

# Analysing Failure Modes in Unsupervised Image-to-Image Translation

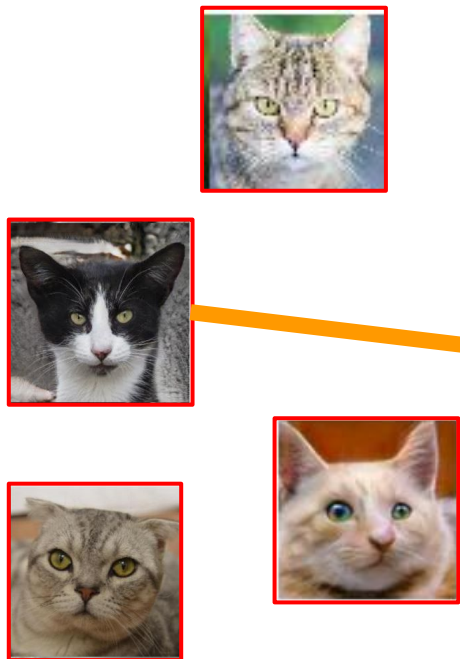
Doctoral Qualifying Oral Exam  
Boston University 2021  
Ben Usman

# Presentation Plan

1. Problem and motivation
2. Existing solutions and possible failure modes
3. Tools for analysing these failure modes in prior work
  - a. **“Generalization and Equilibrium in Generative Adversarial Nets”**  
by Arora et al., PMLR 2017.
  - b. **“Training Generative Adversarial Networks with Limited Data”**  
by Karras et al., NeurIPS 2020.
  - c. **“Risk Bounds for Unsupervised Cross-Domain Mapping with IPMs”**  
by Galanti et al., JMLR 2021. // 2017-2021

# Task

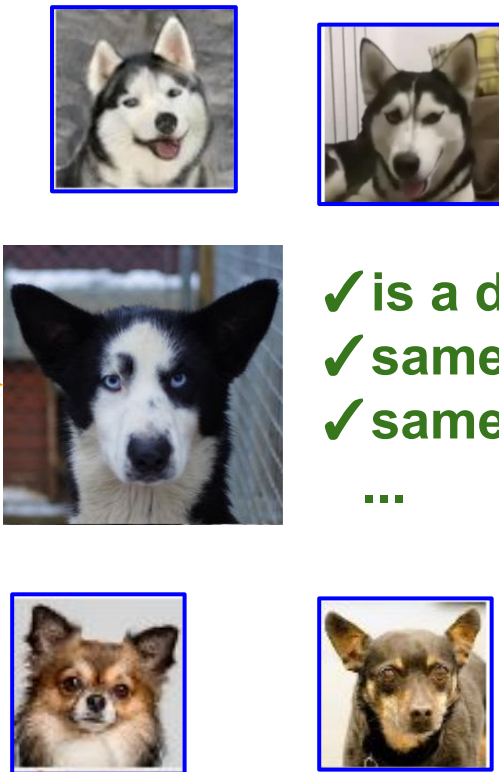
Source Samples (Cats)



F



Target Samples (Dogs)



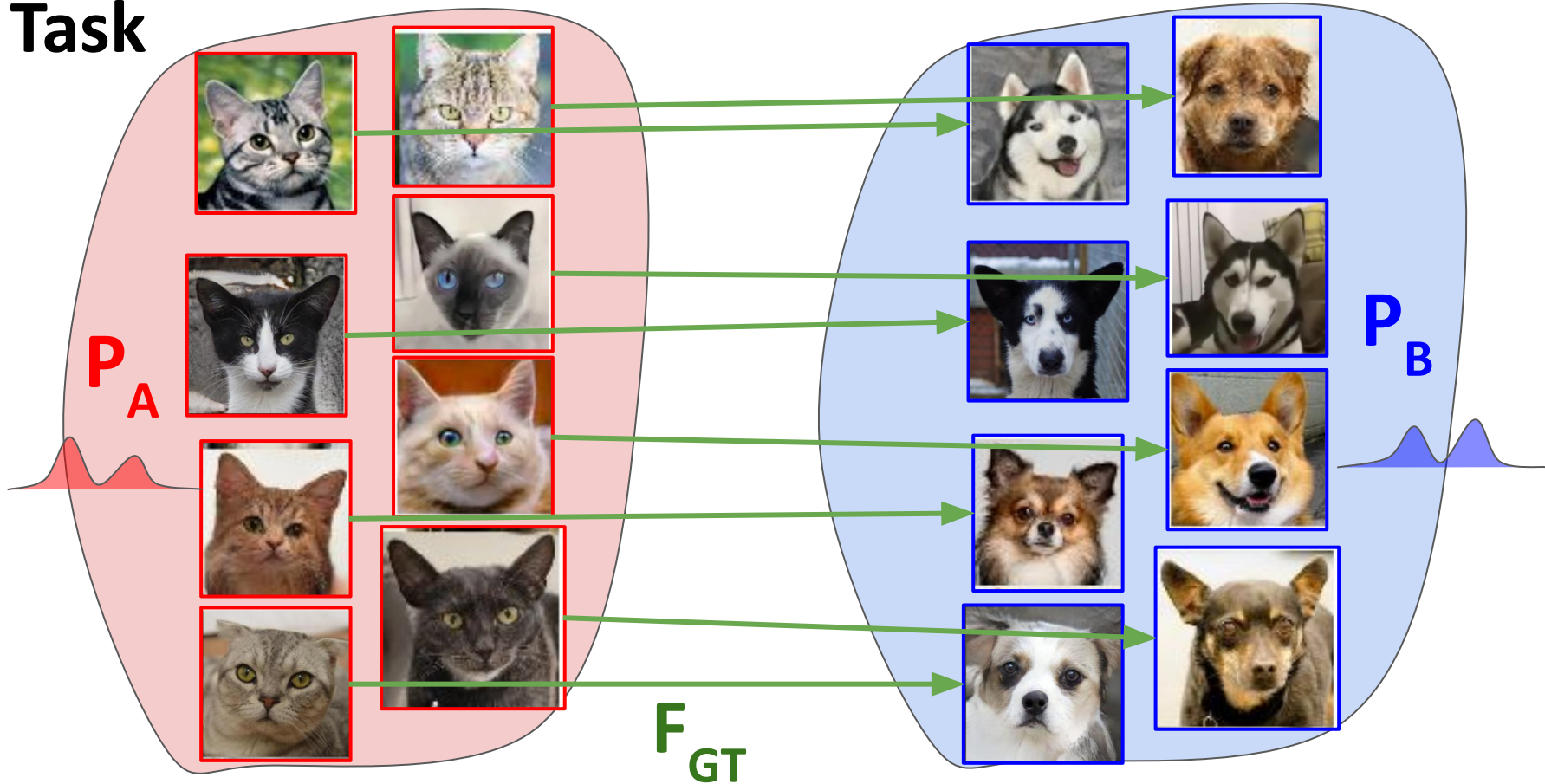
✓ is a dog  
✓ same coat color  
✓ same pose

...

# Task

Source Distribution

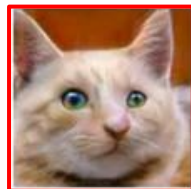
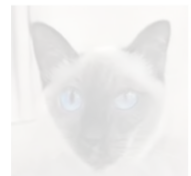
Target Distribution



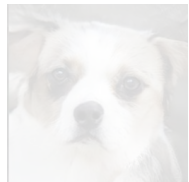
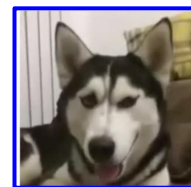
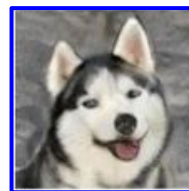
Ground truth 1-to-1 cross-domain mapping.

# Task

## Source Samples

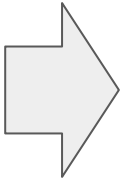


## Target Samples

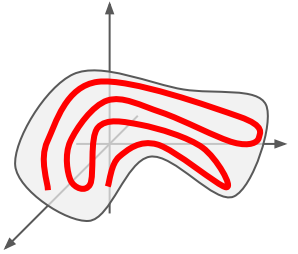
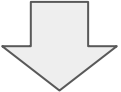
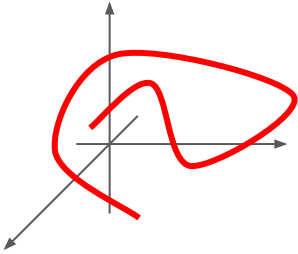


**Goal:**  
reconstruct  $F$  from  
unpaired samples

# Task



1D manifold



2D manifold

# How to find a good F?

Source Samples



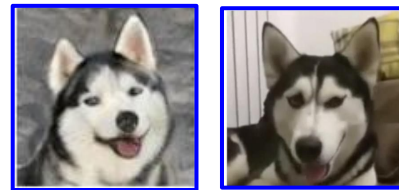
$$A = \{a_i\}$$

Translated Source Samples



$$F(A) = \{F(a_i) : a_i \in A\}$$

Target Samples



$$B = \{b_j\}$$



minimize  
statistical  
distance  
 $d(F(A), B)$



$$\min_{F \in \mathcal{F}} d(F(A), B) + R(F)$$

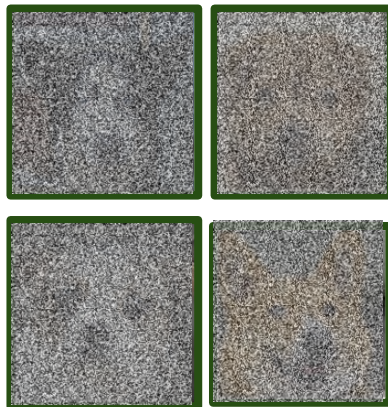
# How to find a good F? - what we expect

Source Samples



$F$   
 $t=0$   
→

Translated Source Samples



**HIGH**  
statistical  
distance  
 $d(F(A), B)$

Target Samples



... optimizing F ...

$$\min_{F \in \mathcal{F}} d(F(A), B) + R(F)$$

$F$   
 $t=T$   
→



**LOW**  
statistical  
distance  
 $d(F(A), B)$





# Why care about this problem?

This is an **unsupervised generative** problem that has **GT outputs!**

As a result, we are learning a **neural data model**, but can reason about its **correctness** and the **prediction error vs GT outputs** (e.g. L2).

In contrast,

- in GANs - there are **no expected outputs**
- in classification/regression - often no need to **model data**.

# Presentation Plan

1. Problem and motivation
2. Existing solutions and possible failure modes
3. Tools for analysing these failure modes in prior work
  - a. **“Generalization and Equilibrium in Generative Adversarial Nets”**  
by Arora et al., PMLR 2017.
  - b. **“Training Generative Adversarial Networks with Limited Data”**  
by Karras et al., NeurIPS 2020.
  - c. **“Risk Bounds for Unsupervised Cross-Domain Mapping with IPMs”**  
by Galanti et al., JMLR 2021. // 2017-2021

# What could go wrong?

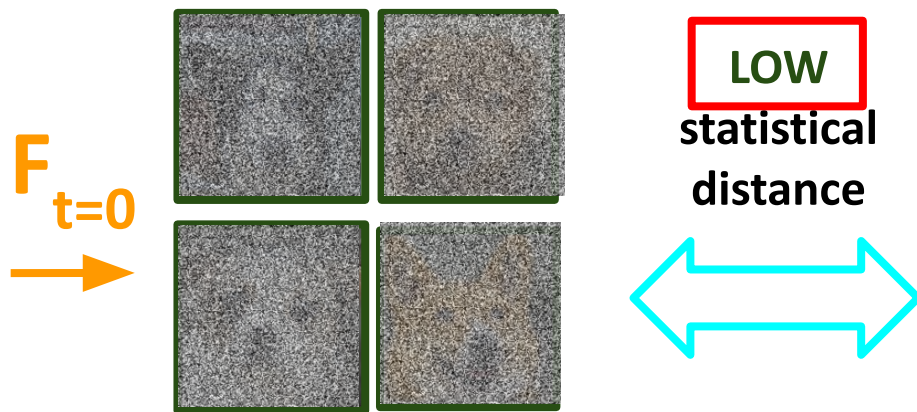
## I: The statistical distance is too weak

