4236-4245).
  - [KSDV18] Kocaoglu, M., Snyder, C., Dimakis, A. G., & Vishwanath, S. (2018). CausalGAN: Learning causal implicit generative models with adversarial training. In *Proceedings of the International Conference on Learning Representations*.
  - [LWZZ18] Li, C., Welling, M., Zhu, J., & Zhang, B. (2018). Graphical generative adversarial networks. In *Advances in Neural Information Processing Systems* (pp. 6069-6080).
  - [DW19] Dai, B., & Wipf, D. (2019). Diagnosing and Enhancing Gaussian VAE Models. In *International Conference on Learning Representations*.

# References

- Probabilistic Graphical Models
  - Directed PGM: Variational Auto-Encoders
    - [LBN+19] Luo, Y., Beatson, A., Norouzi, M., Zhu, J., Duvenaud, D., Adams, R. P., & Chen, R. T. (2019). SUMO: Unbiased Estimation of Log Marginal Probability for Latent Variable Models. In *International Conference on Learning Representations*.
    - [PHN+20] Pang, B., Han, T., Nijkamp, E., Zhu, S. C., & Wu, Y. N. (2020). Learning Latent Space Energy-Based Prior Model. In *Advances in Neural Information Processing Systems*.
    - [LTQ+21] Liu, C., Tang, H., Qin, T., Wang, J., & Liu, T. Y. (2021). On the Generative Utility of Cyclic Conditionals. In *Advances in Neural Information Processing Systems*.

# References

- Probabilistic Graphical Models
  - Directed PGM: Disentanglement
    - [CDH+16] Chen, X., Duan, Y., Houthooft, R., Schulman, J., Sutskever, I., & Abbeel, P. (2016). InfoGAN: interpretable representation learning by information maximizing Generative Adversarial Nets. In *Neural Information Processing Systems*.
    - [HMP+17] Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., & Lerchner, A. (2017). Beta-VAE: Learning basic visual concepts with a constrained variational framework. In *Proceedings of the International Conference on Learning Representations*.
    - [CLG+18] Chen, R. T., Li, X., Grosse, R., & Duvenaud, D. (2018). Isolating sources of disentanglement in VAEs. In *Neural Information Processing Systems*.
    - [HAP+18] Higgins, I., Amos, D., Pfau, D., Racaniere, S., Matthey, L., Rezende, D., & Lerchner, A. (2018). Towards a definition of disentangled representations. *arXiv preprint arXiv:1812.02230*.
    - [LBL+19] Locatello, F., Bauer, S., Lucic, M., Raetsch, G., Gelly, S., Schölkopf, B., & Bachem, O. (2019). Challenging common assumptions in the unsupervised learning of disentangled representations. In *International Conference on Machine Learning*.

# References

- Probabilistic Graphical Models
  - Directed PGM: Disentanglement
    - [LTB+19] Locatello, F., Tschannen, M., Bauer, S., Rätsch, G., Schölkopf, B., & Bachem, O. (2019). Disentangling Factors of Variations Using Few Labels. In *International Conference on Learning Representations*.
    - [CB20] Chen, J., & Batmanghelich, K. (2020). Weakly supervised disentanglement by pairwise similarities. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
    - [LPR+20] Locatello, F., Poole, B., Rätsch, G., Schölkopf, B., Bachem, O., & Tschannen, M. (2020). Weakly-supervised disentanglement without compromises. In *International Conference on Machine Learning*.
    - [SCK+20] Shu, R., Chen, Y., Kumar, A., Ermon, S., & Poole, B. (2020). Weakly Supervised Disentanglement with Guarantees. In *International Conference on Learning Representations*.
    - [KKM+20] Khemakhem, I., Kingma, D., Monti, R., & Hyvarinen, A. (2020). Variational autoencoders and nonlinear ICA: A unifying framework. In *International Conference on Artificial Intelligence and Statistics*.

