

Fast Fréchet Inception Distance

Frederik Hvilshøj

May 5, 2021

Data-Intensive Systems Group, Aarhus University



Paper URL



Slides URL

Today (45 minutes):

1. Comparing Generative Models

2. Metrics

3. Fast Fréchet Inception Distance

4. What is it Good For? 🎵

Comparing Generative Models

Which one is better?

“Real”

$$X_1, \dots, X_n \sim P_{\text{real}}$$



Which one is better?

“Real”

$$x_1, \dots, x_n \sim P_{\text{real}}$$



“Fake 1”

$$x_1^{(1)}, \dots, x_n^{(1)} \sim P_{G_1}$$



Which one is better?

“Real”

$$x_1, \dots, x_n \sim P_{\text{real}}$$



“Fake 1”

$$x_1^{(1)}, \dots, x_n^{(1)} \sim P_{G_1}$$



“Fake 2”

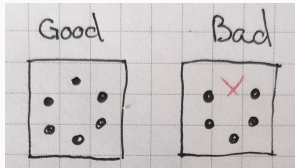
$$x_1^{(2)}, \dots, x_n^{(2)} \sim P_{G_2}$$



1. Fast

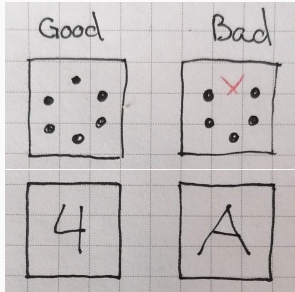
Desiderata

1. Fast
2. Diversity



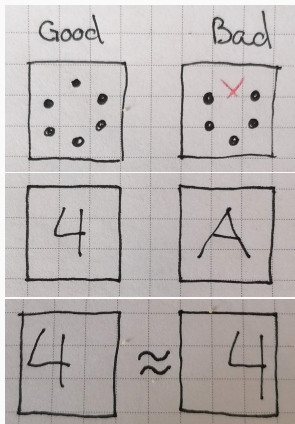
Desiderata

1. Fast
2. Diversity
3. Classifiable



Desiderata

1. Fast
2. Diversity
3. Classifiable
4. Translation invariant



Metrics

First Idea: Inception Score ¹

Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A. and Chen, X., 2016. Improved techniques for training gans. arXiv preprint arXiv:1606.03498.

Idea: Use Inception network $N(x) = p(y|x)$ to “analyze” the generated data.

Idea: Use Inception network $N(x) = p(y|x)$ to “analyze” the generated data.

Easily classifiable If N is confident in predictions \Rightarrow better sample quality.

Idea: Use Inception network $N(x) = p(y|x)$ to “analyze” the generated data.

Easily classifiable If N is confident in predictions \Rightarrow better sample quality.

Diversity If all the classes are represented, the samples are diverse.

Easily classifiable If N is confident in predictions \Rightarrow better sample quality.

Diversity If all the classes are represented, the samples are diverse.

(1)

(2)

Inception Score

Easily classifiable If N is confident in predictions \Rightarrow better sample quality.

Diversity If all the classes are represented, the samples are diverse.

$$IS(X) = \exp \left\{ \mathbb{E}_x \left[\text{KL} (p(y | x) \| p(y)) \right] \right\} \quad (1)$$

(2)

Inception Score

Easily classifiable If N is confident in predictions \Rightarrow better sample quality.

Diversity If all the classes are represented, the samples are diverse.

$$IS(X) = \exp \left\{ \mathbb{E}_x \left[\text{KL} (p(y | x) \| p(y)) \right] \right\} \quad (1)$$

$$= \exp \left\{ H(y) - \mathbb{E}_x [H(y|x)] \right\} \quad (2)$$

Inception Score

Easily classifiable If N is confident in predictions \Rightarrow better sample quality.

Diversity If all the classes are represented, the samples are diverse.