Source Samples



$$A = \{a_i\}$$

Translated Source Samples



$$F(A) = \{F(a_i) : a_i \in A\}$$

Target Samples

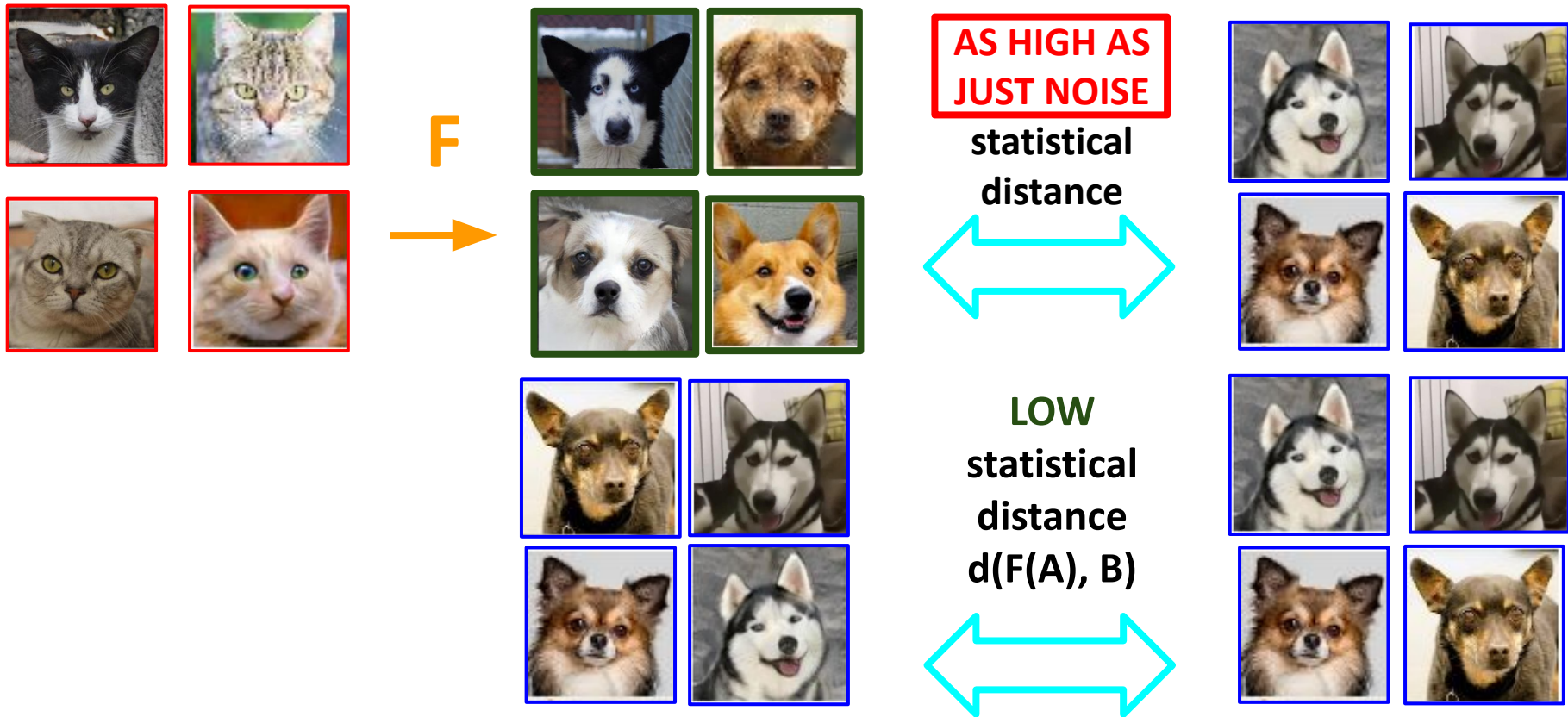


$$B = \{b_i\}$$

$$\min_{F \in \mathcal{F}} d(F(A), B) + R(F)$$

# What could go wrong?

I: The statistical distance is too strong



# What could go wrong?

## II: The stat distance is too sharp (hard to optimize)



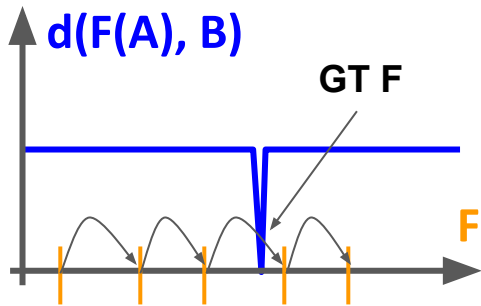
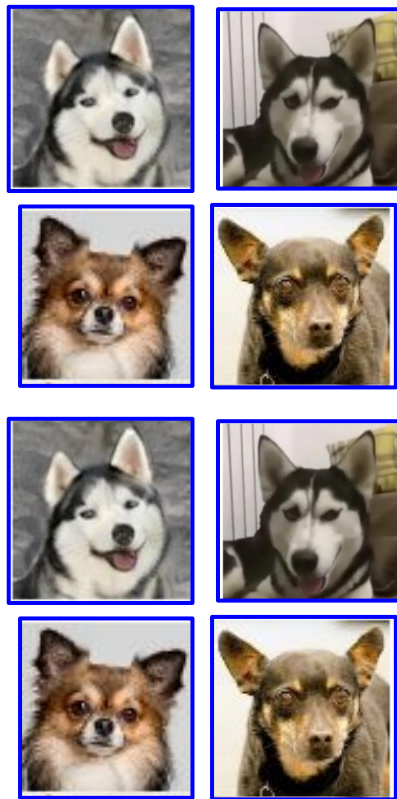
**F**  
→



**EQUALLY  
HIGH**  
statistical  
distance



**EQUALLY  
HIGH**  
statistical  
distance



# What could go wrong?

## III: The final mapping is nonsensical



**NO SEMANTIC  
CORRESPONDENCE**

$$\min_{F \in \mathcal{F}} d(F(A), B) + R(F)$$

# Selected prior work

1. **“Generalization and Equilibrium in Generative Adversarial Nets”** by Sanjeev Arora, Rong Ge, Yingyu Liang, Tengyu Ma, Yi Zhang, Proceedings International Conference on Machine Learning (PMLR) 2017.
2. **“Training Generative Adversarial Networks with Limited Data”** by Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, Timo Aila; Advances in Neural Information Processing Systems (NeurIPS) 2020.
3. **“Risk Bounds for Unsupervised Cross-Domain Mapping with IPMs”** by Tomer Galanti, Sagie Benaim, Lior Wolf; JMLR 2021. // *paper series (2017-2021)*

# Why these papers?

Introduce important theoretical tools to understand related problems:

- **Chernoff bound and  $\epsilon$ -cover method** to estimate sample complexity of adversarial statistical distances
- **$\epsilon$ -approximate Nash equilibrium** to analyse the existence of the solution to the adversarial alignment problem
- **Markov operators and group structure of augmentations** to estimate statistical distances between distributions under data augmentations
- **unsupervised bias-variance tradeoff and Rademacher complexity** to relate the prediction error with the alignment error and the complexity of the function class



# Other background papers

I also use there papers / books to provide context / refer for proofs

- “Simple Strategies for Large Zero-Sum Games with Applications to Complexity Theory” by Lipton & Young, STOC’94
- “Foundations of Machine Learning” Mohri, Rostamizadeh, Talwalkar, 2nd Edition, 2018
- “Towards Principled Methods for Training Generative Adversarial Networks”, Arjovsky & Bottou, ICLR’17
- “Stabilizing Training of Generative Adversarial Networks through Regularization”, Roth et al, NeurIPS’17
- “Which Training Methods for GANs do actually Converge?”, Mescheder et al., ICML’18
- “Kernel of CycleGAN as a Principle homogeneous space”, Moriakov et al., ICLR’20
- “Guiding the One-to-One Mapping in CycleGAN via Optimal Transport’, Lu et al., AAAI’19

# Other background papers

I have “backup slides” covering these papers as well in the end:

- "Table for estimating the goodness of fit of empirical distributions", Smirnov, Annals of Mathematical Statistics '48 - introduces KS-test
- "A Kernel Two-Sample Test", Gretton et al, JMLR'12 - introduces MMD test
- "Permutation tests for equality of distributions in high-dimensional settings", Hall & Tajvidi, Biometrika'02; "Multivariate Generalizations of the Wald-Wolfowitz and Smirnov Two-Sample Tests", Friedman & Rafsky, Ann Stat 79 - multivariate extensions of non-parametric tests
- "Wasserstein GAN" Arjovsky et al, ICML'17 - introduces WGAN objective
- "Are GANs Created Equal? A Large-Scale Study", Lucic et al., NeurIPS'18; "Pros and Cons of GAN Evaluation Measures", Ali Borji, arxiv'18; "Improved Precision and Recall Metric for Assessing Generative Models", Kynkäänniemi et al., NeurIPS'19 - introduces FID, KID, IS, GAN-F1 score and compares them


# Other background papers

I have “backup slides” covering these papers as well in the end:

- “On the Decreasing Power of Kernel and Distance based Nonparametric Hypothesis Tests in High Dimensions”, Ramdas et al., AAAI'15 - shows that with a “fair alternative” MMD test has exponentially low power in higher dims
- “Revisiting Classifier Two-Sample Tests”, Lopez-Paz et al., ICLR'17 - compares the test power of the GAN-like objective to MMD/KS/other test
- “Reducing Noise in GAN Training with Variance Reduced Extragradient”, Chavdarova et al., NeuIPS'19

# How to choose the statistical distance?

**Def 1:** the statistical distance  $d(A, B)$  “generalizes”

$$\left| d(\mathcal{D}_{real}, \mathcal{D}_G) - d(\hat{\mathcal{D}}_{real}, \hat{\mathcal{D}}_G) \right| \leq \varepsilon \quad \text{(with exponentially high probability over the choice of } m \text{ samples as the number of samples increases)}$$


**Lem 1:** JSD and Wasserstein distances “do not generalize”!

$$\mathcal{N}(0, \frac{1}{d}I) = \mu \quad \Rightarrow \quad d_{JS}(\mu, \hat{\mu}) = \log 2, \quad d_W(\mu, \hat{\mu}) \geq 1.1$$

$$JS(p; q) = \frac{1}{2} \int \left( p \log \frac{2p}{p+q} + q \log \frac{2q}{p+q} \right) d\mu$$

(for  $q = 0$  almost everywhere)

$$\Pr[\forall i \in [m] \|y - x_i\| \geq 1.2] \geq 1 - m \exp(-\Omega(d)) \geq 1 - o(1)$$

$$d_W(\mu, \hat{\mu}) \geq 1.2 \Pr[\forall i \in [m] \|y - x_i\| \geq 1.2] \geq 1.1$$

(prob of all pairs of points in  $A$  and  $B$  being at least 1.2 away from each other does not decay fast enough with  $m$ )

# How to choose the statistical distance?

## Def 2: $F$ -divergence wrt $\phi$

$$d_{\mathcal{F},\phi}(\mu, \nu) = \sup_{D \in \mathcal{F}} \mathbb{E}_{x \sim \mu} [\phi(D(x))] + \mathbb{E}_{x \sim \nu} [\phi(1 - D(x))] - 2\phi(1/2)$$

**Lem 2:**  $m \geq \frac{cp\Delta^2 \log(LL\phi p/\epsilon)}{\epsilon^2}$ , we have with probability at least  $1 - \exp(-p)$

$D(X)$  outputs for different weights

$$|d_{\mathcal{F},\phi}(\hat{\mu}, \hat{\nu}) - d_{\mathcal{F},\phi}(\mu, \nu)| \leq \epsilon$$

**Proof:**

- all **D weights** can be approx. with err.  $< \mathbf{e}$  using a “covert set” of size  $\sim K = \lceil \log(1/e) / e \rceil$   
 $\Rightarrow$  worst cast D approx. error  $< \mathbf{e} * \mathbf{L}$
- for any “single fixed” D the est. error  $< \mathbf{A}$

# How to choose the statistical distance?

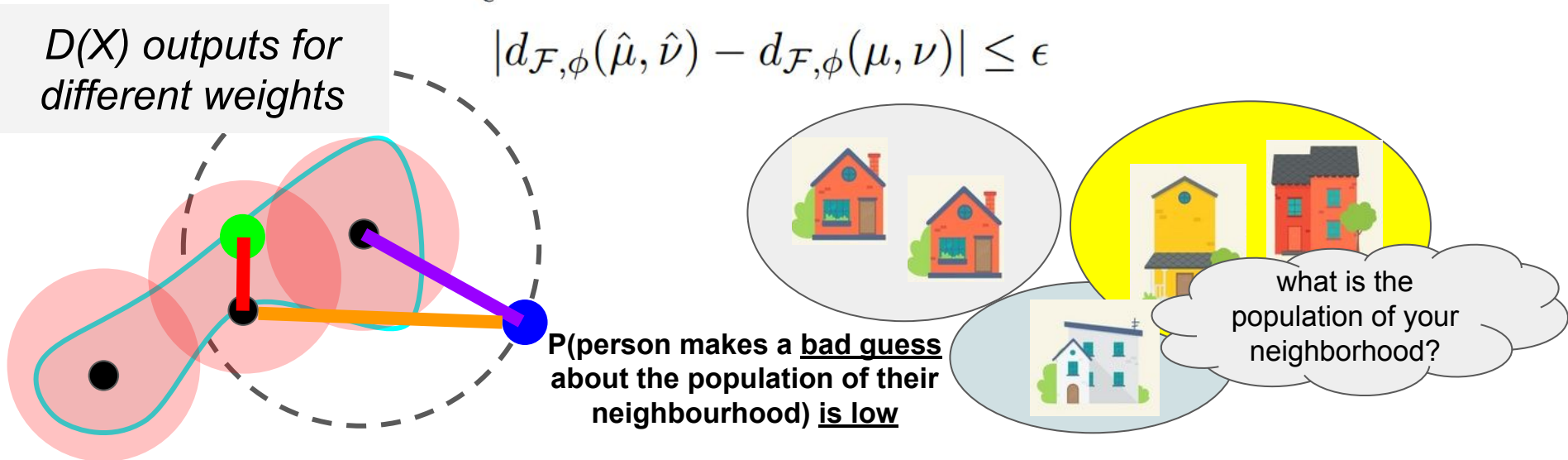
## Def 2: $F$ -divergence wrt $\phi$

$$d_{\mathcal{F},\phi}(\mu, \nu) = \sup_{D \in \mathcal{F}} \mathbb{E}_{x \sim \mu} [\phi(D(x))] + \mathbb{E}_{x \sim \nu} [\phi(1 - D(x))] - 2\phi(1/2)$$

