

# Variational Annealing of GANs: A Langevin Perspective

Chenyang Tao<sup>†</sup>  
chenyang.tao@duke.edu

Electrical & Computer Engineering, Duke University

Jun 12, 2019 @ ICML, Long Beach, CA, USA

Joint work with S Dai, L Chen, K Bai, J Chen, C Liu, G Bobashev and L Carin

# Outline

1. GAN Training & Likelihood Regularization
2. Variational Annealing From a Langevin Perspective
3. Experimental Results



# General Formulation of GANs

## Adversarial distribution matching

- A generator  $G(z), z \sim p(z)$ , a critic  $D(x)$
- A variational objective  $\mathbb{V}(\mu_d, \rho_G; D)$ 
  - computed using samples of **data**  $\mu_d$  and **model**  $\rho_G$
  - $d(\mu_d, \rho_G) = \max_D \mathbb{V}(\mu_d, \rho_G; D)$  defines a discrepancy metric
- Solve the minimax game

$$\min_G \max_D \mathbb{V}(\mu_d, \rho_G; D)$$

- ☺ **No explicit specification of likelihoods**
- ☹ **Brittle training, mode collapsing**

# A Concrete Example (That Is of Particular Interest)

## RKL GAN

- Let  $V_{\text{RKL}}(\rho, \mu; D) \triangleq \mathbb{E}_{X \sim \mu}[D(X)] + \mathbb{E}_{X' \sim \rho}[\log(-D(X'))]$
- $\text{KL}(\rho \parallel \mu) = \mathbb{E}_{X \sim \rho}[\log \frac{\rho(X)}{\mu(X)}] \Leftrightarrow \max_D \{V_{\text{RKL}}(\mu, \rho; D)\}$

# Regularizing GANs with Likelihoods

$$\rho^* = \arg \min_{\rho} \{ \max_D \{ V_{\text{RKL}}(\rho, \mu; D) \} - \lambda \mathcal{R}_{\mu}(\rho) \}$$

With data likelihoods  $\mathcal{R}_{\mu}(\rho) = -\mathbb{E}_{\rho}[\log \mu]$  [Li, 2018]

- Promoting **plausible samples** (concentrate)

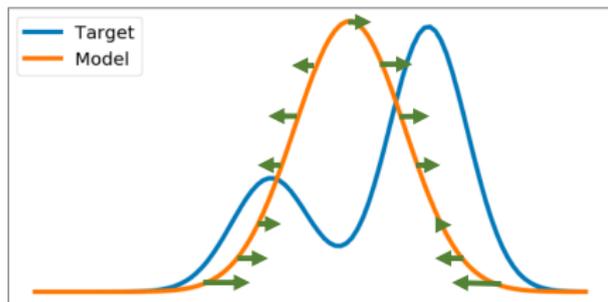
With model likelihoods  $\mathcal{R}_{\mu}(\rho) = \mathbb{E}_{\rho}[\log \rho]$  [Warde-Farley, 2017]

- Encouraging **sample diversity** (disperse)

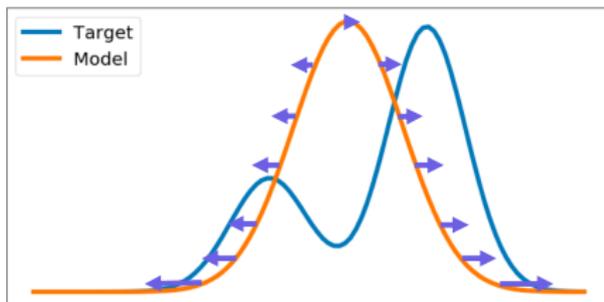
# Regularizing GANs with Likelihoods: Gradient View

$$\rho^* = \arg \min_{\rho} \{ \max_D \{ V_{\text{RKL}}(\rho, \mu; D) \} - \lambda \mathcal{R}_{\mu}(\rho) \}$$

## Likelihood Regularization



## Entropy Regularization



**We aim to provide theoretical groundings for such practices!**

# Preliminary

## Gibbs distribution

- $\mu_\beta(x) \propto \exp(-\beta\psi(x))$  is called a Gibbs distribution
  - $\psi(x)$  is the potential function
  - $\beta$  is the inverse temperature
- **Annealing** approaches the target distribution by gradually tuning  $\beta$  to avoid numerical difficulties



**Figure:** Illustration of annealed Gibbs distribution in 1-D.  $\beta = 1$  (green) is the target distribution,  $\beta < 1$ , mode covering and  $\beta > 1$ , mode seeking.

