
Stabilizing Training of Generative Adversarial Networks through Regularization

Kevin Roth

Department of Computer Science
ETH Zürich
kevin.roth@inf.ethz.ch

Aurelien Lucchi

Department of Computer Science
ETH Zürich
aurelien.lucchi@inf.ethz.ch

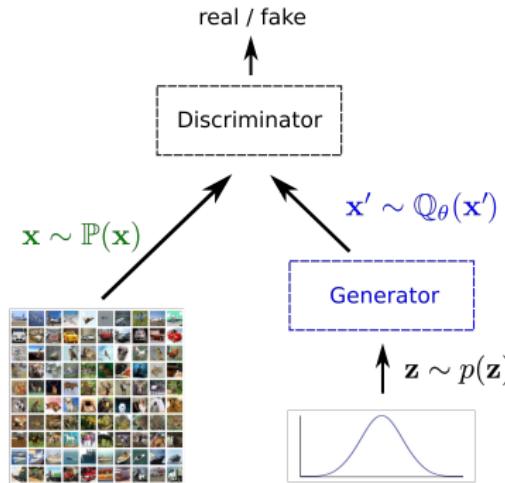
Sebastian Nowozin

Microsoft Research
Cambridge, UK
sebastian.Nowozin@microsoft.com

Thomas Hofmann

Department of Computer Science
ETH Zürich
thomas.hofmann@inf.ethz.ch

GANs



GAN objective [Goodfellow et al.]

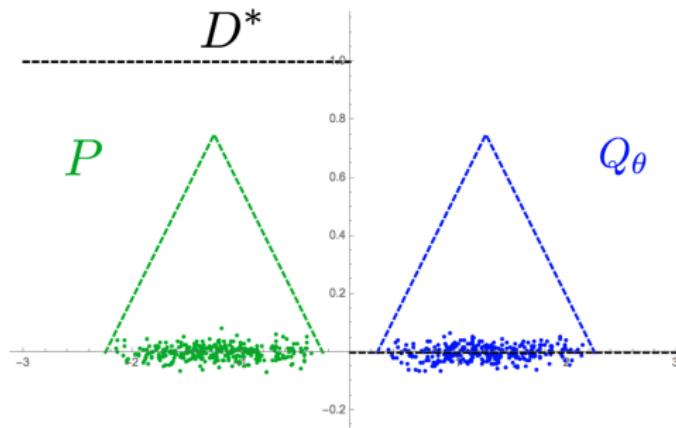
$$\min_G \max_D \mathbb{E}_{x \sim P} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))]$$

$$= \min_G D_{JS}(P \| Q_\theta) \text{ for Bayes-optimal } D^*(x)$$

GANs are exciting, but ...
they're also notoriously hard to train!

Problem: *dimensional mismatch* or

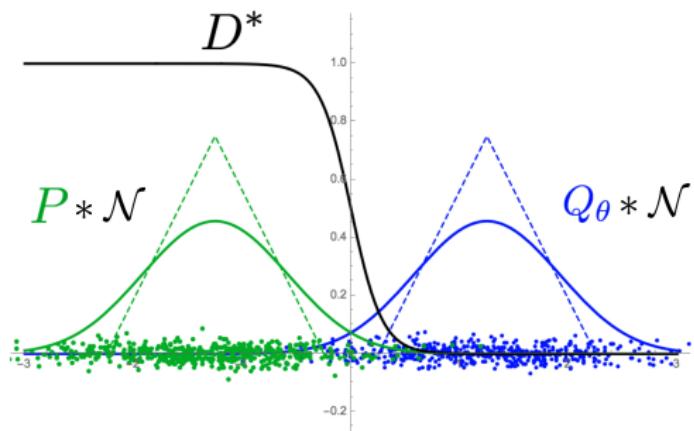
$$\text{supp}(\textcolor{green}{P}) \cap \text{supp}(\textcolor{blue}{Q}_\theta) = \emptyset$$



=> undefined f -divergence $D_f(\textcolor{green}{P} || \textcolor{blue}{Q}_\theta) = ???$