**Lem 2:**  $m \geq \frac{cp\Delta^2 \log(LL_{\phi}p/\epsilon)}{\epsilon^2}$ , we have with probability at least  $1 - \exp(-p)$

$D(X)$  outputs for different weights

$$|d_{\mathcal{F},\phi}(\hat{\mu}, \hat{\nu}) - d_{\mathcal{F},\phi}(\mu, \nu)| \leq \epsilon$$



# How to choose the statistical distance?

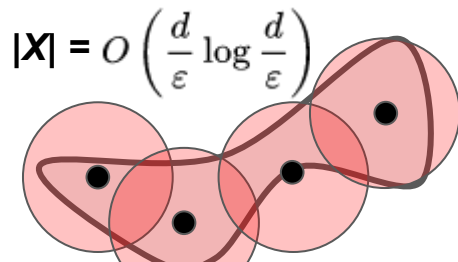
## Def 2: $F$ -divergence wrt $\phi$

$$d_{\mathcal{F},\phi}(\mu, \nu) = \sup_{D \in \mathcal{F}} \mathbb{E}_{x \sim \mu} [\phi(D(x))] + \mathbb{E}_{x \sim \nu} [\phi(1 - D(x))] - 2\phi(1/2)$$

**Lem 2:**  $m \geq \frac{cp\Delta^2 \log(LL_\phi p/\epsilon)}{\epsilon^2}$ , we have with probability at least  $1 - \exp(-p)$

$$|d_{\mathcal{F},\phi}(\hat{\mu}, \hat{\nu}) - d_{\mathcal{F},\phi}(\mu, \nu)| \leq \epsilon$$

$\epsilon$ -net method on discr weights



$$|\mathcal{X}| = O\left(\frac{d}{\epsilon} \log \frac{d}{\epsilon}\right)$$

$\Pr[x \geq E[x] + t] \leq e^{-2t^2/n}$  - Chernoff bound for  $x_i \sim [0, 1]$

applied to a single discriminator from the  $\epsilon$ -net (bounded outputs!)

$$\Pr\left[\left| \mathbb{E}_{x \sim \mu} [\phi(D_v(x))] - \mathbb{E}_{x \sim \hat{\mu}} [\phi(D_v(x))] \right| \geq \frac{\epsilon}{4}\right] \leq 2 \exp\left(-\frac{\epsilon^2 m}{2\Delta^2}\right)$$

$\Rightarrow$  (+ union bound) for  $m \geq \frac{Cp\Delta^2 \log(LL_\phi p/\epsilon)}{\epsilon^2}$  <sup>\*  $|\mathcal{X}|$</sup>  the error is  $< \epsilon/4$  and

within distance  $\epsilon/8LL_\phi$

$$\log |\mathcal{X}| \leq O(p \log(LL_\phi p/\epsilon))$$

$$\|v - v'\| \leq \epsilon/8LL_\phi \Rightarrow \left| \mathbb{E}_{x \sim \hat{\mu}} [\phi(D_{v'}(x))] - \mathbb{E}_{x \sim \hat{\mu}} [\phi(D_v(x))] \right| \leq \epsilon/8$$

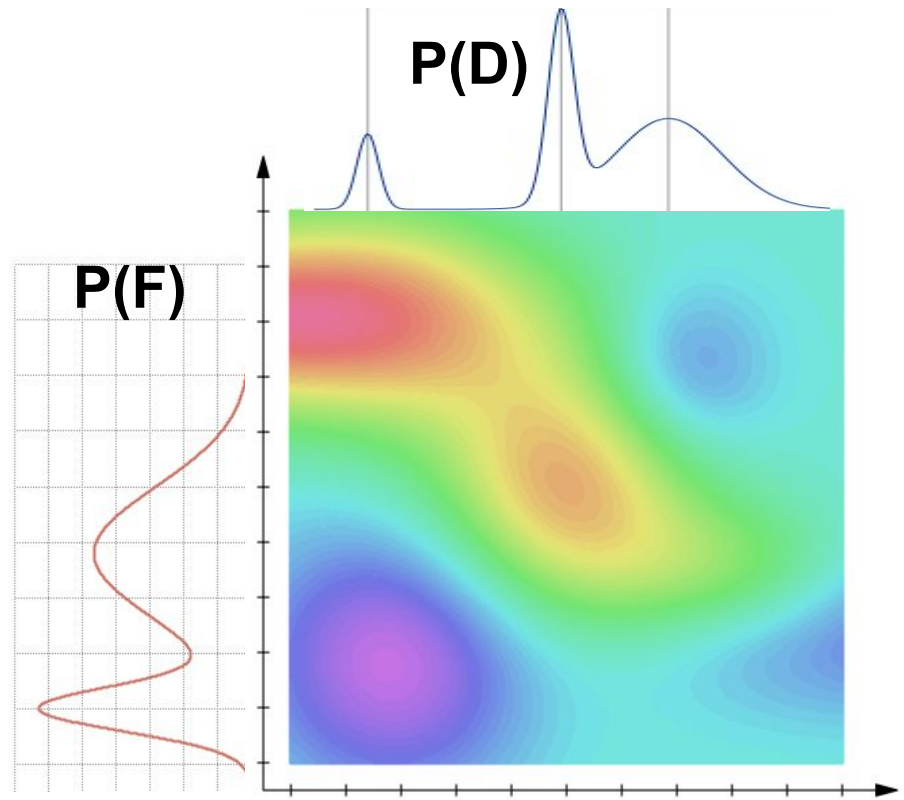
# Does the minimum exist?

$$\min_F \max_D D(X) - D(F(Y))$$

$\min h(F)$

	$D=D_A$		$D=D_B$	
$F=F_A$	2	←	-2	.3
↓			↑	
$F=F_B$	-1	→	1	.7
↑				

The “pure” equilibrium might not always exist, but a **mixed strategy** that yields an equilibrium **always exist!**  
 [Nash '50; Glicksberg '52]





**Def 3:**  $\epsilon$ -approximate equilibrium

$$\forall v \in \mathcal{V}, \quad \mathbb{E}_{u \sim \mathcal{S}_u} [F(u, v)] \leq V + \epsilon;$$

$$\forall u \in \mathcal{U}, \quad \mathbb{E}_{v \sim \mathcal{S}_v} [F(u, v)] \geq V - \epsilon.$$

**Th:** if  $p$ -parameter  $(k-1)$ -layer network can generate/discriminate each sample  $\Rightarrow \exists$   $k$ -layer  $G$  and  $D$  with  $A$  parameters that are in  $\epsilon$ -eq

**Proof:**

- 1) an infinite mixture of  $G_i(z) = x_i$ ,  $x_i \sim P(X)$  is mixed Nash eq.
- 2)  $K$ -sized epsilon-net over samples  $x_i$  and params of  $D$  gives small error  $\Rightarrow$  “subsampled”  $G'$  and  $D'$  are in  $\epsilon/2$ -eq
- 3) can approximate “subsampled”  $G'(x)$  with a neural  $G''(x)$  that “mixes” outputs of  $G_i$  with weights produced by a neural  $\epsilon/2$ -approx. “1-vs- $K$  indicator”  $h(z)$ , i.e.  $G''(x) = \sum_i G_i(z) * h_i(z)$

## Proof:

- 1) an infinite mixture of  $G_i(z) = x_i$  is Nash equilibrium
- 2)  $\epsilon/4LL'L_\phi$ -net (of size  $T$ ) over params of  $G$  and  $D$  gives (with high prob) error  $< \epsilon/2 \Rightarrow$  “subsampling”  $G'$  and  $D'$  are in  $\epsilon$ -eq
- 3) a 2-layer network  $h(z)$  can  $\delta$ -approx. a “multi-way step fn”
- 4) we build new  $G$  that “mixes” outputs of  $G'_i$  with weights produced by  $h(z)$ , it is  $\epsilon/2$ -away from “true mixture of  $G'_i$ ’s”

$$\begin{aligned} F^*(G, D') &\geq \mathbb{E}_{i \in [T], v \in D'} F(u_i, v) \\ &\quad - |F^*(G, D') - \mathbb{E}_{i \in [T], v \in D'} F(u_i, v)| \\ &\geq V - \epsilon/2 - 2\Delta \frac{\epsilon}{4\Delta} \\ &\geq V - \epsilon. \end{aligned} \quad \begin{aligned} F^*(G', D) &\leq \mathbb{E}_{i \in [T], u \in G'} F(u, v_i) \\ &\quad + |F^*(G', D) - \mathbb{E}_{i \in [T], u \in G'} F(u, v_i)| \\ &\leq V + \epsilon/2 + 2\Delta \frac{\epsilon}{4\Delta} \\ &\leq V + \epsilon. \end{aligned}$$

second half of the proof is based on

["Simple Strategies for Large Zero-Sum Games with Applications to Complexity Theory" by Lipton & Young, STOC'94]

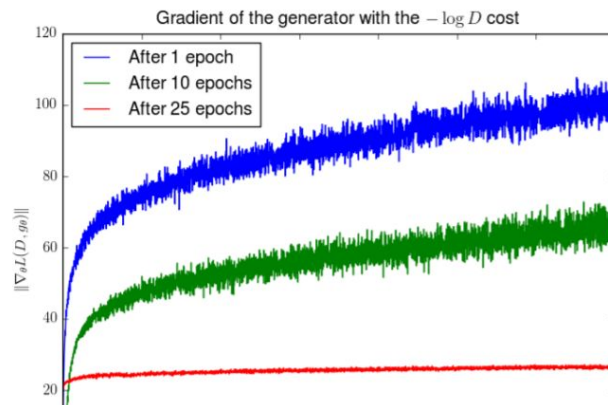
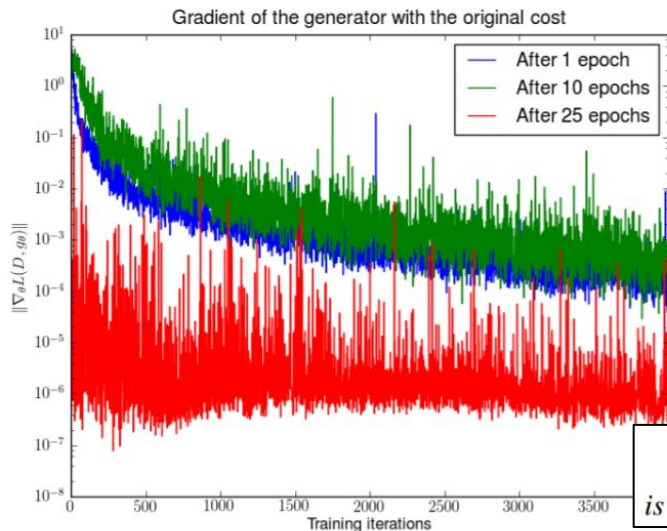
# Takeaway

1. The “sample GAN loss” reasonably quickly converges to its “true” value.
2. Jensen-Shannon and Wasserstein distances **do not**.
3. For large networks  $\epsilon$ -approximate equilibriums **exists**.

## Comments

1. No point in approximating JSD, Wasserstein (and MMD) precisely because their sample estimates are too far from actual values!
2. No point in using them for evaluation either (in higher dimensions)!
3. We are still optimizing for “exact” not “ $\epsilon$ -approximate” equilibriums.
4. Those equilibriums might also be very hard to get into!

