
Supplementary Material for “Variational Annealing of GANs: A Langevin Perspective”

C. Tao¹, S. Dai¹, L. Chen¹, K. Bai¹, J. Chen^{1,2}, C. Liu^{1,3}, R. Zhang¹, G. Bobashev⁴, L. Carin¹
¹Duke, ²Fudan, ³Tsinghua, ⁴RTI

A Proof of Theorem 2.1

Proof.

The Fokker-Planck equation of Langevin diffusion. Since the derivation is quite involved, we refer the readers to the work of Garcia-Palacios (2007) for a complete treatment.

The stationary solution of Langevin system. One can easily verify that

$$\partial_t \rho_s(x, t) = \nabla \cdot \rho_s \nabla \psi + \beta^{-1} \Delta \rho_s = 0. \quad (\text{S1})$$

The stationary solution solves the variational problem. Let $\rho^*(x)$ be a minimum of functional $F_\mu(\rho; \beta)$, and $\epsilon(x)$ is an arbitrary function that vanishes on the boundary $\partial\mathcal{X}$, then for any number s close to 0, $F_\mu(\rho; \beta) \leq F_\mu(\rho + s\epsilon; \beta)$. The term $s\epsilon$ is the variation of the function ρ . Substituting $\rho + s\epsilon$ for y in the functional $F_\mu(y; \beta)$, the result is a function of s , $\Psi(s) = F_\mu(\rho + s\epsilon; \beta)$. Since the functional $F_\mu(\rho; \beta)$ has a minimum for $\rho = \rho^*$, the function $\Psi(s)$ has a minimum at $s = 0$ and thus,

$$\Psi'(0) = \left. \frac{d}{ds} F_\mu(\rho^* + s\epsilon; \beta) \right|_{s=0} = 0$$

Differentiating under the integral sign, we find

$$\begin{aligned} & \left. \frac{d}{ds} F_\mu(\rho^* + s\epsilon; \beta) \right|_{s=0} \\ &= \int_{\mathcal{X}} \frac{d}{ds} \{ \beta \psi(x) [\rho^*(x) + s\epsilon(x)] \\ & \quad + [\rho(x) + s\epsilon(x)] \log(\rho^*(x) + s\epsilon(x)) \} \Big|_{s=0} dx \\ &= \beta \int_{\mathcal{X}} \epsilon(x) \psi(x) dx + \int_{\mathcal{X}} \epsilon(x) [\log \rho^*(x) + 1] dx \\ &= \int_{\mathcal{X}} \epsilon(x) [\beta \psi(x) + \log \rho^*(x) + 1] dx = 0 \end{aligned}$$

According to the multidimensional version of the fundamental lemma of calculus of variations (Gelfand and Fomin (1963), pp.9 Lemma 1), $f(x) = \beta \psi(x) + \log \rho^*(x) + 1 \equiv 0$. Hence, $\rho^*(x) \propto \exp(-\beta \psi(x))$. \square

B Derivation for Eqn (4-5)

$$\begin{aligned} F_\mu(\rho; \beta) &= \text{KL}(\rho \parallel \mu) + (1 - \beta) \mathbb{E}_\rho[\log \mu] \\ &= \beta \text{KL}(\rho \parallel \mu) \\ & \quad + (1 - \beta) \mathbb{E}_\rho[\log \rho - \log \mu] \\ & \quad + (1 - \beta) \mathbb{E}_\rho[\log \mu] \\ &= \beta \text{KL}(\rho \parallel \mu) + (1 - \beta) \mathbb{E}_\rho[\log \rho] \\ &= \beta \text{KL}(\rho \parallel \mu) + (1 - \beta) S(\rho) \end{aligned}$$

