Generative Models



Generative Adversarial Nets



Implicit Generative Model

- Goal: a sampler $g(\cdot)$ to generate images
- A simple generator $g(z; \theta)$:
 - $z \sim N(0,I)$
 - $x = g(z; \theta)$ deterministic transformation
- Likelihood-free training:
 - Given a dataset from some distribution p_{data}
 - Goal: $g(z;\theta)$ defines a distribution, we want this distribution $pprox p_{data}$
 - Training: minimize $D(g(z;\theta),p_{data})$
 - *D* is some distance metric (not likelihood)
 - ullet Key idea: **Learn a differentiable** D

GAN (Goodfellow et al., '14)

- ullet Parameterize the discriminator $D(\ \cdot\ ; \phi)$ with parameter ϕ
- Goal: learn ϕ such that $D(x;\phi)$ measures how likely x is from p_{data}
 - $D(x, \phi) = 1$ if $x \sim p_{data}$
 - $D(x, \phi) = 0$ if $x! \sim p_{data}$
 - a.k.a., a binary classifier
- GAN: use a neural network for $D(\cdot;\phi)$
- Training: need both negative and positive samples
 - Positive samples: just the training data
 - Negative samples: use our sampler $g(\cdot;z)$ (can provide infinite samples).
- Overall objectives:
 - Generator: $\theta^* = \max_{\theta} D(g(z; \theta); \phi)$
 - Discriminator uses MLE Training:

$$\phi^* = \max_{\phi} \mathbb{E}_{x \sim p_{data}} [\log D(x; \phi)] + \mathbb{E}_{\hat{x} \sim g(\cdot)} [\log(1 - D(\hat{x}; \phi))]$$

GAN (Goodfellow et al., '14)

- Generator $G(z; \theta)$ where $z \sim N(0,I)$
 - Generate realistic data
- Discriminator $D(x; \phi)$
 - Classify whether the data is real (from p_{data}) or fake (from G)
- Objective function:

$$L(\theta, \phi) = \min_{\theta} \max_{\phi} \mathbb{E}_{x \sim p_{data}} \left[\log D(x; \phi) \right] + \mathbb{E}_{\hat{x} \sim G} \left[\log(1 - D(\hat{x}; \phi)) \right]$$

- Training procedure:
 - Collect dataset $\{(x,1) | x \sim p_{data}\} \cup \{(\hat{x},0) \sim g(z;\theta)\}$
 - Train discriminator

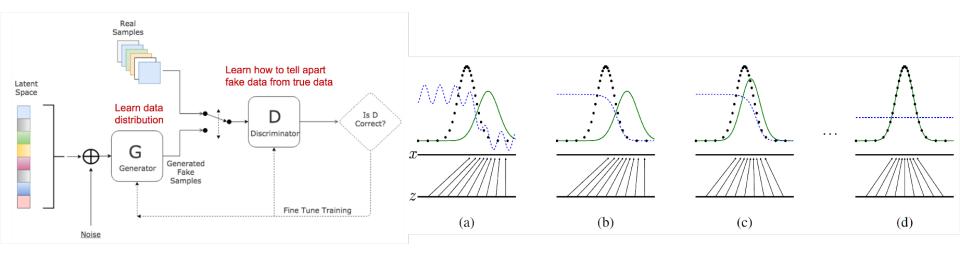
$$D: L(\phi) = \mathbb{E}_{x \sim p_{data}} \left[\log D(x; \phi) \right] + \mathbb{E}_{\hat{x} \sim G} \left[\log(1 - D(\hat{x}; \phi)) \right]$$

- Train generator $G: L(\theta) = \mathbb{E}_{z \sim N(0,I)} \left[\log D(G(z;\theta),\phi) \right]$
- Repeat

GAN (Goodfellow et al., '14)

Objective function:

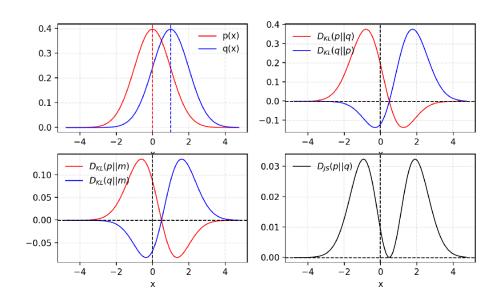
$$L(\theta, \phi) = \min_{\theta} \max_{\phi} \mathbb{E}_{x \sim p_{data}} \left[\log D(x; \phi) \right] + \mathbb{E}_{\hat{x} \sim G} \left[\log(1 - D(\hat{x}; \phi)) \right]$$



Math Behind GAN

Math Behind GAN

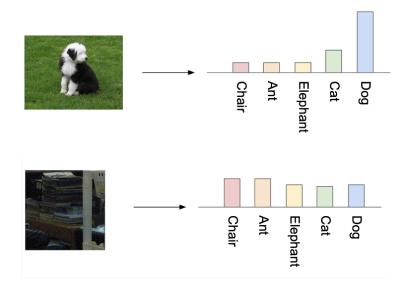
KL-Divergence and JS-Divergence



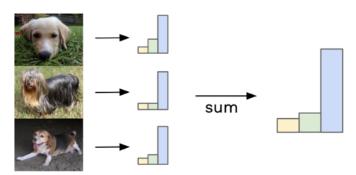
Math Behind GAN

Evaluation of GAN

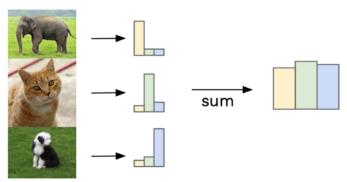
- No p(x) in GAN.
- Idea: use a trained classifier $f(y \mid x)$:
- If $x \sim p_{data}$, f(y | x) should have low entropy
 - Otherwise, $f(y \mid x)$ close to uniform.
- Samples from *G* should be diverse:
 - $p_f(y) = \mathbb{E}_{x \sim G}[f(y \mid x)]$ close to uniform.



Similar labels sum to give focussed distribution



Different labels sum to give uniform distribution



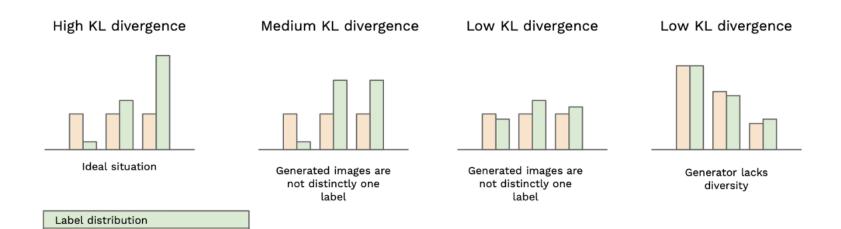
Evaluation of GAN

- Inception Score (IS, Salimans et al. '16)
 - Use Inception V3 trained on ImageNet as f(y | x)

•
$$IS = \exp\left(\mathbb{E}_{x \sim G}\left[KL(f(y|x)||p_f(y)))\right]\right)$$

Higher the better

Marginal distribution



Comments on GAN

- Other evaluation metrics:
 - Fréchet Inception Distance (FID): Wasserstein distance between Gaussians
- Mode collapse:
 - The generator only generate a few type of samples.
 - Or keep oscillating over a few modes.
- Training instability:
 - Discriminator and generator may keep oscillating
 - Example: -xy, generator x, discriminatory. NE: x = y = 0 but GD oscillates.
 - No stopping criteria.
 - Use Wsserstein GAN (Arjovsky et al. '17):

$$\min_{G} \max_{f: \mathsf{Lip}(f) \le 1} \mathbb{E}_{x \sim p_{data}} \left[f(x) \right] - \mathbb{E}_{\hat{x} \sim p_{G}} [f(\hat{x})]$$

And need many other tricks...