# GAN-loss is pretty bad optimization-wise



$$\mathbb{E}_{z \sim p(z)} [-\nabla_{\theta} \log D(g_{\theta}(z))]$$

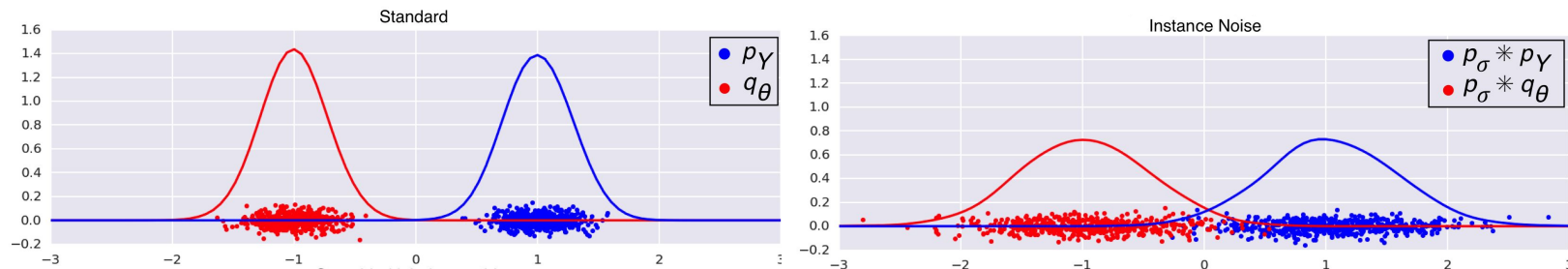
is a centered Cauchy distribution with *infinite expectation and variance*.<sup>4</sup>

(with generator fixed after X epochs)

But noise might help:

$$W(\mathbb{P}_r, \mathbb{P}_g) \leq 2V^{\frac{1}{2}} + 2C \sqrt{JSD(\mathbb{P}_{r+\epsilon} \| \mathbb{P}_{g+\epsilon})}$$

# Instance noise in the discriminator might help. Closed-form regularizer exist.



## Regularized Jensen-Shannon GAN

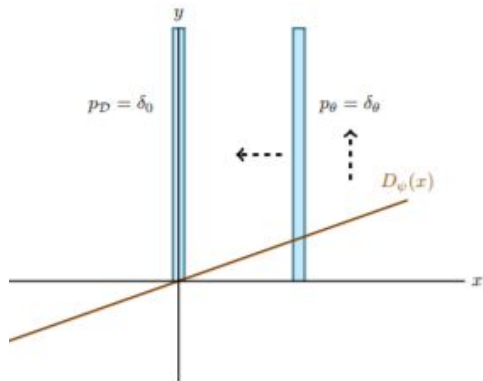
$$F_\gamma(\mathbb{P}, \mathbb{Q}; \varphi) = \mathbf{E}_{\mathbb{P}} [\ln(\varphi)] + \mathbf{E}_{\mathbb{Q}} [\ln(1 - \varphi)] - \frac{\gamma}{2} \Omega_{JS}(\mathbb{P}, \mathbb{Q}; \varphi)$$
$$\Omega_{JS}(\mathbb{P}, \mathbb{Q}; \varphi) := \mathbf{E}_{\mathbb{P}} [(1 - \varphi(\mathbf{x}))^2 \|\nabla \phi(\mathbf{x})\|^2] + \mathbf{E}_{\mathbb{Q}} [\varphi(\mathbf{x})^2 \|\nabla \phi(\mathbf{x})\|^2]$$

but requires figuring out a good annealing schedule

["Stabilizing Training of Generative Adversarial Networks through Regularization", Roth et al, NeurIPS'17]

["Instance Noise: A trick for stabilising GAN training", Ferenc Huszár, inference.vc]

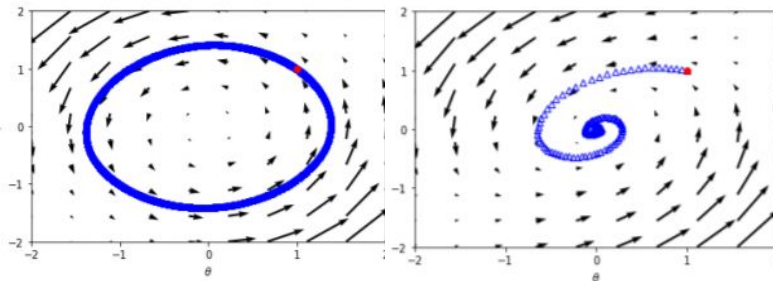
# A “toy GAN problem” confirms it.



$$D_\psi(x) = \psi \cdot x$$

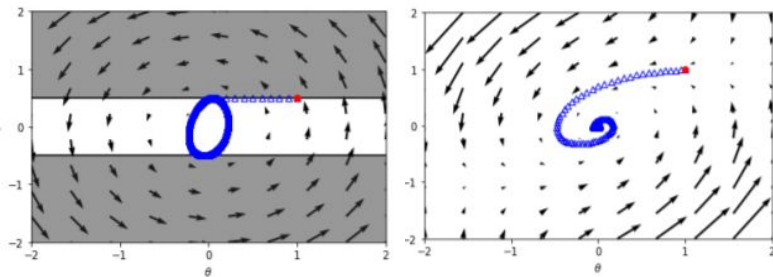
$$p_\theta = \delta_\theta$$

$$p_{\mathcal{D}} = \delta_0$$



(a) Standard GAN

(f) Instance noise



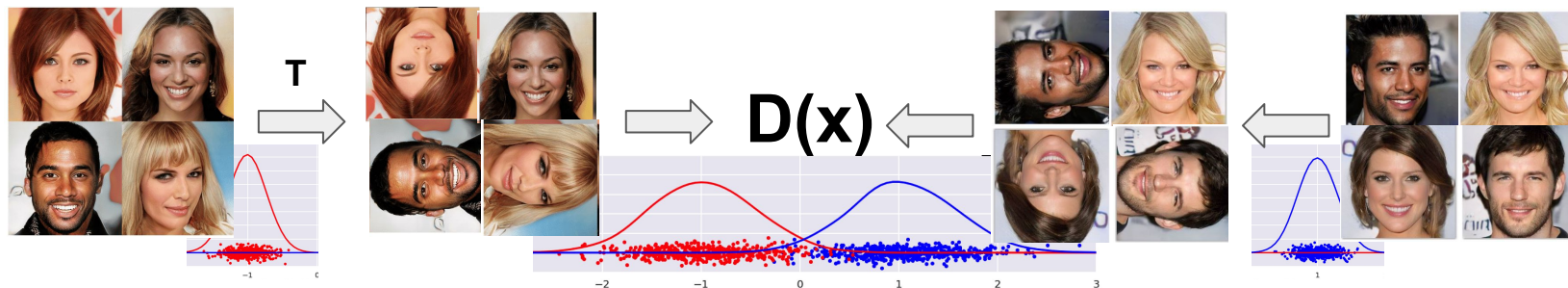
(c) WGAN ( $n_d = 5$ )

(g) Gradient penalty

# Let's extend to arbitrary augmentations.

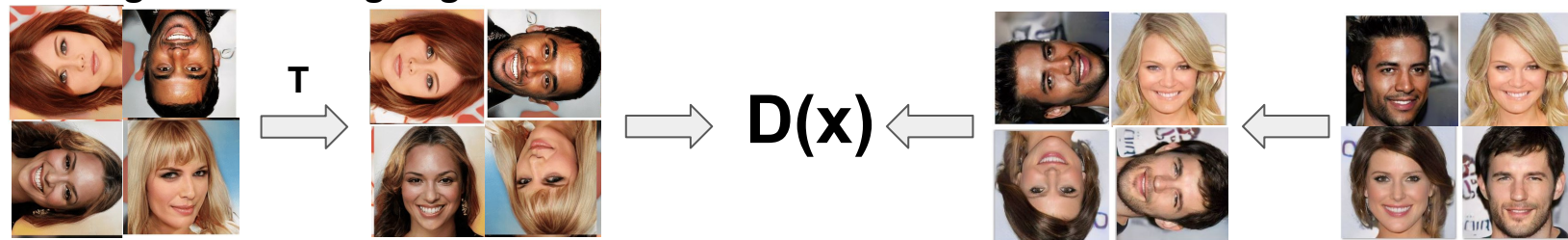
Assume augmentation  $T(x)$  randomly flips an image by  $[0, 90, 180, 270]$  and we apply  $T(x)$  “as instance noise” before passing them to  $D(x)$  to make images “less separable”.

“good” generated images



real images

generated images with wrong original orientation



Here is what you get - “leaking augmentation”.

$T(x)$  is flip

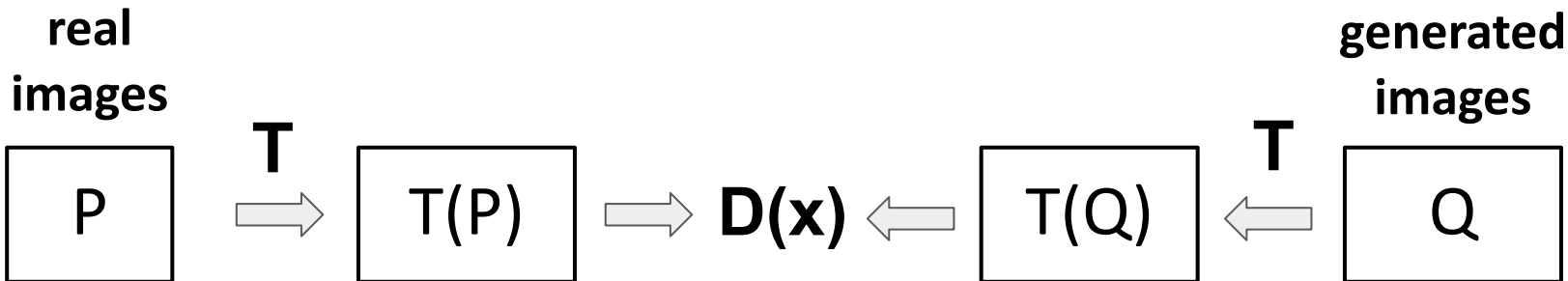


$T(x)$  is color shift





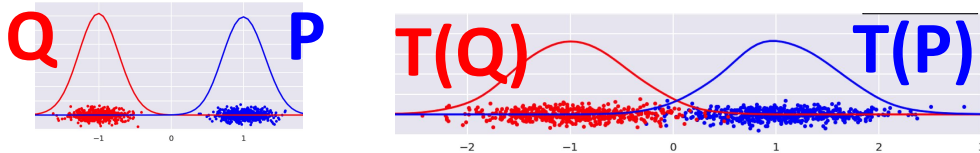
# How to avoid “leaking augmentation”?



We want  $T(x)$  such that  $T(P) = T(Q) \Leftrightarrow P = Q$ ,  
i.e. we want an *invertible* operator “ $T$ : distribution  $\square$  distribution”.

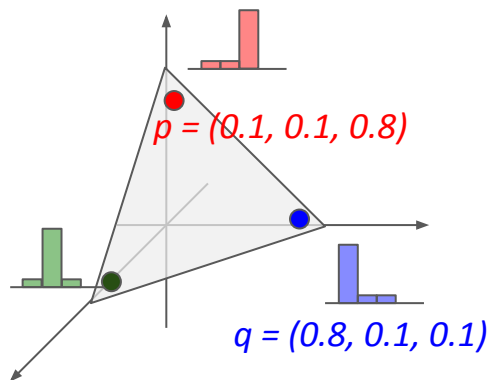
**Not** same as an invertible augmentation  $T(x)$ !

Example:  $T(P) = P * \text{Gaussian}(0, 1)$ , i.e.  $T(x) = x + \varepsilon$ ,  $\varepsilon \sim N(0, 1)$ .



# Markov operator

$X = 1\text{D}$  random variable,  $\text{supp}(X) = \{0, 1, 2\}$   
 $P(X) =$  a vector in  $\mathbb{R}^3$  that lies inside  $\Delta_3$



In this case, the  
 “T: distribution  $\rightarrow$  distribution”  
 is just a linear operator “T:  $\Delta_3 \rightarrow \Delta_3$ ”.

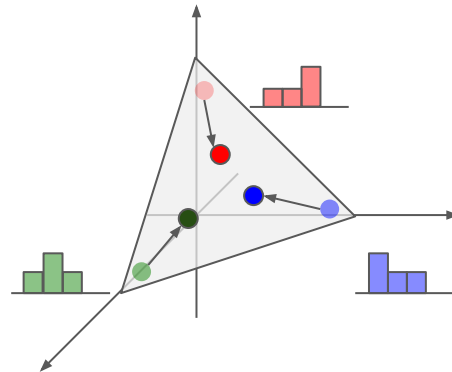
$$T = \begin{pmatrix} 0.6 & 0.2 & 0.2 \\ 0.2 & 0.6 & 0.2 \\ 0.2 & 0.2 & 0.6 \end{pmatrix}$$

**invertible!** sums  
over rows to 1

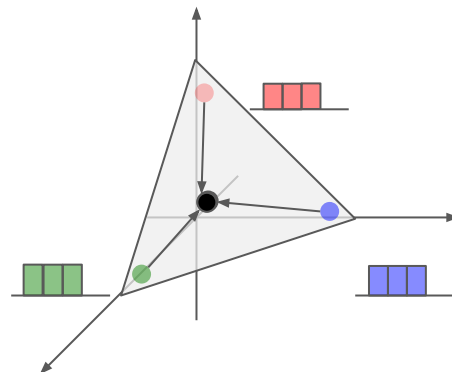
$$T = \begin{pmatrix} 1.0 & 1.0 & 1.0 \\ 1.0 & 1.0 & 1.0 \\ 1.0 & 1.0 & 1.0 \end{pmatrix} * \frac{1}{3}$$

**not invertible!** sums  
over rows to 1

in **observation space** the augmentation function is **random**  
 e.g.  $f(1) = \{=0 \text{ with } p=0.2, =1 \text{ with } p=0.6, \text{ and } =2 \text{ with } p=0.2\}$



deterministic linear  
operator in the  
distribution space



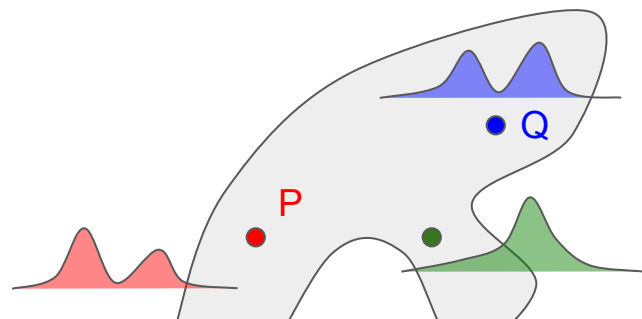
in **observation space**  
 $f(x)$  is either  $\{0, 1, 2\}$   
 with equal probabilities

not invertible!