C Proof of Corollary 2.2

Proof. We only need to $\beta_{\text{lik}} = 1 + \lambda$ and $\beta_{\text{ent}} = 1/(1 + \lambda)$ respectively for the two types of regularizations. For the likelihood regularization, based on (4) we have

$$\beta_{\text{lik}}(1 - \beta_{\text{lik}}^{-1}) = \lambda \Rightarrow \beta_{\text{lik}} = \lambda + 1,$$

and similarly for entropy regularization we have

$$\beta_{\text{ent}}^{-1} - 1 = \lambda \Rightarrow \beta_{\text{ent}} = 1/(1 + \lambda)$$

from (5). This concludes our proof. \square

D Continuity Equation

In the following, we omit the dependency on space-time to avoid notational clutter. By the divergence theorem, a general continuity equation takes the following differential form:

$$\partial_t \rho + \nabla \cdot \mathbf{j} = \sigma, \quad (\text{S2})$$

where ρ is the quantity of substance per unit volume, and \mathbf{j} is the flux of the substance and σ is the generation/absorption rate per unit volume per unit time. When \mathbf{v} is a velocity field that describes the flow \mathbf{j} , then $\mathbf{j} = \rho \mathbf{v}$. Since the mass of probability distribution is conserved, we have $\sigma(x, t) \equiv 0$. Plugging these terms into the continuity equation gives us

$$\partial_t \rho + \nabla \cdot (\rho \mathbf{v}) = 0. \quad (\text{S3})$$

E Proof of Theorem 2.3

Proof. Notice

$$\begin{aligned} & \text{KL}(\rho_t \parallel \mu) \\ &= \text{KL}(\rho_t \parallel \mu_\beta) + \mathbb{E}_{X' \sim \rho_t} [\log \mu_\beta(X') - \log \mu(X')] \\ &\leq \text{KL}(\rho_t \parallel \mu_\beta) + \delta_\beta. \end{aligned}$$

To get the $\mathcal{O}(t^{-1})$ convergence rate for first term, simply apply Theorem 4.5 from Liu (2017). This result can be improved to exponential decay assuming stronger regularity conditions on μ . For instance, when $\text{KL}(\cdot \parallel \mu)$ is geodesically strongly convex (*e.g.*, when the density of μ is strongly log-concave on \mathcal{X} (Villani (2008), Theorem 17.15)) on the 2-Wasserstein space \mathcal{P}_2 (Villani (2008), Definition 6.4), the gradient flow ρ_t of $\text{KL}(\cdot \parallel \mu)$ will converge exponentially (Villani (2008), Theorem 23.25 & 24.7). \square

F Further Notes on RKL-GAN as Gradient Flow

It is helpful to further the understanding of the dynamics from a gradient flow perspective. Consider minimizing $F(q) : \mathcal{P} \rightarrow \mathbb{R}$ be a function on the space of probability measure $\mathcal{P}(\mathcal{X})$, in our case, the anneal RKL. At each step, we would like to find a proper perturbation $\partial_t q_t$ for updating q . Formally, $\partial_t q_t$ should be an element in the tangent space of $\mathcal{P}(\mathcal{X})$ at q .

To practically describe $\partial_t q_t$, people notice that a $\partial_t q_t$ is related to a bunch of velocity fields $\{v_t\}$ on \mathcal{X} through the continuity equation (6), and a one-to-one relation can be established by choosing a particular v_t with the minimal norm in the bunch, when a proper norm is defined for velocity fields on \mathcal{X} (*e.g.*, when the norm is taken as the one of $L^2_q(\mathcal{X}; \mathbb{R}^n)$, the one-to-one relation (unique existence of v_t) is guaranteed by *e.g.*, Villani (2008), Theorem 13.8; Ambrosio et al. (2008), Theorem 8.3.1, Proposition 8.4.5.) Also note that the description with v_t automatically satisfies the restriction on $\partial_t q_t$: $\int_{\mathcal{X}} \partial_t q_t dx = -\int_{\mathcal{X}} \nabla \cdot (q_t v_t) dx = -\int_{\partial \mathcal{X}} q_t v_t \cdot d\vec{S} = 0$, where \vec{S} is the infinitesimal directed surface area on the boundary $\partial \mathcal{X}$. With this description, we can write the directional derivative of the RKL F with respect to $\partial_t q_t$, or v_t : $\frac{d}{ds} F(q + s \partial_t q_t)|_{s=0} = \mathbb{E}_q[-\nabla \cdot v_t + \Psi \cdot v_t]$, as is shown in Theorem 3.1 of Liu and Wang (2016).