Energy-Based Models



Energy-based Models

- Goal of generative models:
 - a probability distribution of data: P(x)
- Requirements
 - $P(x) \ge 0$ (non-negative)

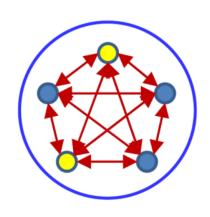
$$\int_{x} P(x)dx = 1$$

- Energy-based model:
 - Energy function: $E(x; \theta)$, parameterized by θ
 - $P(x) = \frac{1}{z} \exp(-E(x;\theta))$ (why exp?) $z = \int_{z} \exp(-E(x;\theta))dx$

Boltzmann Machine

- Generative model

 - $\bullet \ E(y) = -\frac{1}{2} y^{\top} W y$ $\bullet \ P(y) = \frac{1}{z} \exp(-\frac{E(y)}{T}) \text{, T: temperature hyper-parameter}$
 - W: parameter to learn
- ullet When y_i is binary, patterns are affecting each other through W



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

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

Boltzmann Machine: Training

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

$$P(y) = \frac{\exp(\frac{1}{2}y^{\mathsf{T}}y)}{\sum_{y'} \exp(y'^{\mathsf{T}}Wy')}$$

• Maximum log-likelihood:

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

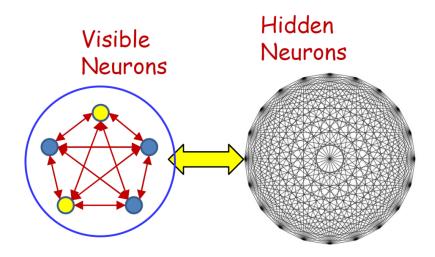
Boltzmann Machine: Training

Boltzmann Machine: Training

Boltzmann Machine with Hidden Neurons

- Visible and hidden neurons:
 - y: visible, h: hidden

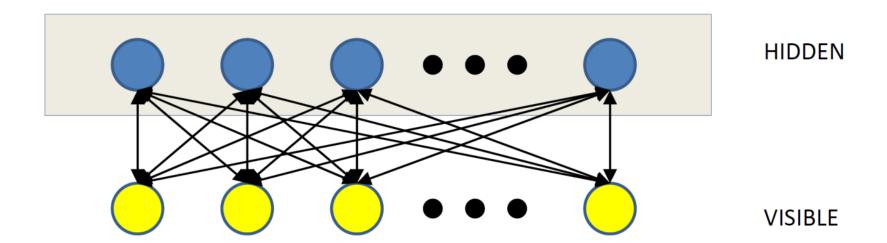
$$P(y) = \sum_{h} P(y, v)$$



Boltzmann Machine with Hidden Neurons: Training

Boltzmann Machine with Hidden Neurons: Training

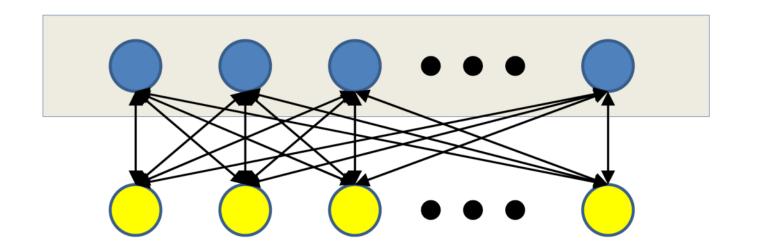
- A structured Boltzmann Machine
 - Hidden neurons are only connected to visible neurons
 - No intra-layer connections
 - Invented by Paul Smolensky in '89
 - Became more practical after Hinton invested fast learning algorithms in mid
 2000