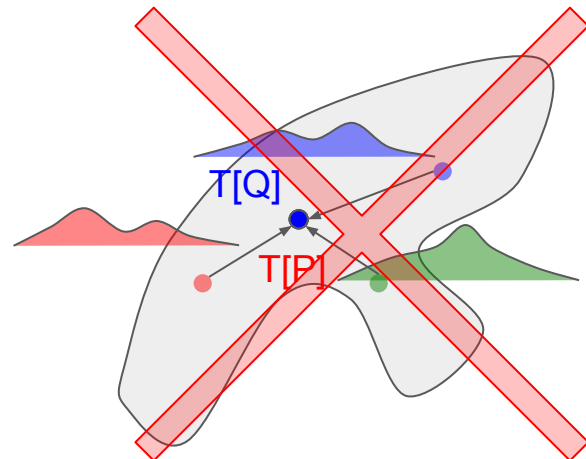
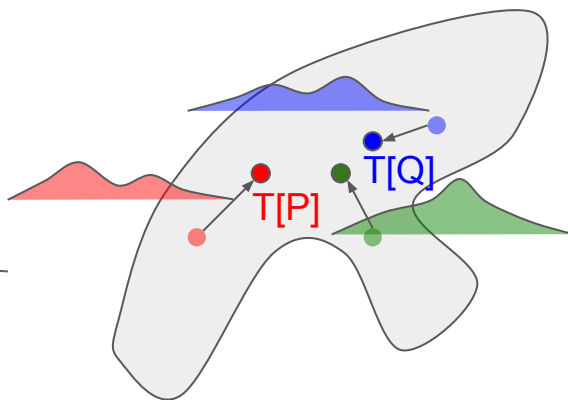
# Markov operator

We want  $T(x)$  such that  $T(P) = T(Q) \Leftrightarrow P = Q$ ,  
i.e. we want an *invertible* operator “ $T$ : distribution  $\square$  distribution”.

Example:  $T(P) = P * N(0, 1)$ , i.e.  $T(x) = x + \varepsilon$ ,  $\varepsilon \sim N(0, 1)$ ,  
i.e.  $T[P](x) = [P * N(0, 1)](x)$ ,  $T^{-1}(T(P)) = P$ ,  $T^{-1}(W) = \text{deconv}(W)$



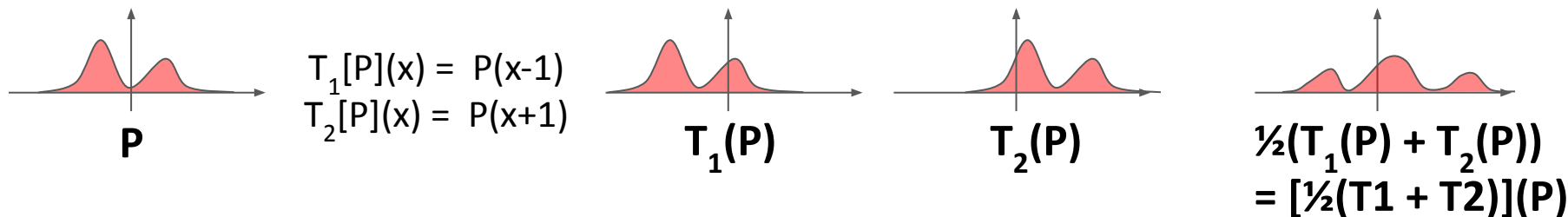
this is like an **infinite dimensional** vector  
and the operator  $T$  is also infinite dimensional



# How to test invertibility of an infinite-dimensional operator?

General statements:

1. A **composition** of invertible operators is **invertible** (i.e. a sequence of “good”/”non-leaking” augmentations is still good)
2. A **linear combination** of invertible operators is **not necessarily** invertible [ $\frac{1}{2}(T_1 + T_2)$ ](P) means randomly choosing between augmentations  $T_1$  and  $T_2$  and applying it to a single sample from P)



# How to test invertibility of an infinite-dimensional operator?

$$\mathcal{T} = \sum_{i=0}^{N-1} p_i \mathcal{G}^i$$

$$\mathcal{U} = \sum_{j=0}^{N-1} q_j \mathcal{G}^j$$

$$\mathcal{UT} = \left( \sum_{i=0}^{N-1} p_i \mathcal{G}^i \right) \left( \sum_{j=0}^{N-1} q_j \mathcal{G}^j \right) = \sum_{i,j=0}^{N-1} p_i q_j \mathcal{G}^{i+j}$$

$$= \sum_{k=0}^{N-1} [p \otimes q]_k \mathcal{G}^k = \sum_{k=0}^{N-1} [\mathbf{F}^{-1}(\hat{p} \odot \hat{q})]_k \mathcal{G}^k$$

if we set  
 $\hat{q}_i = \frac{1}{\hat{p}_i}$

$$= \sum_{k=0}^{N-1} [\mathbf{F}^{-1}(\hat{p} \odot \hat{p}^{-1})]_k \mathcal{G}^k = \sum_{k=0}^{N-1} [\mathbf{F}^{-1} \mathbf{1}]_k \mathcal{G}^k = \mathcal{G}^0 = \mathcal{I}$$

cyclic group  $\{\mathcal{G}^i\}_{i=0}^{N-1}$   
 $\mathcal{G}^0$  is the identity mapping  
 e.g.  $\{0, 90, 180, 270\}$  img rots  
 = not just “any” stochastic matrix

**Inverse operator exists if there are no zeros in operator’s Fourier spectre.**

**Solution:** uniform but with higher probability of  $\mathcal{G}^0$  like  $[0.28 \ 0.24 \ 0.24 \ 0.24]$  - has no zeros in spectrum!

Essentially  $T = [(1-\alpha) * \text{Uniform} + \alpha * \text{Identity}]$  - e.g. almost like a regularization.

# How to test invertibility of an infinite-dimensional operator?

*Inverse operator exists if there are no zeros in operator's Fourier spectre.*

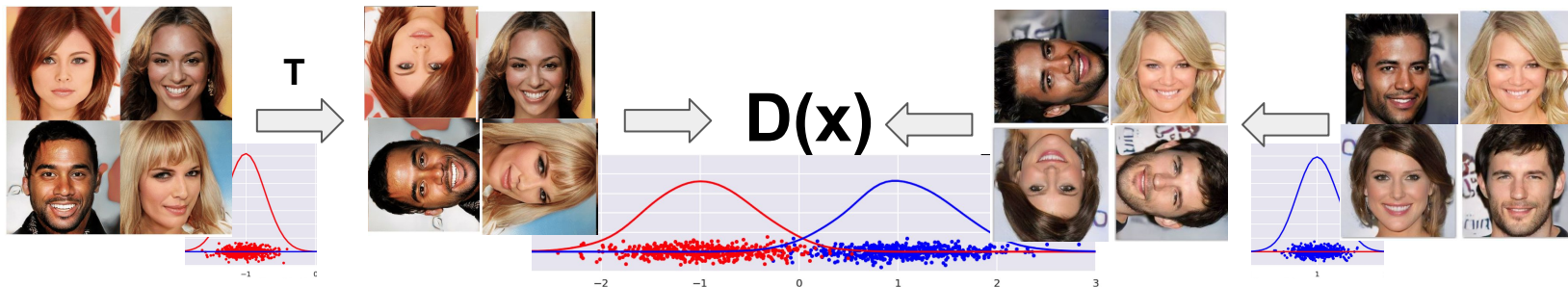
**Solution:** uniform **but** with higher probability of  $G^0$  like  $[0.28 \ 0.24 \ 0.24 \ 0.24]$  - has no zeros in spectrum!

Essentially  $T = [(1-\alpha) * \text{Uniform} + \alpha * \text{Identity}]$  - e.g. almost like a regularization.

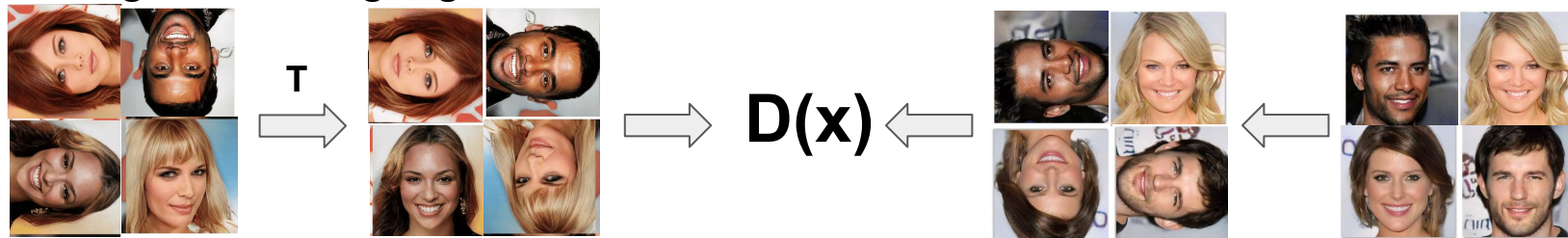
“good” generated images

both have slightly more “upright” faces!

real images



generated images with wrong orientation



left does not have slightly more “upright” faces!

# How to test invertibility of an infinite-dimensional operator?

Other cases:

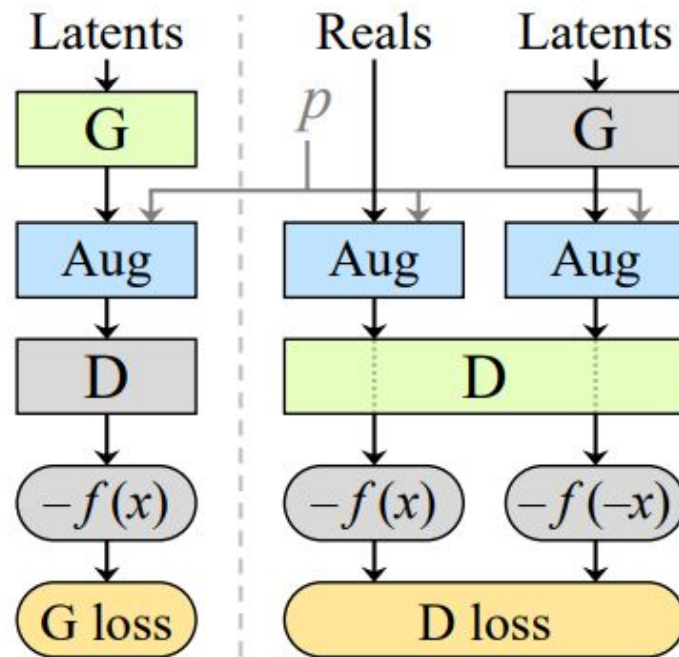
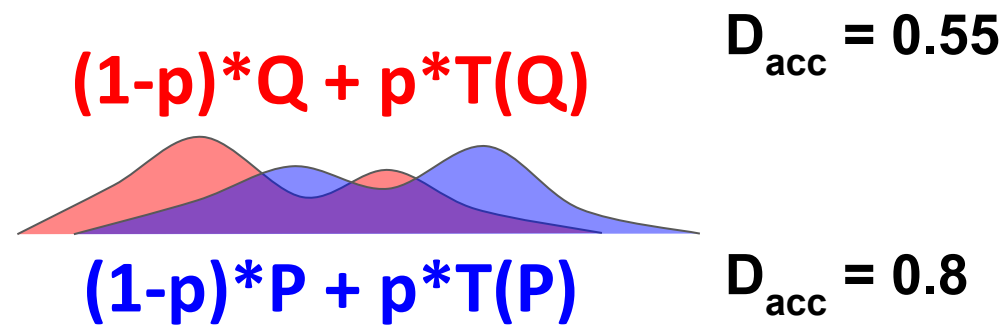
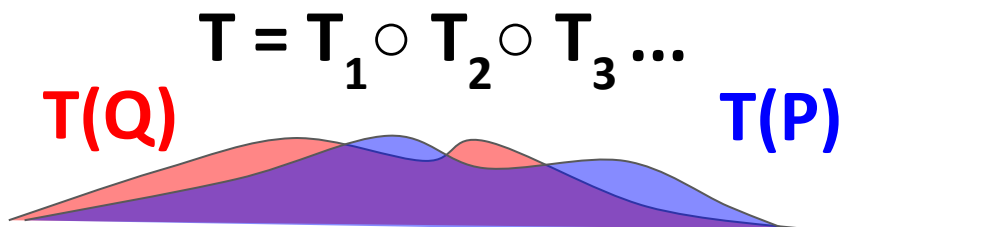
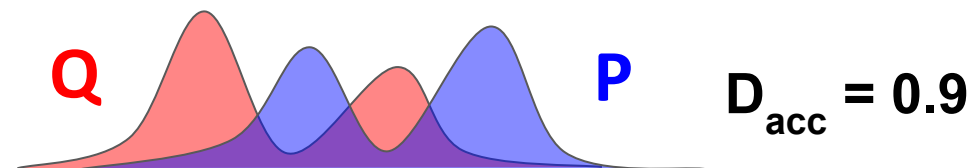
1. Non-compact discrete groups (integer shift): also “non-zero Fourier”
2. For continuous groups (e.g. rotations): also “non-zero Fourier” (use Haar measure over that groups under the integral);
3. Additive pixel noise: “non-zero Fourier” of the noise kernel
4. Cropping / blitting / “projection”: requires  $P(\text{identity}) > 0$