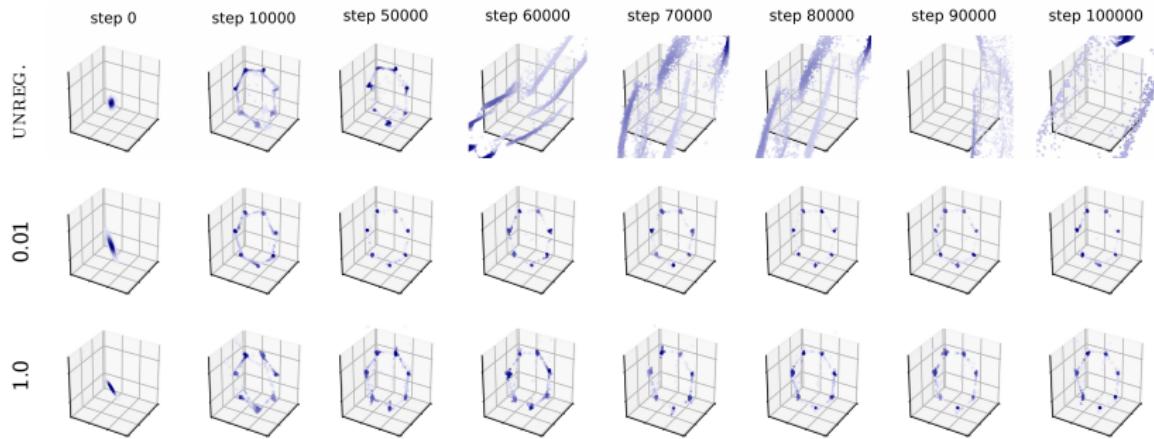
Solution: Adding Noise (Convolving Densities)

$$x \rightarrow x + \xi, \quad \xi \sim \mathcal{N}(0, \gamma)$$



=> Regularizing Discriminator

Submanifold Mixture



Gradient-norm regularized GANs for different levels of regularization γ

Background

From f -divergences ...

Measure “similarity” between distributions P and Q (that possess densities p and q , absolutely continuous w.r.t. base measure $\mathrm{d}x$),

$$D_f(P\|Q) := \int_{\mathcal{X}} q(x) f\left(\frac{p(x)}{q(x)}\right) \mathrm{d}x,$$

where $f : \mathbb{R}_+ \rightarrow \mathbb{R}$ is convex and satisfies $f(1) = 0$.

Legendre - Fenchel duality:

$$f(u) = \sup_{t \in \mathrm{dom}_{f^c}} \{tu - f^c(t)\}$$

parametrizing f through its tangents or supporting hyperplanes.

... to f -GANs [Nowozin et al.]

$$\begin{aligned}
 D_f(P\|Q) &= \int_{\mathcal{X}} q(x) \sup_{t \in \text{dom}_{f^c}} \left\{ t \frac{p(x)}{q(x)} - f^c(t) \right\} dx \\
 &\geq \sup_{\psi \in \Psi} \left(\int_{\mathcal{X}} p(x) \psi(x) dx - \int_{\mathcal{X}} q(x) f^c(\psi(x)) dx \right) \\
 &= \sup_{\psi \in \Psi} (\mathbf{E}_{x \sim P} [\psi(x)] - \mathbf{E}_{x \sim Q_\theta} [f^c(\psi(x))]),
 \end{aligned}$$

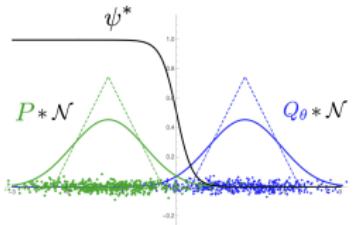
where Ψ is an arbitrary class of discriminants $\psi : \mathcal{X} \rightarrow \mathbb{R}$.