- Computation Rules
 - Iterative sampling

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

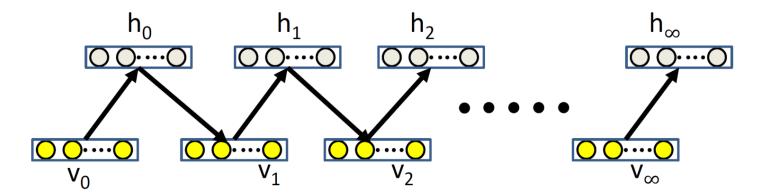
• Visible neurons
$$v_j$$
: $z_j = \sum_i w_{ij} h_i$, $P(v_j | h) = \frac{1}{1 + \exp(-z_j)}$



HIDDEN

VISIBLE

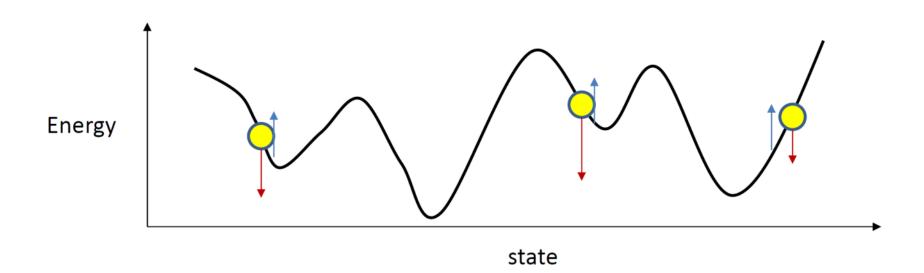
- Sampling:
 - Randomly initialize visible neurons v_0
 - Iterative sampling between hidden neurons and visible neurons
 - Get final sample (v_{∞}, h_{∞})



Maximum likelihood estimated:

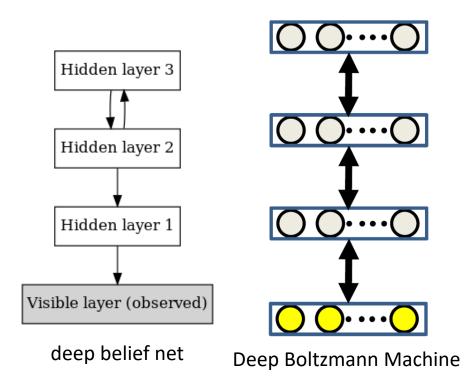
$$\quad \nabla_{w_{ij}} L(W) = \frac{1}{N_P K} \sum_{v \in P} v_{0i} h_{0j} - \frac{1}{M} \sum_{v \in P} v_{\infty i} h_{\infty j}$$

- No need to lift up the entire energy landscape!
 - Raising the neighborhood of desired patterns is sufficient

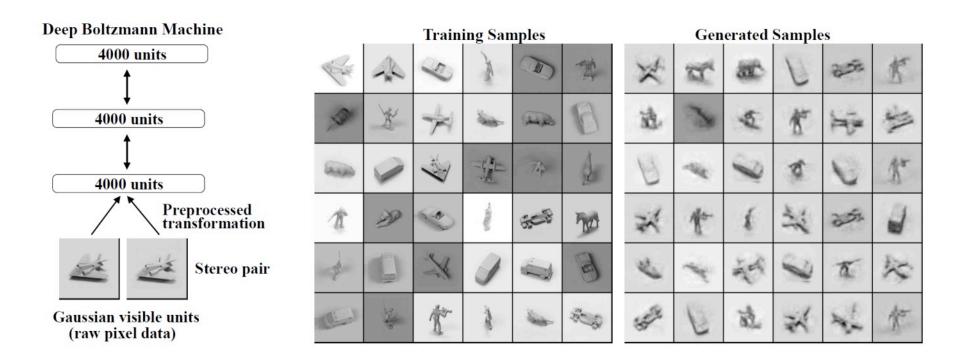


Deep Bolzmann Machine

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



Deep Bolzmann Machine



Summary

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