Assume  $\exists \mathbf{y} \neq \mathbf{z}$  s.t.  $T\mathbf{y} = T\mathbf{z}$ , e.g.  $T(\mathbf{y}-\mathbf{z}) = \mathbf{0}$ , e.g.  $T\mathbf{x} = \mathbf{0}$ .

$$T = p_0 \mathcal{I} + \sum_{j=1}^N p_j \mathcal{P}_j \quad \Bigg| \quad 0 = T\mathbf{x} = p_0 \mathbf{x} + \sum_{j=1}^N p_j \mathcal{P}_j \mathbf{x} \quad \Bigg| \quad \sum_{j=1}^N p_j \langle \mathbf{x}, \mathcal{P}_j \mathbf{x} \rangle = -p_0 \langle \mathbf{x}, \mathbf{x} \rangle$$

$\geq 0$   $\geq 0$   
*invertible if  $p_0 \neq 0$*

# So how do we use it?

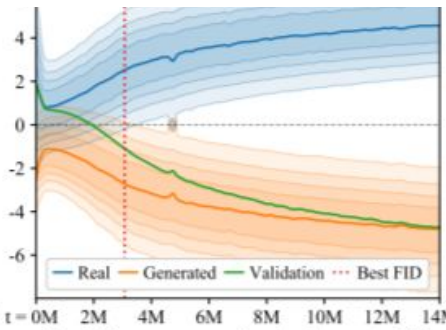
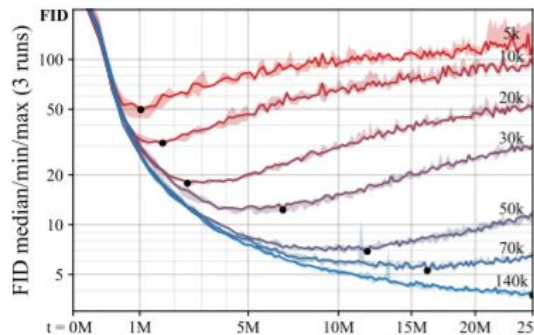




# Does it help? - yes!

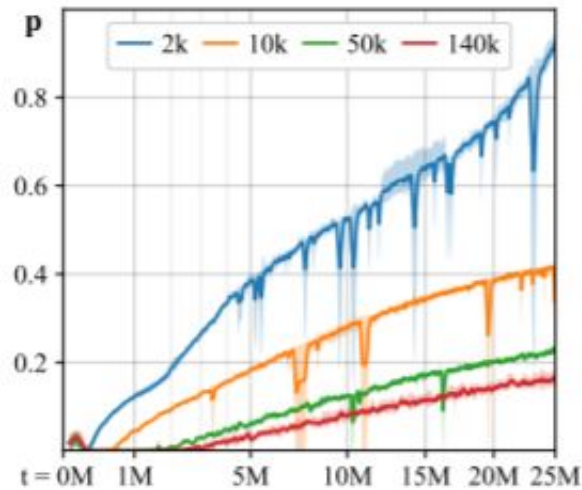
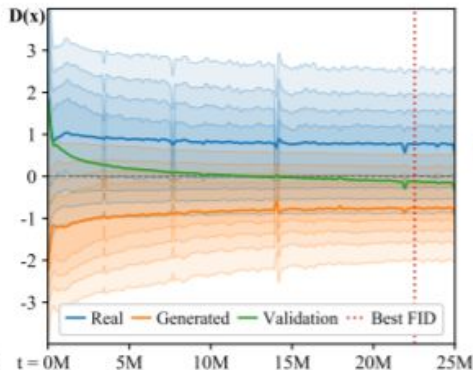
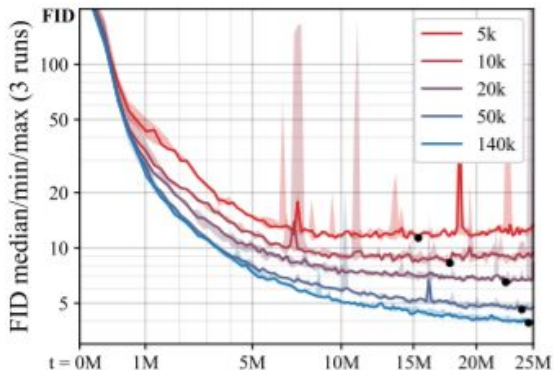
real train / real val / generated  
 D(x) scores trained on 20k samples

w/o augmentation



$$r_v = \frac{\mathbb{E}[D_{\text{train}}] - \mathbb{E}[D_{\text{validation}}]}{\mathbb{E}[D_{\text{train}}] - \mathbb{E}[D_{\text{generated}}]}$$

with aug.

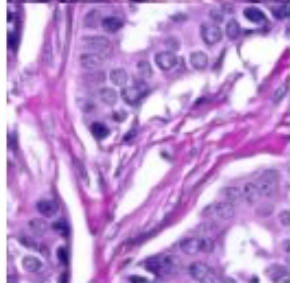
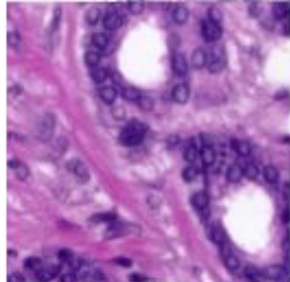
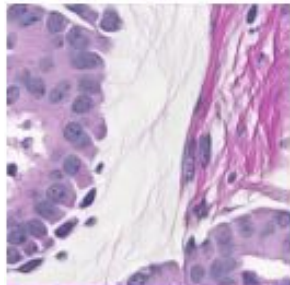


(c) Evolution of  $p$  over training

FID(train step) for  
 different amount of real data

Does it help? - yes!

BRECAHAD  
1944 img,  $512^2$



AFHQ CAT, DOG, WILD ( $512^2$ )  
5153 img      4739 img      4738 img



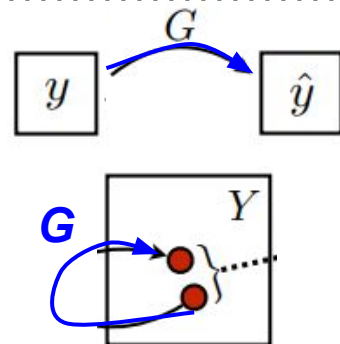
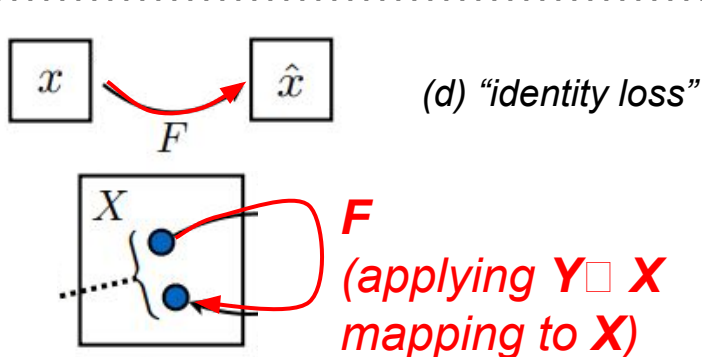
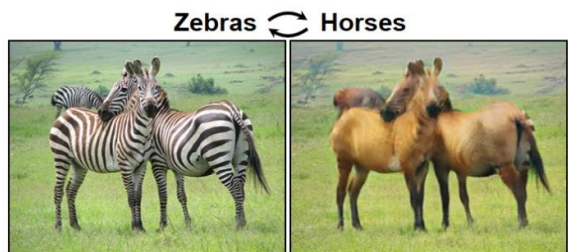
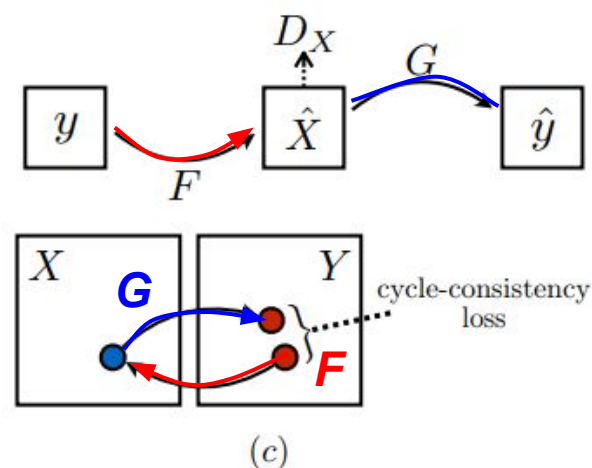
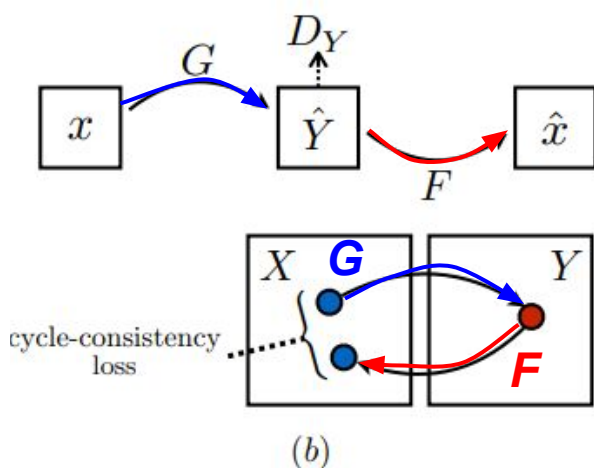
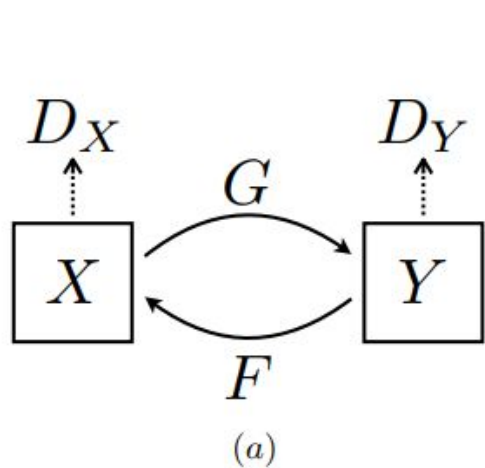
# Takeaway

1. Regularizing the discriminator with augmentations helps.
2. But it has to be done in a way that does not “leak” into generated images.
3. For a wide variety of transformations, applying them with a fixed probability “does not leak” into generated examples.

## Comments

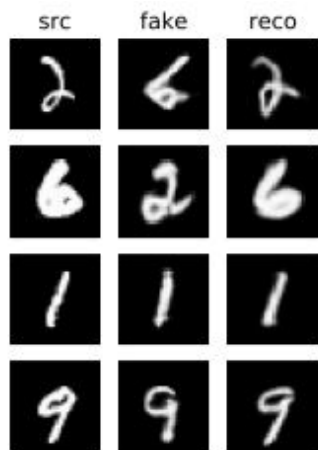
1. All these methods still require careful parameter annealing.
2. As a result we can not reason about the convergence of an objective because there is no single objective! (we change it as we train)

# CycleGAN overview



# How to reason about the complexity of the CycleGAN?

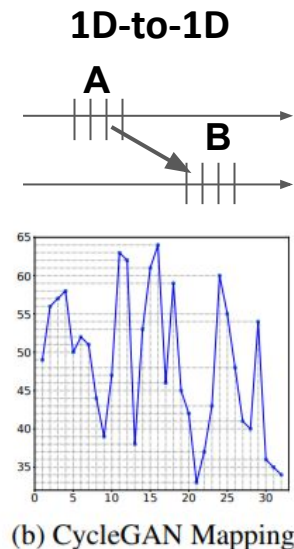
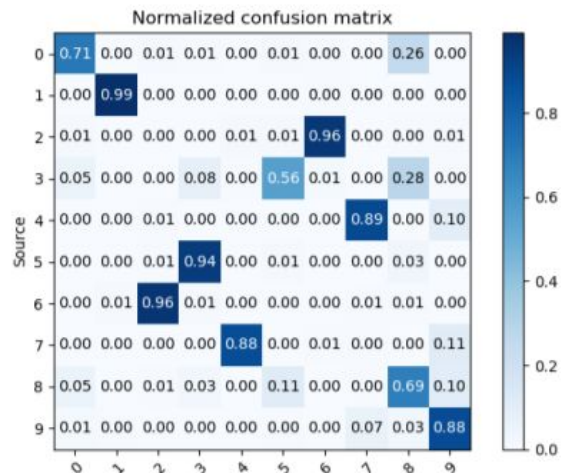
CycleGAN trained to map MNIST train split to the MNIST test split.