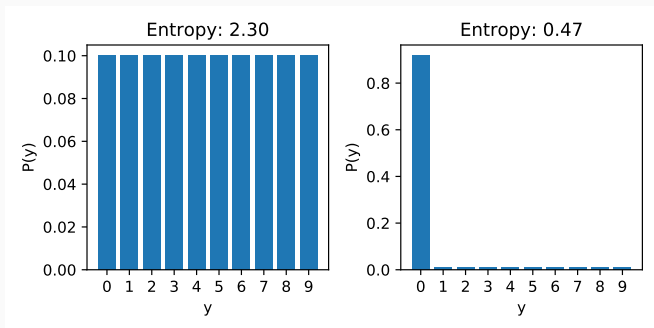
$$IS(X) = \exp \left\{ \mathbb{E}_x \left[\text{KL} (p(y | x) \| p(y)) \right] \right\} \quad (1)$$

$$= \exp \left\{ H(y) - \mathbb{E}_x [H(y|x)] \right\} \quad (2)$$

$$IS(X) = \exp\{ H(y) - \mathbb{E}_x[H(y|x)] \}$$

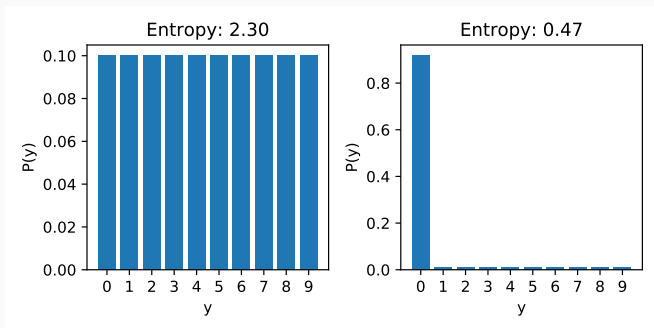
Inception Score

$$IS(X) = \exp\{ H(y) - \mathbb{E}_x[H(y|x)] \}$$



Inception Score

$$IS(X) = \exp \left\{ \underbrace{H(y)}_{\text{Maximize}} - \underbrace{\mathbb{E}_x[H(y|x)]}_{\text{Minimize}} \right\} \quad (2)$$



Correlates with human judgement **but** doesn't
take P_d into account!

Second Idea: Fréchet Inception Distance ²

² Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B. and Hochreiter, S., 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. arXiv preprint arXiv:1706.08500.

Fréchet Inception Distance

Idea: Compare Inception network encodings between P_d and P_f .

Fréchet Inception Distance

Idea: Compare Inception network encodings between P_d and P_f .

Given means μ_r, μ_f and covariances Σ_r, Σ_f of Inception encodings, the Fréchet Inception Distance (FID) is defined as

Fréchet Inception Distance

Idea: Compare Inception network encodings between P_d and P_f .

Given means μ_r, μ_f and covariances Σ_r, Σ_f of Inception encodings, the Fréchet Inception Distance (FID) is defined as

$$\text{FID}(X_r, X_f) = W_2^2(\mathcal{N}\{\mu_r, \Sigma_r\}, \mathcal{N}\{\mu_f, \Sigma_f\}) \quad (3)$$

(4)

Fréchet Inception Distance

Idea: Compare Inception network encodings between P_d and P_f .

Given means μ_r, μ_f and covariances Σ_r, Σ_f of Inception encodings, the Fréchet Inception Distance (FID) is defined as

$$\text{FID}(X_r, X_f) = W_2^2(\mathcal{N}\{\mu_r, \Sigma_r\}, \mathcal{N}\{\mu_f, \Sigma_f\}) \quad (3)$$

$$= \underbrace{\|\mu_r - \mu_f\|_2^2}_{\mathcal{O}(d)} + \quad + \quad - 2 \quad (4)$$

Fréchet Inception Distance

Idea: Compare Inception network encodings between P_d and P_f .

Given means μ_r, μ_f and covariances Σ_r, Σ_f of Inception encodings, the Fréchet Inception Distance (FID) is defined as

$$\text{FID}(X_r, X_f) = W_2^2(\mathcal{N}\{\mu_r, \Sigma_r\}, \mathcal{N}\{\mu_f, \Sigma_f\}) \quad (3)$$

$$= \underbrace{\|\mu_r - \mu_f\|_2^2}_{\mathcal{O}(d)} + \underbrace{\text{Tr}[\Sigma_r]}_{\mathcal{O}(d)} + \underbrace{\text{Tr}[\Sigma_f]}_{\mathcal{O}(d)} - 2 \quad (4)$$

Fréchet Inception Distance

Idea: Compare Inception network encodings between P_d and P_f .