There is more that the velocity field description could provide. Let \mathcal{T} be the set of all representative v_t 's, which is a linear space. With a proper inner product so that \mathcal{T} is a Hilbert space, $T_q \mathcal{P}(\mathcal{X})$ will be a Hilbert space due to the one-to-one relationship (which is now an isometric isomorphism), which (roughly) means that $\mathcal{P}(\mathcal{X})$ is a Riemannian manifold. With the

Riemannian structure, gradient of F on $\mathcal{P}(\mathcal{X})$ can be defined, which is characterized by $\frac{d}{ds} F(q + s \partial_t q_t)|_{s=0} = \langle \text{grad } F, \partial_t q_t \rangle_{T_q \mathcal{P}(\mathcal{X})}$, or

$$\text{grad } F = \max \cdot \arg \max \left\{ \frac{d}{ds} F(q + s \partial_t q_t)|_{s=0} \right\}. \quad (\text{S4})$$

When \mathcal{T} is taken as $L^2_q(\mathcal{X}; \mathbb{R}^n)$ and $\mathcal{P}(\mathcal{X})$ as 2-Wasserstein space, we have $\text{grad } F(q) = -\Psi - \nabla \log q$ (Villani (2008), Theorem 23.18; Ambrosio et al. (2008), Example 11.1.2), *i.e.*, the tangent vector on the gradient flow of $F(q)$. Note that the Langevin dynamics (1) is also known as the dynamics of the gradient flow of $F(q)$ on 2-Wasserstein space Jordan et al. (1998), since it produces the same $\partial_t q_t$ through the Fokker-Planck equation as the dynamics of $\text{grad } F$ through the continuity equation (6). Since the $\log q_t$ term here is intractable, it is estimated in (annealed) RKL-GAN through a function approximator (*e.g.* neural net) updated by stochastic gradient descent (the critic updates).

The quality of the approximation affects the empirical convergence rate, as the asymptotic convergence rate is an upper bound on the improvement, and it only holds if the updates are exactly aligned with functional gradients. So the maximizing step can be understood as searching for the best descent direction for the generator update. Another choice of \mathcal{T} is \mathcal{H}^n where \mathcal{H} is the RKHS of a kernel k , as is done in Liu and Wang (2016); Liu (2017). (*Note that in this case, $\mathcal{T} \subset T_q \mathcal{P}(\mathcal{X})$, which means that the existence of v_t in \mathcal{T} is not guaranteed for any $\partial_t q_t \in T_q \mathcal{P}(\mathcal{X})$! Moreover, $\mathcal{P}(\mathcal{X})$ as a set is not defined in this case!*) The gradient in this case is then $\text{grad } F(q) = \mathbb{E}_q(x)[- \Psi(x)k(x, \cdot) + \nabla_x k(x, \cdot)]$, which is the tangent vector to the gradient flow of $F(q)$ on $\mathcal{P}(\mathcal{X})$ defined in Liu (2017). With the help of a kernel, the gradient here is tractable.

G Score Function Estimator (Fig. 2)

We compare different score function estimators with toy examples. We train a DAE with samples from the *banana distribution* $\log u(x_1, x_2) = -\frac{1}{2}(x_2^2/4 + (x_1 - x_2^2/4)^2)$. We set the noise level to $\sigma = 0.2$. In Figure 2 from main text, we compare the model samples drawn from the standard Langevin system using the DAE-score estimator $s_\sigma(x) = \sigma^{-2}(\phi_\sigma(x) - x)$ and the DAE-residual estimator $s'_\sigma(x) = -\nabla \|\phi_\sigma(x) - x\|_2^2$. The residual estimator is unable to faithfully recover the underlying distribution, for practical training setups, often collapsing samples on modes that not even necessarily align with the real data modes.