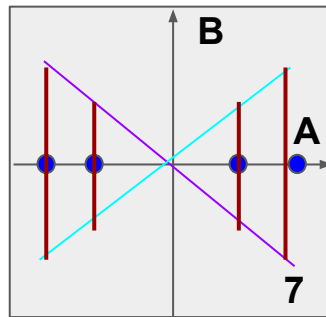
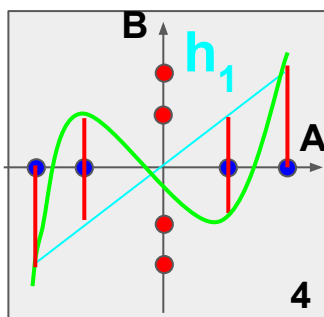
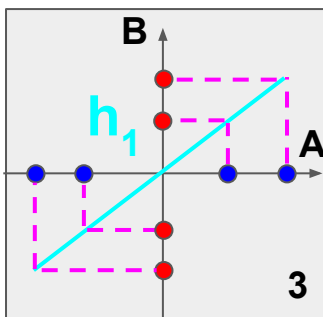
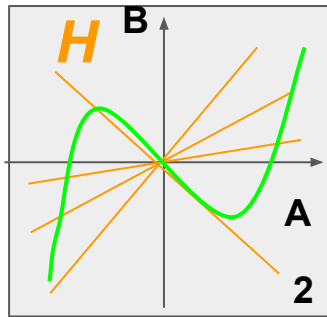
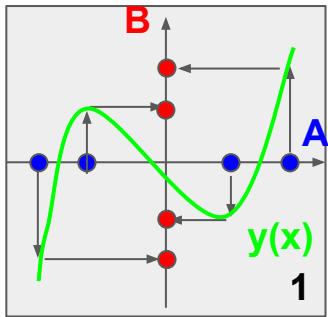
X



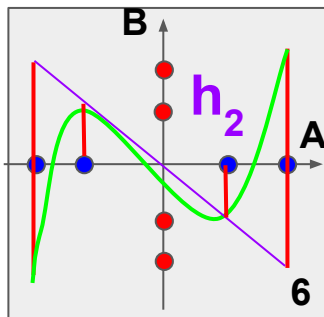
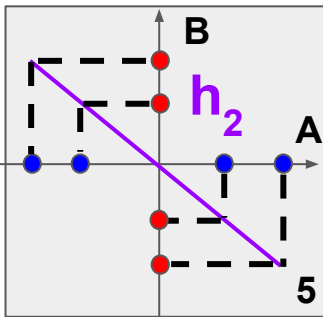
Y



# How to bound the unsupervised alignment error?



**prediction error** <  
smallest unsupervised alignment error  
+ smallest approximation error in H  
+ the variance between functions  
minimizing the alignment loss.



[“Risk Bounds for Unsupervised Cross-Domain Mapping with IPMs”, Galanti et al., JMLR’21]

[“Estimating the Success of Unsupervised Image to Image Translation”, Benaim et al., ECCV’18]

[“The role of minimal complexity functions in unsupervised learning of semantic mappings”, Galanti et al., ICLR’18]

# How to bound the unsupervised alignment error?

$$R_{D_A}[h_1, y] \lesssim \underbrace{\sup_{h_2 \in \mathcal{P}_\omega} R_{S_A}[h_1, h_2]}_{\text{red line}} + c \underbrace{\inf_{h \in \mathcal{P}_\omega} \rho_C(h \circ S_A, S_B)}_{\text{blue line}} + \dots$$

$R$  - pred error;  $\rho$  - alignment error;  $D_A$  - distribution;  $S_A$  - dataset;  $\mathcal{P}_k$  - hypotheses with “low” alignment error

## How to use this bound?

A given choice of hyperparameters is evaluated as follows:

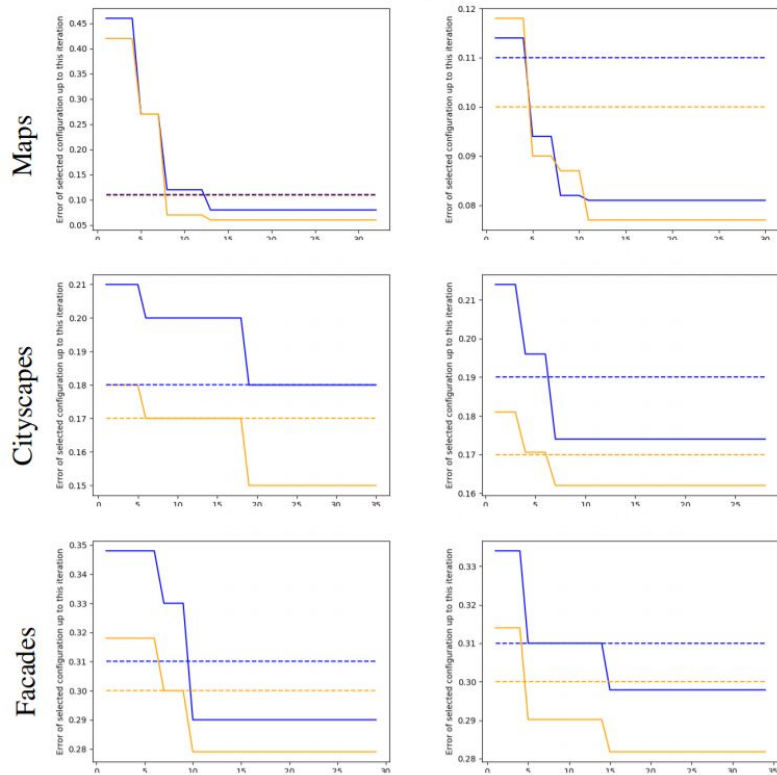
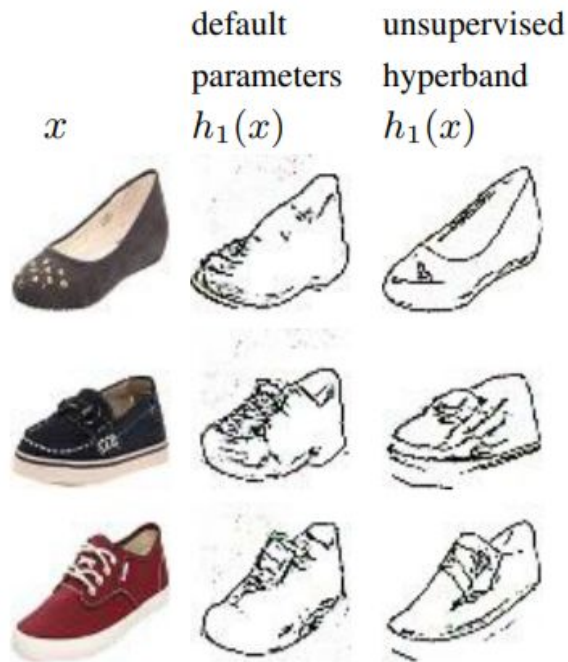
1.  $\inf_{h \in \mathcal{P}_k} \rho_C(h \circ S_A, S_B)$  - we minimize the “GAN loss” to get “the first best”  $h_1$
2.  $\min_{h_2 \in \mathcal{H}_k} \left\{ \rho_C(h_2 \circ S_A, S_B) - \lambda R_{S_A}[h_1, h_2] \right\}$  - then pick the “second best”  $h_2$
3.  $R_{S_A}[h_1, h_2] + \rho_C(h_1 \circ S_A, S_B)$  - and then bound the unknown GT error

[“Risk Bounds for Unsupervised Cross-Domain Mapping with IPMs”, Galanti et al., JMLR’21]

[“Estimating the Success of Unsupervised Image to Image Translation”, Benaim et al., ECCV’18]

[“The role of minimal complexity functions in unsupervised learning of semantic mappings”, Galanti et al., ICLR’18]

# Does it help? - yes!



The **ground truth error** and the **theoretical bound** as a function of hyper-parameter optimization steps varying

- encoder and decoder #layers
- batch size
- learning rate
- ...

[“Risk Bounds for Unsupervised Cross-Domain Mapping with IPMs”, Galanti et al., JMLR’21]

[“Estimating the Success of Unsupervised Image to Image Translation”, Benaim et al., ECCV’18 ]

[“The role of minimal complexity functions in unsupervised learning of semantic mappings”, Galanti et al., ICLR’18]



translation  
functions

discriminator  
class

domains

**Theorem 1 (Cross-Domain Mapping with IPMs)** Assume that  $\mathcal{X}_A \subset \mathbb{R}^N$  and  $\mathcal{X}_B \subset \mathbb{R}^M$  are convex and bounded sets. Let  $\mathcal{H}$  be the hypothesis class and  $\mathcal{C}$  the class of discriminators. Assume that  $\mathcal{C} \subset C^2$  and  $\sup_{d \in \mathcal{C}} \|d\|_{\infty, \mathcal{X}_A \cup \mathcal{X}_B} < \infty$ . Then, for any  $\delta \in (0, 1)$  and  $c \geq 1$ , with probability at least  $1 - \delta$  over the selection of  $\mathcal{S}_A \sim D_A^{m_1}$  and  $\mathcal{S}_B \sim D_B^{m_2}$ , for every  $\omega \in \Omega$  and  $h_1 \in \mathcal{P}_\omega := \mathcal{P}_\omega(\mathcal{S}_A, \mathcal{S}_B)$ , we have:

single  
discriminator

finite datasets

hyperparams

a mapping  
produced by  
the learning  
algorithm

[“Risk Bounds for Unsupervised Cross-Domain Mapping with IPMs”, Galanti et al., JMLR’21]

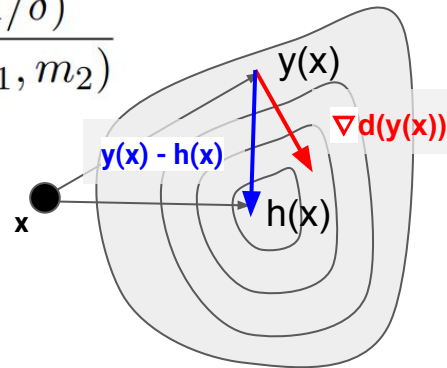
[“Estimating the Success of Unsupervised Image to Image Translation”, Benaim et al., ECCV’18]

[“The role of minimal complexity functions in unsupervised learning of semantic mappings”, Galanti et al., ICLR’18]

**Theorem 1 (Cross-Domain Mapping with IPMs)** Assume that  $X_A \subset \mathbb{R}^N$  and  $X_B \subset \mathbb{R}^M$  are convex and bounded sets. Let  $\mathcal{H}$  be the hypothesis class and  $\mathcal{C}$  the class of discriminators. Assume that  $\mathcal{C} \subset C^2$  and  $\sup_{d \in \mathcal{C}} \|d\|_{\infty, X_A \cup X_B} < \infty$ . Then, for any  $\delta \in (0, 1)$  and  $c \geq 1$ , with probability at least  $1 - \delta$  over the selection of  $S_A \sim D_A^{m_1}$  and  $S_B \sim D_B^{m_2}$ , for every  $\omega \in \Omega$  and  $h_1 \in \mathcal{P}_\omega := \mathcal{P}_\omega(S_A, S_B)$ , we have:

$$R_{D_A}[h_1, y] \lesssim \sup_{h_2 \in \mathcal{P}_\omega} R_{S_A}[h_1, h_2] + c \inf_{h \in \mathcal{P}_\omega} \rho_{\mathcal{C}}(h \circ S_A, S_B) + \inf_{h \in \mathcal{P}_\omega} \inf_{\substack{d \in \mathcal{C} \\ \beta(d) \leq 1}} \mathcal{K}(h, d; y) \\ + \hat{\mathcal{R}}_{S_A}(\ell_{\mathcal{H}}) + \hat{\mathcal{R}}_{S_A}(\mathcal{C} \circ \mathcal{H}) + \hat{\mathcal{R}}_{S_B}(\mathcal{C}) + \sqrt{\frac{\log(1/\delta)}{\min(m_1, m_2)}} \quad (9)$$

where,  $\mathcal{K}(h, d; y) := \mathbb{E}_{x \sim D_A} [\|\nabla_{y(x)} d(y(x)) - (h(x) - y(x))\|_2]$ .