# References

- Probabilistic Graphical Models
  - Undirected PGM
    - [HS83] Hinton, G., & Sejnowski, T. (1983). Optimal perceptual inference. In *IEEE Conference on Computer Vision and Pattern Recognition*.
    - [Smo86] Smolensky, P. (1986). Information processing in dynamical systems: Foundations of harmony theory. In *Parallel Distributed Processing*, volume 1, chapter 6, pages 194-281. MIT Press.
    - [Hin02] Hinton, G. E. (2002). Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8), 1771-1800.
    - [LCH+06] LeCun, Y., Chopra, S., Hadsell, R., Ranzato, M., & Huang, F. (2006). A tutorial on energy-based learning. *Predicting structured data*, *1*(0).
    - [HOT06] Hinton, G. E., Osindero, S., & Teh, Y. W. (2006). A fast learning algorithm for deep belief nets. *Neural computation*, 18(7), 1527-1554.
    - [SH09] Salakhutdinov, R., & Hinton, G. (2009, April). Deep Boltzmann machines. *In AISTATS* (pp. 448-455).
    - [Sal15] Salakhutdinov, R. (2015). Learning deep generative models. *Annual Review of Statistics and Its Application*, 2, 361-385.
    - [KB16] Kim, T., & Bengio, Y. (2016). Deep directed generative models with energy-based probability estimation. *arXiv preprint arXiv:1606.03439*.
    - [DM19] Du, Y., & Mordatch, I. (2019). Implicit generation and generalization in energy-based models. In *Advances in Neural Information Processing Systems*.

# References

- Probabilistic Graphical Models
  - Score-based methods
    - [Hyv05] Hyvärinen, A., & Dayan, P. (2005). Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4).
    - [SE19] Song, Y., & Ermon, S. (2019). Generative Modeling by Estimating Gradients of the Data Distribution. In *Advances in Neural Information Processing Systems*.

# References

- Probabilistic Graphical Models
  - Diffusion-based models
    - [SWMG15] Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., & Ganguli, S. (2015). Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning* (pp. 2256-2265).
    - [HJA20] Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*.
    - [SSK+21] Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. (2021). Score-Based Generative Modeling through Stochastic Differential Equations. In *International Conference on Learning Representations*.
    - [BTHD21] De Bortoli, V., Thornton, J., Heng, J., & Doucet, A. (2021). Diffusion Schrödinger bridge with applications to score-based generative modeling. In *Advances in Neural Information Processing Systems*.
    - [DVK22] Dockhorn, T., Vahdat, A., & Kreis, K. (2022). Score-Based Generative Modeling with Critically-Damped Langevin Diffusion. In *International Conference on Learning Representations*.

# References

- Others
  - Probabilistic Graphical Models
    - [KYD+18] Kim, T., Yoon, J., Dia, O., Kim, S., Bengio, Y., & Ahn, S. (2018). Bayesian model-agnostic meta-learning. In *Advances in Neural Information Processing Systems* (pp. 7332-7342).
    - [LST15] Lake, B.M., Salakhutdinov, R., & Tenenbaum, J.B. (2015). Human-level concept learning through probabilistic program induction. *Science*, 350(6266), pp. 1332-1338.
  - Bayesian Neural Network
    - [LG17] Li, Y., & Gal, Y. (2017). Dropout inference in Bayesian neural networks with alpha-divergences. In *Proceedings of the International Conference on Machine Learning* (pp. 2052-2061).
  - Related References
    - [Wil92] Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4), 229-256.
    - [HV93] Hinton, G., & Van Camp, D. (1993). Keeping neural networks simple by minimizing the description length of the weights. In *in Proc. of the 6th Ann. ACM Conf. on Computational Learning Theory*.
    - [NJ01] Ng, A. Y., & Jordan, M. I. (2002). On discriminative vs. generative classifiers: A comparison of logistic regression and naive Bayes. In *Advances in neural information processing systems* (pp. 841-848).

# The End