Table S1: Network architectures used for Cifar10 experiments. BN: batch normalization, lReLU: Leaky ReLU.

Denoiser	Generator	Discriminator
Input feature f	Input z	Input x
MLP output 2048, lReLU, BN	MLP output 2048, ReLU, BN	3 × 3 conv. 32 lReLU, stride 1, BN
MLP output 2048, lReLU, BN		4 × 4 conv. 64 lReLU, stride 2, BN
MLP output 2048, lReLU, BN	4 × 4 deconv. 256 ReLU, stride 2, BN	4 × 4 conv. 128 lReLU, stride 2, BN
MLP output 2048, lReLU, BN	4 × 4 deconv. 128 ReLU, stride 2, BN	4 × 4 conv. 256 lReLU, stride 2, BN
MLP output 2048, lReLU, BN	4 × 4 deconv. 64 ReLU, stride 2, BN	4 × 4 conv. 512 lReLU, stride 2, BN
MLP output 2048, lReLU, BN	3 × 3 deconv. 3 ReLU, stride 1, BN	Output feature f with shape 512 × 2 × 2
MLP output 2048, Reshape to 512 × 2 × 2		MLP output 1

Table S2: Network architectures used for CelebA experiments. BN: batch normalization, lReLU: Leaky ReLU.

Denoiser	Generator	Discriminator
Input feature f	Input z	Input x
MLP output 2048, lReLU, BN	MLP output 2048, ReLU, BN	4 × 4 conv. 16 lReLU, stride 2, BN
MLP output 2048, lReLU, BN		4 × 4 conv. 32 lReLU, stride 2, BN
MLP output 2048, lReLU, BN	4 × 4 deconv. 256 ReLU, stride 2, BN	4 × 4 conv. 64 lReLU, stride 2, BN
MLP output 2048, lReLU, BN	4 × 4 deconv. 128 ReLU, stride 2, BN	4 × 4 conv. 128 lReLU, stride 2, BN
MLP output 2048, lReLU, BN	4 × 4 deconv. 64 ReLU, stride 2, BN	4 × 4 conv. 256 lReLU, stride 2, BN
MLP output 2048, lReLU, BN	4 × 4 deconv. 32 ReLU, stride 2, BN	4 × 4 conv. 512 lReLU, stride 2, BN
MLP output 2048, lReLU, BN	4 × 4 deconv. 3 ReLU, stride 2, BN	Output feature f with shape 512 × 2 × 2
MLP output 2048, Reshape to 512 × 2 × 2		MLP output 1

H More on generative flows

Typically, the shift-scale flow are given as

$$z_{k+1} = t(z_k) + s(z_k) \odot z_k, \quad (S5)$$

where $t(z), s(z)$ respects the causal dependency of an auto-regressive flow, *i.e.*, $[f(z)]_d = f_d([z]_{<d})$. This ensures Jacobians are triangular by design. Equation (S5) is a so-called *forward flow* which is fast in generation but slow in evaluation, because the backward process (*e.g.*, $x \rightarrow z$) resembles a Gaussian elimination process. Such trade-offs are common in more expressive generative flows. MAF instead exploits the *backward flow* described in the main text.

While it is tempting to train a flow ν directly towards μ , there is one caveat. To optimize $\text{KL}(\nu \parallel \mu)$, one must sample from ν , evaluate the log-likelihood $\log \nu(x)$ and then back-propagate the gradient through it. As discussed above, computational efficiency of sampling and evaluation for popular choices of generative flows are generally at odds with each other. Coupling the two will always invoke the back-propagation of the more costly part. On the other hand, the scheme we have developed avoids back-propagating the more challenging part via amortizing it through a free-form generator.

I Additional Remarks & Discussions

I.1 Comments on U-GAN

While the U-GAN can be much more flexible than the paired generative flow, we found that more expressive flow allows much quicker learning for U-GAN.