["Risk Bounds for Unsupervised Cross-Domain Mapping with IPMs", Galanti et al., JMLR'21]

["Estimating the Success of Unsupervised Image to Image Translation", Benaim et al., ECCV'18]

["The role of minimal complexity functions in unsupervised learning of semantic mappings", Galanti et al., ICLR'18]

$$R_{D_A}[h_1, y] \lesssim \sup_{h_2 \in \mathcal{P}_\omega} R_{S_A}[h_1, h_2] + c \inf_{h \in \mathcal{P}_\omega} \rho_C(h \circ S_A, S_B) + \inf_{h \in \mathcal{P}_\omega} \inf_{\substack{d \in \mathcal{C} \\ \beta(d) \leq 1}} \mathcal{K}(h, d; y) \\ + \hat{\mathcal{R}}_{S_A}(\ell_{\mathcal{H}}) + \hat{\mathcal{R}}_{S_A}(C \circ \mathcal{H}) + \hat{\mathcal{R}}_{S_B}(C) + \sqrt{\frac{\log(1/\delta)}{\min(m_1, m_2)}}$$

**Lem 3:**

$$R_{D_A}[h_1, y] \leq 3 \sup_{h_2 \in \mathcal{P}} R_{D_A}[h_1, h_2] + 3 \inf_{h \in \mathcal{P}} R_{D_A}[h, y]$$

**Lem 4:**

$$R_{D_A}[h, y] \leq \frac{2\rho_C(h \circ D_A, D_B)}{2 - \beta(d)} + \frac{2 \sup_{u \in X_A} \|h(u) - y(u)\|_2 \leq L}{2 - \beta(d)} \cdot \mathcal{K}(h, d; y)$$

**Lem 7:**

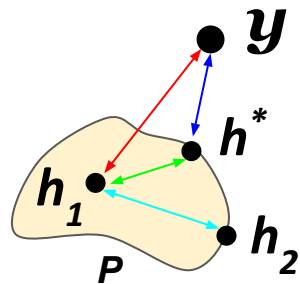
$$R_{D_A}[h_1, h_2] \leq R_{S_A}[h_1, h_2] + 2\hat{\mathcal{R}}_{S_A}(\ell_{\mathcal{H}}) + 9K^2 \sqrt{\frac{\log(6/\delta)}{2m_1}}$$

**Lem 3** (Triangle inequality and the “set diameter”)

$$R_{D_A}[h_1, y] \leq 3 \sup_{h_2 \in \mathcal{P}} R_{D_A}[h_1, h_2] + 3 \inf_{h \in \mathcal{P}} R_{D_A}[h, y]$$

$$h^* \in \arg \inf_{h \in \mathcal{P}} R_{D_A}[h, y]$$

$$\begin{aligned} \|a - c\|_2^2 &\leq (\|a - b\|_2 + \|b - c\|_2)^2 \\ &\leq \|a - b\|_2^2 + \|b - c\|_2^2 + 2 \max(\|a - b\|_2^2, \|b - c\|_2^2) \\ &\leq 3(\|a - b\|_2^2 + \|b - c\|_2^2) \end{aligned}$$



$$\begin{aligned} R_{D_A}[h_1, y] &= \mathbb{E}_{x \sim D_A} [\|h_1(x) - y(x)\|_2^2] \\ &\leq \mathbb{E}_{x \sim D_A} [3\|h_1(x) - h^*(x)\|_2^2 + 3\|h^*(x) - y(x)\|_2^2] \\ &= 3 \left[ R_{D_A}[h_1, h^*] + \inf_{h \in \mathcal{P}} R_{D_A}[h, y] \right] \end{aligned}$$

$$R_{D_A}[h_1, h^*] \leq \sup_{h_2 \in \mathcal{P}} R_{D_A}[h_1, h_2]$$

**Lem 4** (Prediction Error via Stat. Distance and Discriminator Capacity)

$$R_{D_A}[h, y] \leq \frac{\text{avg prediction error}}{2 - \beta(d)} + \frac{\text{stat. distance}}{2 - \beta(d)} \cdot \frac{\text{max prediction error} < L}{2 - \beta(d)} \cdot \text{discriminator error}$$

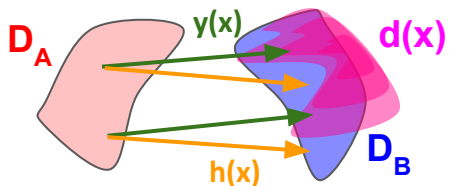
$$\text{where, } \mathcal{K}(h, d; y) := \mathbb{E}_{x \sim D_A} [\|\nabla_{y(x)} d(y(x)) - (h(x) - y(x))\|_2]$$

$$\rho_C(h \circ D_A, D_B) = \sup_{d \in \mathcal{C}} \left\{ \mathbb{E}_{u \sim h \circ D_A} [d(u)] - \mathbb{E}_{v \sim D_B} [d(v)] \right\} = \sup_{d \in \mathcal{C}} \left\{ \mathbb{E}_{x \sim D_A} [d \circ h(x) - d \circ y(x)] \right\}$$

stat. distance between  $h(P_A)$  and  $P_B$

$$\geq \mathbb{E}_{x \sim D_A} [d(h(x)) - d(y(x))] = \mathbb{E}_{x \sim D_A} [\|h(x) - y(x)\|_2^2]$$

$$\begin{aligned} f(x) - f(y) &= \langle \nabla f(x), x - y \rangle + \dots \\ &= \|x - y\|^2 + \langle \nabla f(x) - (x - y), x - y \rangle + \dots \end{aligned}$$



$$\begin{aligned} &+ \mathbb{E}_{x \sim D_A} [\langle \nabla_{y(x)} d(y(x)) - (h(x) - y(x)), h(x) - y(x) \rangle] \\ &+ \frac{1}{2} \mathbb{E}_{x \sim D_A} [\langle (h(x) - y(x))^\top \cdot H_d(u_{d,x}^*), h(x) - y(x) \rangle] \end{aligned}$$

$$\begin{aligned} &= \mathbb{E}_{x \sim D_A} [\|h(x) - y(x)\|_2^2] - \frac{1}{2} \mathbb{E}_{x \sim D_A} [\beta(d) \cdot \|h(x) - y(x)\|_2^2] \\ &\quad - \sup_{u \in \mathcal{X}_A} \|h(u) - y(u)\|_2 \cdot \mathbb{E}_{x \sim D_A} [\|\nabla_{y(x)} d(y(x)) - (h(x) - y(x))\|_2] \end{aligned}$$

discriminator approx error (capacity)

$$= \left(1 - \frac{\beta(d)}{2}\right) R_{D_A}[h, y] - \sup_{u \in \mathcal{X}_A} \|h(u) - y(u)\|_2 \cdot \mathcal{K}(h, d; y)$$

prediction error

max pred error < L

**Definition 3.1 (Empirical Rademacher complexity)**

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^m \sigma_i g(z_i) \right]$$

**Theorem 3.3** Let  $\mathcal{G}$  be a family of functions mapping from  $\mathcal{Z}$  to  $[0, 1]$ , with probability at least  $1 - \delta$

$$\mathbb{E}[g(z)] \leq \frac{1}{m} \sum_{i=1}^m g(z_i) + 2\widehat{\mathfrak{R}}_S(\mathcal{G}) + 3\sqrt{\frac{\log \frac{2}{\delta}}{2m}}$$

$$\mathbb{E}_S[\Phi(S)] = \mathbb{E}_S \left[ \sup_{g \in \mathcal{G}} (\mathbb{E}[g] - \widehat{\mathbb{E}}_S(g)) \right] = \mathbb{E}_{S, S'} \left[ \sup_{a \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^m (g(z'_i) - g(z_i)) \right]$$

$$= \mathbb{E}_{\sigma, S, S'} \left[ \sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^m \sigma_i (g(z'_i) - g(z_i)) \right]$$

**sup g** - a function with the **largest** deviation of its value from average across possible splits  $S, S'$

$$\leq \mathbb{E}_{\sigma, S'} \left[ \sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^m \sigma_i g(z'_i) \right] + \mathbb{E}_{\sigma, S} \left[ \sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^m -\sigma_i g(z_i) \right]$$

$$= 2 \mathbb{E}_{\sigma, S} \left[ \sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^m \sigma_i g(z_i) \right] = 2\mathfrak{R}_m(\mathcal{G}).$$

with probability at least  $1 - \delta$   
by McDiarmid's inequality

$$\Phi(S) \leq \mathbb{E}_S[\Phi(S)] + \sqrt{\frac{\log \frac{2}{\delta}}{2m}} \quad \text{and} \quad \mathfrak{R}_m(\mathcal{G}) \leq \widehat{\mathfrak{R}}_S(\mathcal{G}) + \sqrt{\frac{\log \frac{2}{\delta}}{2m}}$$

**Lem 7** (Sample complexity)

$$\mathbb{E}[g(z)] \leq \frac{1}{m} \sum_{i=1}^m g(z_i) + 2\hat{\mathfrak{R}}_S(\mathcal{G}) + 3\sqrt{\frac{\log \frac{2}{\delta}}{2m}}$$

$$R_{D_A}[h_1, h_2] \leq R_{S_A}[h_1, h_2] + 2\hat{\mathfrak{R}}_{S_A}(\ell_{\mathcal{H}}) + 9K^2 \sqrt{\frac{\log(6/\delta)}{2m_1}}$$

---

$$\rho_C(D_B, S_B) = \sup_{d \in C} \left\{ \mathbb{E}_{x \sim D_B}[d(x)] - \frac{1}{m_2} \sum_{x \in S_B} d(x) \right\} \lesssim \hat{\mathfrak{R}}_{S_B}(C) + \sqrt{\frac{\log(1/\delta)}{m_2}}$$

---

$$\rho_C(h \circ D_A, h \circ S_A) \lesssim \hat{\mathfrak{R}}_{S_A}(C \circ \mathcal{H}) + \sqrt{\frac{\log(1/\delta)}{m_1}}$$

$$R_{D_A}[h_1, y] \lesssim \sup_{h_2 \in \mathcal{P}_\omega} R_{S_A}[h_1, h_2] + c \inf_{h \in \mathcal{P}_\omega} \rho_C(h \circ S_A, S_B) + \inf_{h \in \mathcal{P}_\omega} \inf_{\substack{d \in \mathcal{C} \\ \beta(d) \leq 1}} \mathcal{K}(h, d; y) \\ + \hat{\mathcal{R}}_{S_A}(\ell_{\mathcal{H}}) + \hat{\mathcal{R}}_{S_A}(C \circ \mathcal{H}) + \hat{\mathcal{R}}_{S_B}(C) + \sqrt{\frac{\log(1/\delta)}{\min(m_1, m_2)}}$$

**Lem 3:**

$$R_{D_A}[h_1, y] \leq 3 \sup_{h_2 \in \mathcal{P}} R_{D_A}[h_1, h_2] + 3 \inf_{h \in \mathcal{P}} R_{D_A}[h, y]$$

**Lem 4:**

$$R_{D_A}[h, y] \leq \frac{2\rho_C(h \circ D_A, D_B)}{2 - \beta(d)} + \frac{2 \sup_{u \in X_A} \|h(u) - y(u)\|_2 \leq L}{2 - \beta(d)} \cdot \mathcal{K}(h, d; y)$$

**Lem 7:**

$$R_{D_A}[h_1, h_2] \leq R_{S_A}[h_1, h_2] + 2\hat{\mathcal{R}}_{S_A}(\ell_{\mathcal{H}}) + 9K^2 \sqrt{\frac{\log(6/\delta)}{2m_1}}$$



# Takeaway

The **prediction error** of the **unsupervised** image translation method (wrt the ground truth output) can be bounded via

- **minimal statistical distance** attainable by the network and
- **variance between solutions** that attain that lowest statistical distance.

And this bound actually works in practice!

## Comments

1. Regression CNNs can fit almost random  $(x,y)$  pairs - can I2I networks fit random  $(x_1, x_2)$  pairs? if so why doesn't  $R[h_1, h_2]$  explode?
2. The empirical **Rademacher complexity** of the discriminator class seems related to the “**expected statistical distance between random splits** of that dataset”?

# Recap

1. “Neural GAN distances” between datasets seem to have have **better sample complexity** than “classical” distances. We used an  $\epsilon$ -net over NN weights and Chernoff bound on each element of  $\epsilon$ -net to show that.
2. These neural distances can be “smoothened” via **instance noise and augmentations** to make gradient descent iterations more stable. By treating random augmentations as Markov operators we showed that in most cases **skipping augmentations with fixed probability** ensures that the neural distance remain “non-leaking” even under augmentations.
3. The **prediction error** of the unsupervised alignment method can be bounded via the **variance between solutions** attaining similar GAN loss.

# Thank you for your time!

Main papers:

1. **“Generalization and Equilibrium in Generative Adversarial Nets”** by Sanjeev Arora, Rong Ge, Yingyu Liang, Tengyu Ma, Yi Zhang, Proceedings International Conference on Machine Learning (PMLR) 2017.
2. **“Training Generative Adversarial Networks with Limited Data”** by Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, Timo Aila; Advances in Neural Information Processing Systems (NeurIPS) 2020.
3. **“Risk Bounds for Unsupervised Cross-Domain Mapping with IPMs”** by Tomer Galanti, Sagie Benaim, Lior Wolf; JMLR 2021. // *paper series (2017-2021)*