Bound is tight for the **optimal discriminant** [Nguyen et al.]

$$\psi^*(x) = f' \left(\frac{p(x)}{q(x)} \right)$$

Regularization

Training with Noise



- Practitioner:

Explicitly adding noise $\xi \sim \Lambda = \mathcal{N}(0, \gamma \mathbb{I})$ to $\mathbf{x} \sim \mathbb{P}, \mathbb{Q}$

- Theory: Convolving Distributions

$$\begin{aligned}\mathbf{E}_{\mathbb{P}} \mathbf{E}_{\Lambda} [h(\psi(\mathbf{x} + \xi))] &= \int h(\psi(\mathbf{x})) \int p(\mathbf{x} - \xi) \lambda(\xi) d\xi d\mathbf{x} \\ &= \mathbf{E}_{\mathbb{P} * \Lambda} [h \circ \psi]\end{aligned}$$

=> Convolved *f*-GAN Objective

$$F(\mathbb{P} * \Lambda_\gamma, \mathbb{Q} * \Lambda_\gamma; \psi) = \mathbf{E}_{\mathbb{P} * \Lambda_\gamma} [\psi] - \mathbf{E}_{\mathbb{Q} * \Lambda_\gamma} [f^c \circ \psi]$$

Analytic Approximation

- For small noise variance γ we can Taylor expand ψ around $\xi = 0$ [Bishop 95]:

$$\psi(\mathbf{x} + \xi) = \psi(\mathbf{x}) + [\nabla \psi(\mathbf{x})] \xi + \frac{1}{2} \xi^T [\nabla^2 \psi(\mathbf{x})] \xi + \mathcal{O}(\xi^3)$$

- Third-order approximation

$$F_\gamma = F_0 + \frac{\gamma}{2} \{ \mathbf{E}_{\mathbb{P}} [\triangle \psi] - \mathbf{E}_{\mathbb{Q}} [\triangle (f^c \circ \psi)] \} + \mathcal{O}(\gamma^2)$$

- Interpret: Laplace $\triangle = \text{Tr}(\nabla^2)$ measures how much ψ and $f^c \circ \psi$ differ from their local average

Simplification

- Chain-rule:

$$\Delta(f^c \circ \psi) = (f^{c''} \circ \psi) \cdot \|\nabla \psi\|^2 + (f^{c'} \circ \psi) \Delta \psi$$

- Property of **optimal discriminant ψ^***

$$(f^{c'} \circ \psi^*) d\mathbb{Q} = d\mathbb{P}$$

=> **Convenient cancellation** at $\psi = \psi^* + \mathcal{O}(\gamma)$ [Bishop 95]:

$$\mathbf{E}_{\mathbb{P}} [\Delta \psi^*] - \mathbf{E}_{\mathbb{Q}} [\Delta(f^c \circ \psi^*)] = -\mathbf{E}_{\mathbb{Q}} [(f^{c''} \circ \psi^*) \cdot \|\nabla \psi^*\|^2]$$

Regularized f -GAN

$$F_\gamma(\mathbb{P}, \mathbb{Q}_\theta; \psi) = \mathbf{E}_{\mathbb{P}}[\psi] - \mathbf{E}_{\mathbb{Q}_\theta}[f^c \circ \psi] - \frac{\gamma}{2} \Omega_f(\mathbb{Q}_\theta; \psi)$$

$$\Omega_f(\mathbb{Q}_\theta; \psi) := \mathbf{E}_{\mathbb{Q}_\theta} \left[(f^{c''} \circ \psi) \cdot \|\nabla \psi\|^2 \right]$$

- “soft Lipschitz” constraint
- non-negative weighting function $f^{c''} \geq 0$
- lower variance compared to explicitly adding noise
- **easy to implement**
- **can train indefinitely without collapse!**

Algorithm 1 Regularized f -GAN (with annealing). Default values $\gamma_0 = 2.0$, $\alpha = 0.005$, $n_\psi = 1$

Require: Initial noise variance γ_0 , annealing decay rate α , number of discriminator update steps n_ψ per generator iteration, minibatch size m , number of training iterations T

Require: Initial discriminator parameters ω_0 , initial generator parameters θ_0

for $t = 1, \dots, T$ **do**

$\gamma \leftarrow \gamma_0 \cdot \alpha^{t/T}$ # annealing

for $1, \dots, n_\psi$ **do**

Sample minibatch of real data $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\} \sim \mathbb{P}$.