I.2 Discussion of The "Stabilizing GAN Training with Langevin Dynamics" Paper

During the paper review process, the work of Ramak-ers et al. (2017) has brought to our attention by one of the reviewers, which also leverages Langevin dynamics in GAN training. However, this work used Langevin dynamics to evolve the model parameter rather than regularizing the GAN training objective. As such, their development is orthogonal to the variational annealing perspective discussed here.

I.3 Wasserstein GAN, Optimal Transport and Likelihood

In our experiments we have showed that the proposed VA improved GAN training across the board. While for non-RKL GANs this observation is not understood beyond heuristics, there is some theories relating optimal transport and likelihood-based learning, see Zhang et al. (2018) for details.

J Detailed Experimental Setups and Additional Results

J.1 Variational Annealing

In this experiment, we use the DFM implementation from <https://github.com/pfnet-research/chainer-gan-lib> as our codebase. The original model corresponds to the JSD-GAN results reported, and we changed the objective function to derive the RKL-GAN. For the W-GAN experiments, we adapted the code of the WGAN-GP implementation from the same repository, which uses gradient penalty to enforce the Lipschitz constraint. We have used the same hyper-parameter settings from the DFM repository, which potentially explained the fact that the JSD-GAN delivered the best performance in our experiments, as the released DFM code was fine-tuned for this objective. The Adam optimizer with a learning rate of 10^{-4} was used, and all models are trained for 150k iterations with batch-size 64. Detailed model architecture for Cifar10 (32×32) is summarized in Table S1. Similar architectures are used for high-res CelebA (128×128), with additional deconvolutional layers (Table S2). In Figure S2 we provide CelebA samples trained with and without variational annealing.

Dynamic annealing schemes In the dynamic annealing experiments, we found that the linear annealing

$$\lambda_t = \text{sign}(\lambda_0) \max\{|\lambda_0|(1 - t/T_{\max}), 0\}$$

worked better than the exponential annealing

$$\lambda_t = \text{sign}(\lambda_0)|\lambda_0| \exp(-\tau t)$$

among monotonic schemes, where T_{\max} and τ are hyper-parameters controlling the descent rate. T_{\max} was set to be slightly smaller than the total number of iterations to finalize the training, and τ was set chosen so that at end of training $\lambda_t \rightarrow 0$. As for the oscillatory annealing, we evaluated stepwise (discontinuous) and sinuous (continuous) designs, and the latter worked much better. The performance also depends on $|\lambda_0|$, in our experiments, we have tested $|\lambda_0| \in \{0.1, 1, 10\}$ and reported the best result. During the entire training, λ_t loops over for five cycles.

J.2 Unnormalized GAN

The toy model *kidney distribution* is given by

$$\log u_{\text{kid}}(x) = \frac{1}{2}s_1(x)^2 - \log(s_2(x) + s_3(x)),$$

Table S3: Results for Bayesian Logistic regression.

Dataset	Features	Train	Test
Cancer	32	285	284
Heart	13	135	135
German	20	500	500
Sonar	60	104	104

where $s_1(x) = \frac{\|x\| - 2}{0.4}$, $s_2(x) = \exp(-\frac{1}{2}[\frac{x_1 - 2}{0.6}]^2)$ and $s_3(x) = \exp(-\frac{1}{2}[\frac{x_1 + 2}{0.6}]^2)$. In Figure S1 we also visualize the training dynamics using the *banana distribution*: $p(x_1, x_2) \propto \exp(-\frac{1}{2}((x_1 - (x_2/2))^2) + (x_2^2/4))$ using negative annealing, the sampler distributes mass more evenly at a low temperature (stronger negative annealing), then gradually consolidates to the target distribution as the annealing strength diminishes.

J.2.1 MAF

We have used the *tensorflow-probability* library to implement MAF. More specifically, we stacked 8 *shift_and_scale* MAF blocks with [512, 512] hidden layers, each followed by an order flipping permutation layer. We used the same annealing designs discussed above.

J.2.2 Datasets

Table S3 summarizes the datasets used in our Bayesian Logistic regression experiment.