Given means μ_r, μ_f and covariances Σ_r, Σ_f of Inception encodings, the Fréchet Inception Distance (FID) is defined as

$$\text{FID}(X_r, X_f) = W_2^2(\mathcal{N}\{\mu_r, \Sigma_r\}, \mathcal{N}\{\mu_f, \Sigma_f\}) \quad (3)$$

$$= \underbrace{\|\mu_r - \mu_f\|_2^2}_{\mathcal{O}(d)} + \underbrace{\text{Tr}[\Sigma_r]}_{\mathcal{O}(d)} + \underbrace{\text{Tr}[\Sigma_f]}_{\mathcal{O}(d)} - 2\underbrace{\text{Tr}[\sqrt{\Sigma_r \Sigma_f}]}_{\mathcal{O}(d^3)} \quad (4)$$

Fréchet Inception Distance

Idea: Compare Inception network encodings between P_d and P_f .

Given means μ_r, μ_f and covariances Σ_r, Σ_f of Inception encodings, the Fréchet Inception Distance (FID) is defined as

$$\text{FID}(X_r, X_f) = W_2^2(\mathcal{N}\{\mu_r, \Sigma_r\}, \mathcal{N}\{\mu_f, \Sigma_f\}) \quad (3)$$

$$= \underbrace{\|\mu_r - \mu_f\|_2^2}_{\mathcal{O}(d)} + \underbrace{\text{Tr}[\Sigma_r]}_{\mathcal{O}(d)} + \underbrace{\text{Tr}[\Sigma_f]}_{\mathcal{O}(d)} - 2\underbrace{\text{Tr}[\sqrt{\Sigma_r \Sigma_f}]}_{\mathcal{O}(d^3)} \quad (4)$$

Today this is **state-of-the-art**.

Computing $\text{Tr}[\sqrt{\Sigma_r \Sigma_f}]$

Previous method **explicitly** computes $\sqrt{\Sigma_r \Sigma_f}$ and then computes the trace:

Computing $\text{Tr}[\sqrt{\Sigma_r \Sigma_f}]$

Previous method **explicitly** computes $\sqrt{\Sigma_r \Sigma_f}$ and then computes the trace:

Input : Σ_r, Σ_f

Output: $\text{Tr}[C] = \text{Tr} [\sqrt{\Sigma_r \Sigma_f}]$

Computing $\text{Tr}[\sqrt{\Sigma_r \Sigma_f}]$

Previous method **explicitly** computes $\sqrt{\Sigma_r \Sigma_f}$ and then computes the trace:

Input : Σ_r, Σ_f

Output: $\text{Tr}[C] = \text{Tr} [\sqrt{\Sigma_r \Sigma_f}]$

1 $Q, V \leftarrow \text{SchurDecompose}(A);$ $/* \text{ } QVQ^T = A \text{ } */$

Computing $\text{Tr}[\sqrt{\Sigma_r \Sigma_f}]$

Previous method **explicitly** computes $\sqrt{\Sigma_r \Sigma_f}$ and then computes the trace:

Input : Σ_r, Σ_f

Output: $\text{Tr}[C] = \text{Tr} [\sqrt{\Sigma_r \Sigma_f}]$

1 $Q, V \leftarrow \text{SchurDecompose}(A);$

/ QVQ^T = A */*

2 $U \leftarrow \text{TriangSqrt}(V);$

/ V = U² */*

Computing $\text{Tr}[\sqrt{\Sigma_r \Sigma_f}]$

Previous method **explicitly** computes $\sqrt{\Sigma_r \Sigma_f}$ and then computes the trace:

Input : Σ_r, Σ_f

Output: $\text{Tr}[C] = \text{Tr} [\sqrt{\Sigma_r \Sigma_f}]$

```
1  $Q, V \leftarrow \text{SchurDecompose}(A);$            /*  $QVQ^T = A$  */  
2  $U \leftarrow \text{TriangSqrt}(V);$                /*  $V = U^2$  */  
3  $C \leftarrow QUQ^T;$                        /*  $C = \sqrt{\Sigma_r \Sigma_f}$  */
```

Computing $\text{Tr}[\sqrt{\Sigma_r \Sigma_f}]$

Previous method **explicitly** computes $\sqrt{\Sigma_r \Sigma_f}$ and then computes the trace:

Input : Σ_r, Σ_f

Output: $\text{Tr}[C] = \text{Tr} [\sqrt{\Sigma_r \Sigma_f}]$

```
1  $Q, V \leftarrow \text{SchurDecompose}(A);$            /*  $QVQ^T = A$  */
2  $U \leftarrow \text{TriangSqrt}(V);$                /*  $V = U^2$  */
3  $C \leftarrow QUQ^T;$                        /*  $C = \sqrt{\Sigma_r \Sigma_f}$  */
4 return  $\text{Tr}[C];$ 
```

Computing $\text{Tr}[\sqrt{\Sigma_r \Sigma_f}]$

Previous method **explicitly** computes $\sqrt{\Sigma_r \Sigma_f}$ and then computes the trace:

Input : Σ_r, Σ_f

Output: $\text{Tr}[C] = \text{Tr}[\sqrt{\Sigma_r \Sigma_f}]$

```
1  $Q, V \leftarrow \text{SchurDecompose}(A);$            /*  $QVQ^T = A$  */  
2  $U \leftarrow \text{TriangSqrt}(V);$                /*  $V = U^2$  */  
3  $C \leftarrow QUQ^T;$                        /*  $C = \sqrt{\Sigma_r \Sigma_f}$  */  
4 return  $\text{Tr}[C];$ 
```

Line [1-3] each takes **cubic** time!

Fast Fréchet Inception Distance

Idea 3: Don't compute $\text{Tr} [\sqrt{\Sigma_r \Sigma_f}]$, use *eigenvalues* instead.³

³ Mathiasen, A. and Hvilshøj, F., 2020. Fast Fréchet Inception Distance. arXiv preprint arXiv:2009.14075.

Lemma 1

$$\text{Tr}[\sqrt{A}] = \sum_i |\sqrt{\lambda_i(A)}|. \text{ }^4$$

⁴There are some nuances here, please refer to paper for full details.

Lemma 1

$$\text{Tr}[\sqrt{A}] = \sum_i |\sqrt{\lambda_i(A)}|. \quad ^4$$

Lemma 2

*Computing eigenvalues of $d \times d$ matrix A takes $\mathcal{O}(d^3)$ time.
(similar time to compute \sqrt{A})*

⁴There are some nuances here, please refer to paper for full details.

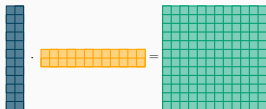
Lemma 3

The nonzero eigenvalues of AB are equal to those of BA , as long as the products are square. ⁵

⁵Nakatsukasa, Y., 2019. The low-rank eigenvalue problem. arXiv preprint arXiv:1905.11490.

Lemma 3

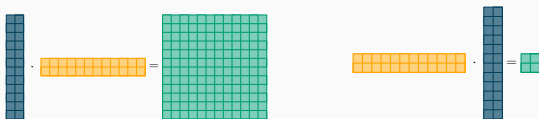
*The nonzero eigenvalues of AB are equal to those of BA , as long as the products are square.*⁵



⁵Nakatsukasa, Y., 2019. The low-rank eigenvalue problem. arXiv preprint arXiv:1905.11490.

Lemma 3

*The nonzero eigenvalues of AB are equal to those of BA , as long as the products are square.*⁵

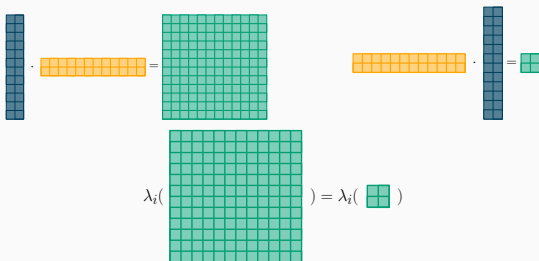


⁵Nakatsukasa, Y., 2019. The low-rank eigenvalue problem. arXiv preprint arXiv:1905.11490.

Lemmas

Lemma 3

*The nonzero eigenvalues of AB are equal to those of BA , as long as the products are square.*⁵



⁵Nakatsukasa, Y., 2019. The low-rank eigenvalue problem. arXiv preprint arXiv:1905.11490.