# The Link



1908

*Ito-Langevin* Diffusion

$$dX_t = -\nabla\psi(X_t)dt + \sqrt{2\beta^{-1}}dW_t$$



1914

*Fokker-Plank* Equation

$$\partial_t \rho = \nabla \cdot (\rho \nabla \psi) + \beta^{-1} \Delta \rho$$



2014

Generative Adversarial Net

$$\min_G \left\{ \max_D \left\{ \mathbb{E}_{X \sim p_d} [D(X)] - \mathbb{E}_{X' \sim p_G} [\ln(-D(X'))] \right\} + \lambda \ln p_d(X) \right\}$$

# The Link



1908

*Ito-Langevin* Diffusion

$$dX_t = -\nabla\psi(X_t)dt + \sqrt{2\beta^{-1}}dW_t$$



1914

*Fokker-Plank* Equation

$$\partial_t \rho = \nabla \cdot (\rho \nabla \psi) + \beta^{-1} \Delta \rho$$

**They All Minimize The *Free Energy*:  $\mathcal{F}_\psi(\rho; \beta) \triangleq \beta \mathbb{E}_\rho[\psi] + \mathbb{E}_\rho[\ln \rho]$**



2014

Generative Adversarial Net

$$\min_G \left\{ \max_D \left\{ \mathbb{E}_{X \sim p_d} [D(X)] - \mathbb{E}_{X' \sim p_G} [\ln(-D(X'))] \right\} + \lambda \ln p_d(X) \right\}$$

# The Link



1908

*Ito-Langevin* Diffusion

$$dX_t = -\nabla\psi(X_t)dt + \sqrt{2\beta^{-1}}dW_t$$



1914

*Fokker-Plank* Equation

$$\partial_t \rho = \nabla \cdot (\rho \nabla \psi) + \beta^{-1} \Delta \rho$$

**The Solution is Given by The Gibbs Distribution  $\mu_{\psi, \beta}(x) \propto e^{-\beta\psi(x)}$**



2014

Generative Adversarial Net

$$\min_G \left\{ \max_D \left\{ \mathbb{E}_{X \sim p_d} [D(X)] - \mathbb{E}_{X' \sim p_G} [\ln(-D(X'))] \right\} + \lambda \ln p_d(X) \right\}$$

# Effect of Likelihood Regularization & Its Implications

## Equivalence of Likelihood & Entropy Regularization

$$\rho^* = \arg \min_{\rho} \{ \max_D \{ V_{\text{RKL}}(\rho, \mu; D) \} - \lambda \mathcal{R}_{\mu}(\rho) \}$$

- For likelihood regularization

$$\mathcal{R}_{\mu}(\rho) = -\mathbb{E}_{\rho}[\log \mu] \Rightarrow \rho_{\text{lik}}^*(x) \propto \exp(-(\lambda + 1)\psi(x))$$

- For entropy regularization

$$\mathcal{R}_{\mu}(\rho) = \mathbb{E}_{\rho}[\log \rho] \Rightarrow \rho_{\text{ent}}^*(x) \propto \exp(-(\lambda + 1)^{-1}\psi(x))$$

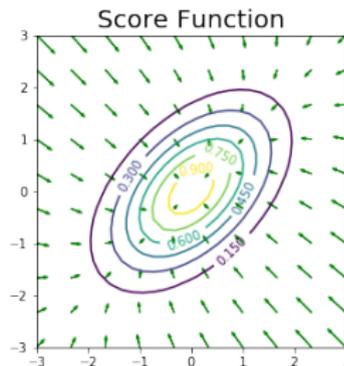
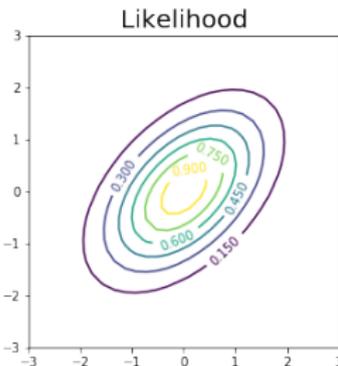
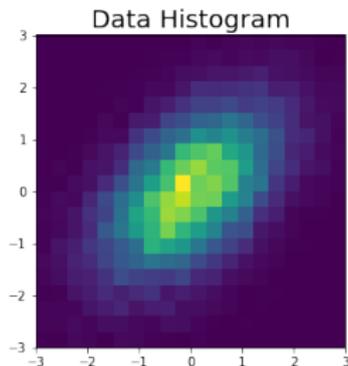
## Implications

- Likelihood/Entropy regularization bias the target distribution!

# Replacing The Likelihood with Score Function Estimate

## Challenges and solutions

- **We don't have the likelihood for data. That said, we know**
- $\log \mu(G_\theta(z)) \Leftrightarrow \mathcal{R}_\mu(z) \triangleq G_\theta(z)^T \text{StopGrad}\{S_\mu(G_\theta(z))\}$ 
  - $S_\mu(x) = \nabla_x \log \mu(x)$  is the **(data) score function**
- Estimating score function is way **easier** than likelihood



# Experiments: Dynamic Annealing

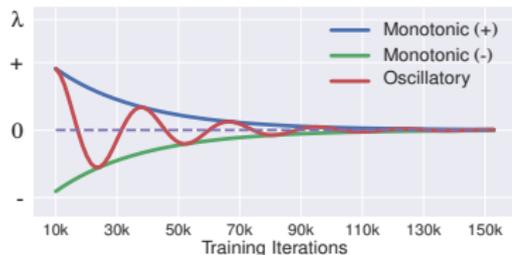


Figure 4. Different dynamic regularization schemes.