Sample minibatch of latent variables from prior $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\} \sim p(\mathbf{z})$.

$$F(\omega, \theta) = \frac{1}{m} \sum_{i=1}^m \left[\psi_\omega(\mathbf{x}^{(i)}) - f^c(\psi_\omega(G_\theta(\mathbf{z}^{(i)}))) \right]$$

$$\Omega_f(\omega, \theta) = \frac{1}{m} \sum_{i=1}^m f^{c''}(\psi_\omega(G_\theta(\mathbf{z}^{(i)}))) \left\| \nabla_{\tilde{\mathbf{x}}} \psi_\omega(\tilde{\mathbf{x}}) \Big|_{\tilde{\mathbf{x}}=G_\theta(\mathbf{z}^{(i)})} \right\|^2$$

$$\omega \leftarrow \omega + \nabla_\omega \left(F(\omega, \theta) - \frac{\gamma}{2} \Omega_f(\omega, \theta) \right) \text{ # gradient ascent}$$

end for

Sample minibatch of latent variables from prior $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\} \sim p(\mathbf{z})$.

$$F(\omega, \theta) = \frac{1}{m} \sum_{i=1}^m -f^c(\psi_\omega(G_\theta(\mathbf{z}^{(i)})))$$

$$\theta \leftarrow \theta - \nabla_\theta F(\omega, \theta) \text{ # gradient descent}$$

end for

The gradient-based updates can be performed with any gradient-based learning rule. We used Adam in our experiments.

tf code → github.com/rothk

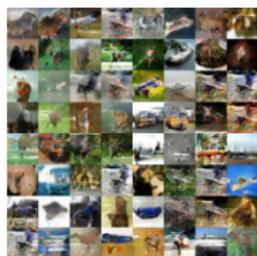
```
# -----
#   JS-Regularizer
# -----
def Discriminator_Regularizer(self, D1, D1_logits, D1_arg, D2, D2_logits, D2_arg):
    grad_D1_logits = tf.gradients(D1_logits, D1_arg)[0]
    grad_D2_logits = tf.gradients(D2_logits, D2_arg)[0]
    grad_D1_logits_norm = tf.norm(tf.reshape(grad_D1_logits, [self.batch_size,-1]), axis=1, keep_dims=True)
    grad_D2_logits_norm = tf.norm(tf.reshape(grad_D2_logits, [self.batch_size,-1]), axis=1, keep_dims=True)
    reg_D1 = tf.multiply(tf.square(1.0-D1), tf.square(grad_D1_logits_norm))
    reg_D2 = tf.multiply(tf.square(D2), tf.square(grad_D2_logits_norm))
    return tf.reduce_mean(reg_D1 + reg_D2)
```

Experimental Results

Regularization vs. Explicitly Adding Noise

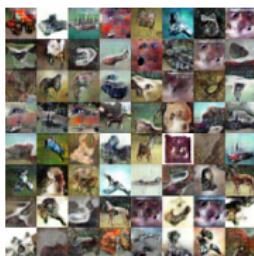
UNREGULARIZED

0.01



EXPLICIT NOISE

0.1



1.0

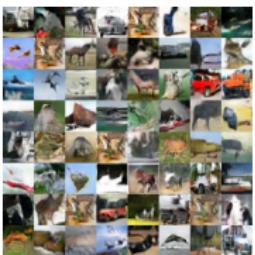


REGULARIZED

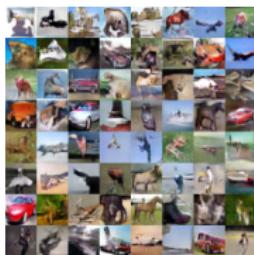
0.001



0.01



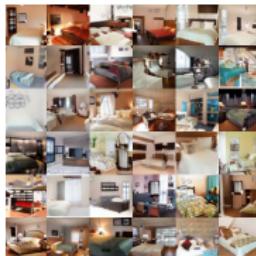
0.1



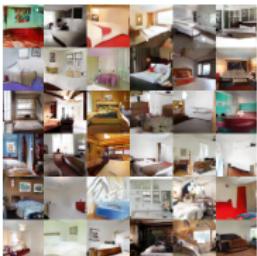
1.0



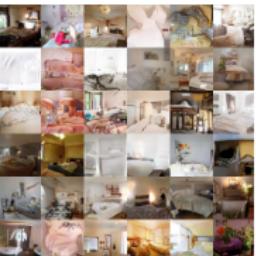
Stability across various architectures



RESNET



DCGAN



NO NORMALIZATION



TANH



UNREG.