Fast Fréchet Inception Distance

High level idea: Construct “small” matrix M such that $\lambda_i(M)$ satisfy $\sum_i |\sqrt{\lambda_i(M)}| = \text{Tr}[\sqrt{\Sigma_r \Sigma_f}]$. When M is sufficiently small, computing eigenvalues will be faster than computing $\sqrt{\Sigma_r \Sigma_f}$ explicitly.

Fast Fréchet Inception Distance

Stack the m fake encoded samples into a $d \times m$ matrix X_f .

Fast Fréchet Inception Distance

Stack the m fake encoded samples into a $d \times m$ matrix X_f .

$$\Sigma_f = C_f C_f^T \quad \text{where } C_r = \frac{1}{\sqrt{n-1}} (X_r - \mu_r \mathbf{1}_n) \quad (5)$$

Fast Fréchet Inception Distance

Stack the m fake encoded samples into a $d \times m$ matrix X_f .

$$\Sigma_f = C_f C_f^T \quad \text{where } C_r = \frac{1}{\sqrt{n-1}} (X_r - \mu_r \mathbf{1}_n) \quad (5)$$

Then

$$\Sigma_r \Sigma_f = \Sigma_r C_f C_f^T \quad (6)$$

$$\Sigma_r \Sigma_f = \Sigma_r C_f C_f^T \quad (6)$$

Fast Fréchet Inception Distance

$$\Sigma_r \Sigma_f = \Sigma_r C_f C_f^T \quad (6)$$

Using Lemma 3:

$$\lambda_i(\underbrace{\Sigma_r C_f C_f^T}_{d \times d}) = \lambda_i(\underbrace{C_f^T \Sigma_r C_f}_{m \times m}) \quad (7)$$

Fast Fréchet Inception Distance

$$\Sigma_r \Sigma_f = \Sigma_r C_f C_f^T \quad (6)$$

Using Lemma 3:

$$\lambda_i(\underbrace{\Sigma_r C_f C_f^T}_{d \times d}) = \lambda_i(\underbrace{C_f^T \Sigma_r C_f}_{m \times m}) \quad (7)$$

Eigenvalue computations go from $\mathcal{O}(d^3)$ to $\mathcal{O}(m^3)$ (Lemma 2).

Fast Fréchet Inception Distance

$$\Sigma_r \Sigma_f = \Sigma_r C_f C_f^T \quad (6)$$

Using Lemma 3:

$$\lambda_i(\underbrace{\Sigma_r C_f C_f^T}_{d \times d}) = \lambda_i(\underbrace{C_f^T \Sigma_r C_f}_{m \times m}) \quad (7)$$

Eigenvalue computations go from $\mathcal{O}(d^3)$ to $\mathcal{O}(m^3)$ (Lemma 2).

Finally due to Lemma 1:

$$\text{Tr} \left[\sqrt{\Sigma_r \Sigma_f} \right] = \sum_{i=1}^{m-1} \left| \sqrt{\lambda_i(C_f^T \Sigma_r C_f)} \right| \quad (8)$$

Fast Fréchet Inception Distance

Overall, we get runningtime

$$\text{FID} = \underbrace{\|\mu_r - \mu_f\|_2^2}_{\mathcal{O}(d)} + \underbrace{\text{Tr}[\Sigma_r + \Sigma_f]}_{\mathcal{O}(d)} - 2 \sum_{i=1}^{m-1} \underbrace{\left| \sqrt{\lambda_i (\mathbf{C}_f^T \Sigma_r \mathbf{C}_f)} \right|}_{\mathcal{O}(d^2 m + m^3)} \quad (9)$$

What is it Good For? 🎵

The Greater Perspective

$$\text{FID} = \underbrace{\|\mu_r - \mu_f\|_2^2 + \text{Tr}[\Sigma_r + \Sigma_f]}_{\mathcal{O}(d)} - 2 \sum_{i=1}^{m-1} \underbrace{\left| \sqrt{\lambda_i (\mathbf{C}_f^T \Sigma_r \mathbf{C}_f)} \right|}_{\mathcal{O}(d^2 m + m^3)} \quad (9)$$