J.3 Reinforcement Learning

In the RL experiment, we consider the *swimmer* problem from *rllab*. We used a three layer feed-forward net as our policy net, both hidden layers has 128 units. For the Q-net, we also used a three layer 128 hidden unit feedforward net. We applied a four block *shift_and_scale* MAF as our proposal density, each has two hidden layers with size 512. We set batch size to 128, learning rate to 0.0004 and maintained a 10^6 replay buffer.

References

- Ambrosio, L., Gigli, N., and Savaré, G. (2008). *Gradient flows: in metric spaces and in the space of probability measures*. Springer Science & Business Media.
- Garcia-Palacios, J. (2007). Introduction to the theory of stochastic processes and brownian motion problems. *arXiv preprint cond-mat/0701242*.
- Gelfand, I. and Fomin, S. (1963). *Calculus of variations*. Prentice-Hall (transl. from Russian).
- Jordan, R., Kinderlehrer, D., and Otto, F. (1998). The variational formulation of the Fokker-Planck equation. *SIAM journal on mathematical analysis*, 29(1):1–17.

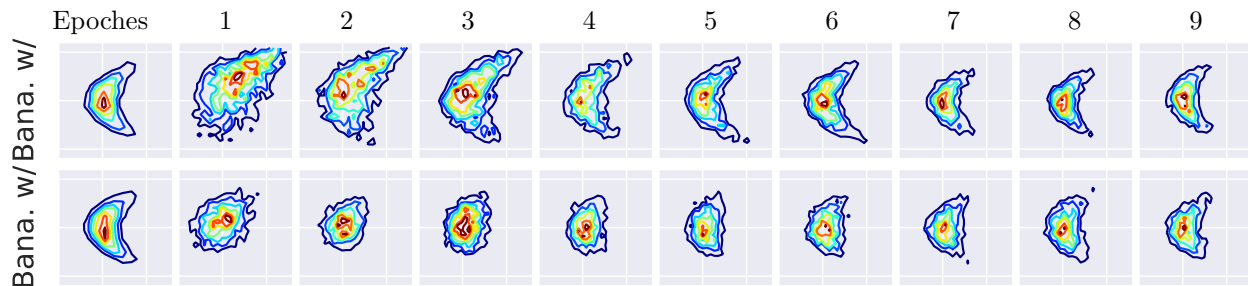


Figure S1: Learning from an unnormalized density to sample the banana distribution.
 With Variational Annealing Without Variational Annealing

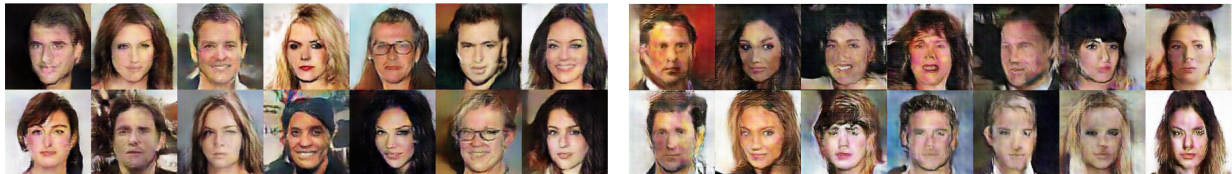


Figure S2: Comparison of CelebA samples trained with and without variational annealing.

- Liu, Q. (2017). Stein variational gradient descent as gradient flow. In *NIPS*, pages 3118–3126.
- Liu, Q. and Wang, D. (2016). Stein variational gradient descent: A general purpose bayesian inference algorithm. In *NIPS*.
- Ramakers, J., Harmeling, S., and Kollmann, M. (2017). Stabilizing generative adversarial networks using langevin dynamics. In *NIPS Workshop*.
- Villani, C. (2008). *Optimal transport: old and new*, volume 338. Springer Science & Business Media.
- Zhang, L., E, W., and Wang, L. (2018). Monge-amp\ere flow for generative modeling. *arXiv preprint arXiv:1809.10188*.