0.5 -

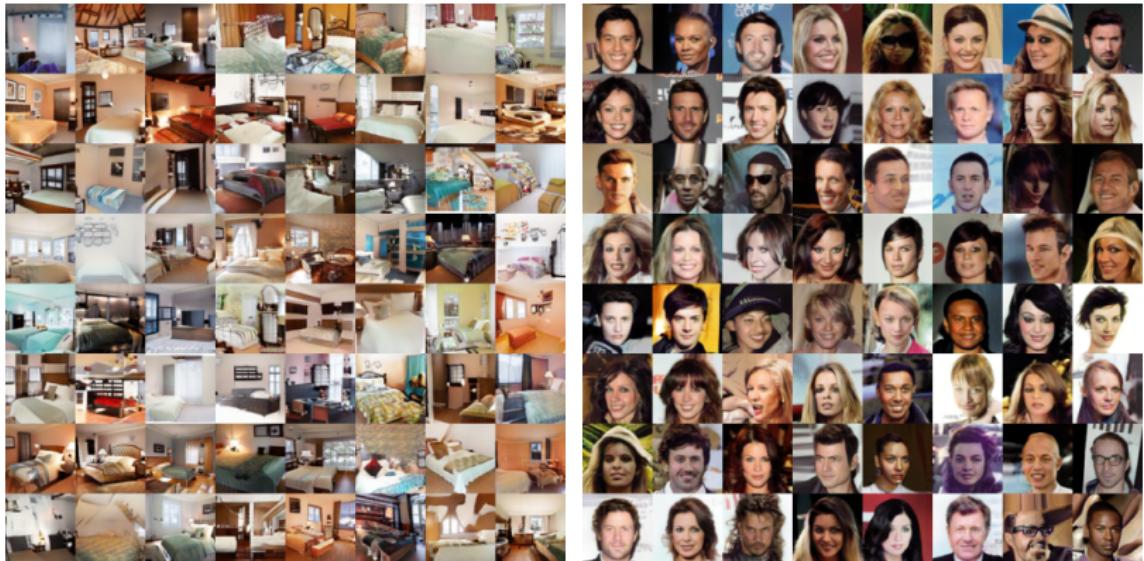


ANNEALING



- 2.0

Sample quality and diversity



Cross-Testing Protocol

Regularized $\gamma = 0.1$

		True Cond.	
		Pos.	Neg.
Pred.	Pos.	0.9688	0.0002
	Neg.	0.0312	0.9998

Cross-testing: FP: 0.0

Unregularized

		True Cond.	
		Pos.	Neg.
Pred.	Pos.	1.0	0.0013
	Neg.	0.0	0.9987

Cross-testing: FP: 1.0

- Regularized D classifies unregularized G's samples as fake
 - Unregularized D classifies samples of regularized G as real
- => Regularized GAN generalizes better!

Outlook

- **Applications:** What to actually do with GANs?
- **Evaluation:** Convergence, generalization, mode collapse, sample quality, latent space semantics, WGAN vs. f -GAN, ...
- **Optimization:** Saddle-points in minimax problems, non-conservative vector-fields, ...

References

- Roth K., Lucchi A., Nowozin S., Hofmann T., Stabilizing Training of Generative Adversarial Networks through Regularization, NIPS '17
- Goodfellow I., Pouget-Abadie J., Mirza M., Xu B., Warde-Farley D., Ozair S., Courville A., Bengio Y., Generative Adversarial Networks, NIPS '14
- Nowozin S., Cseke B., Tomioka R., f-GAN: Training generative neural samplers using variational divergence minimization, NIPS '16
- Nguyen X., Wainwright M.J., Jordan M.I. Estimating divergence functionals and the likelihood ratio by convex risk minimization, IEEE Transactions on Information Theory, 2010
- Bishop C.M., Training with noise is equivalent to Tikhonov regularization, Neural computation, 1995
- An G., The Effects of Adding Noise During Backpropagation Training on Generalization Performance, Neural computation, 1996