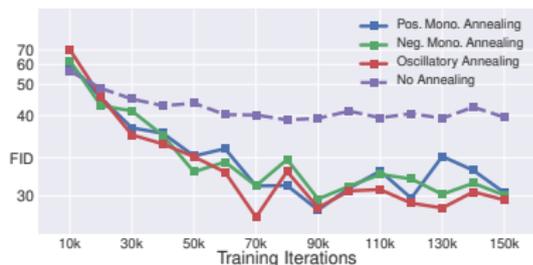


Figure 5. Learning dynamics with dynamic annealing.

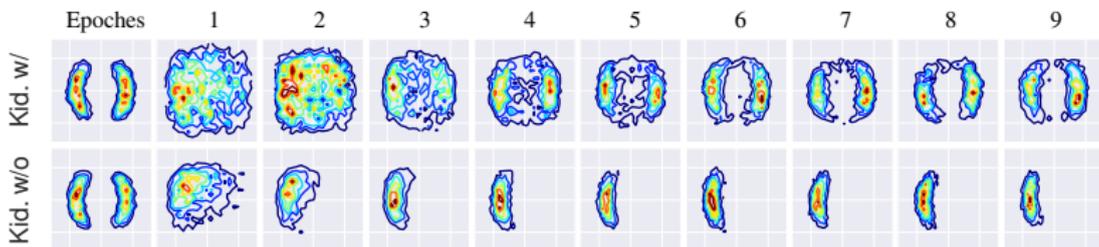


Figure 6. Learning from an unnormalized density to sample the kidney distribution. Top left: target distribution; bottom left: model distribution initialization; 'w/' with variational annealing; 'w/o' without annealing.

# Experiments: Quantitative & Qualitative Results

Table 1. Quantitative results for variational annealing on Cifar10.

$\lambda$	Static Annealing											Dynamic Annealing		
	-50	-10	-1	-0.1	-0.01	0	0.01	0.1	1	10	50	PMA	NMA	OA
Inception score (higher is considered better)														
RKL-GAN	6.24	6.37	6.35	6.33	6.35	6.25	6.24	6.35	6.41	6.19	6.17	6.56	<b>7.08</b>	7.05
JSD-GAN	6.68	6.84	6.64	6.35	6.61	6.29	6.67	6.30	6.93	6.48	6.22	6.80	<b>6.99</b>	6.96
W-GAN	5.77	6.14	6.29	6.86	6.62	5.93	6.22	6.54	5.95	6.00	6.00	<b>6.95</b>	6.92	6.91
FID score (lower is considered better)														
RKL-GAN	38.4	34.5	36.7	36.5	37.0	36.5	37.2	36.1	38.8	36.0	37.3	34.4	29.2	<b>28.9</b>
JSD-GAN	34.9	30.9	35.19	36.6	33.0	37.4	33.5	34.9	30.7	32.75	34.7	30.9	31.0	<b>29.1</b>
W-GAN	44.1	40.6	38.6	31.4	30.4	42.8	39.43	33.6	41.4	41.6	40.2	29.3	29.8	<b>29.0</b>

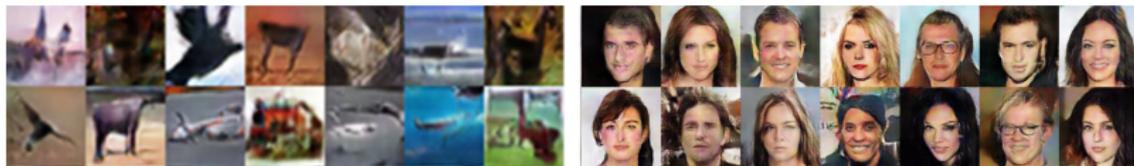


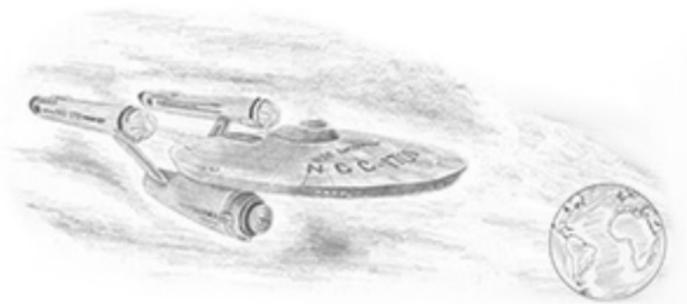
Figure 3. Cifar10 and CelebA generation results with negative static annealing.

Thank you.

*Welcome to our poster #10 @ Pacific Ballroom tonight.*



ICML | 2019  
Long Beach, CA, USA



*The authors would like to thank Prof D Waxman for fruitful discussions.*