The Greater Perspective

$$\text{FID} = \underbrace{\|\mu_r - \mu_f\|_2^2 + \text{Tr}[\Sigma_r + \Sigma_f]}_{\mathcal{O}(d)} - 2 \sum_{i=1}^{m-1} \underbrace{\left| \sqrt{\lambda_i(\textcolor{teal}{C}_f^T \Sigma_r \textcolor{teal}{C}_f)} \right|}_{\mathcal{O}(d^2m+m^3)} \quad (9)$$

During training, we typically have $n \gg d \gg m$

The Greater Perspective

$$\text{FID} = \underbrace{\|\mu_r - \mu_f\|_2^2 + \text{Tr}[\Sigma_r + \Sigma_f]}_{\mathcal{O}(d)} - 2 \sum_{i=1}^{m-1} \underbrace{\left| \sqrt{\lambda_i (\textcolor{teal}{C}_f^T \Sigma_r \textcolor{teal}{C}_f)} \right|}_{\mathcal{O}(d^2 m + m^3)} \quad (9)$$

During training, we typically have $n \gg d \gg m$

Example 4

For GANs on ImageNet, test size (n) is 10 000, encodings (d) are 2048, and batch size (m) is typically 128.

The Greater Perspective

$$\text{FID} = \underbrace{\|\mu_r - \mu_f\|_2^2 + \text{Tr}[\Sigma_r + \Sigma_f]}_{\mathcal{O}(d)} - 2 \sum_{i=1}^{m-1} \underbrace{\left| \sqrt{\lambda_i (C_f^T \Sigma_r C_f)} \right|}_{\mathcal{O}(d^2 m + m^3)} \quad (9)$$

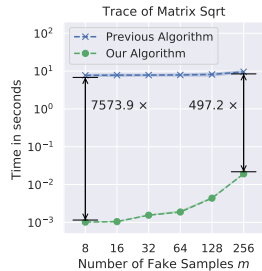
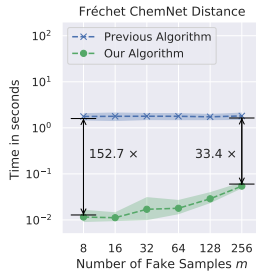
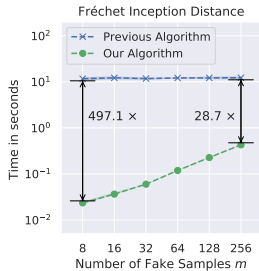
During training, we typically have $n \gg d \gg m$

Example 4

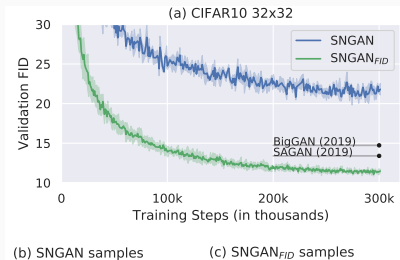
For GANs on ImageNet, test size (n) is 10 000, encodings (d) are 2048, and batch size (m) is typically 128.

💡 Let's use FID for optimizations!

Performance



Minimizing FID



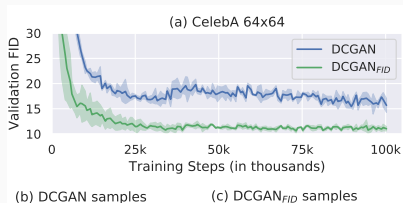
GAN



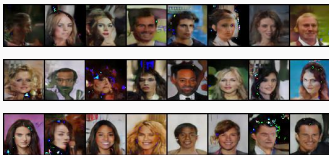
GAN_{FID}



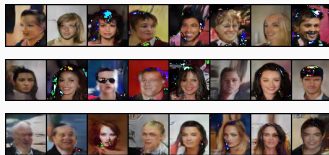
Minimizing FID



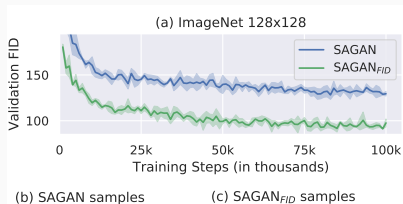
GAN



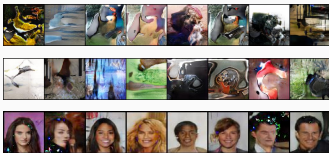
GAN_{FID}



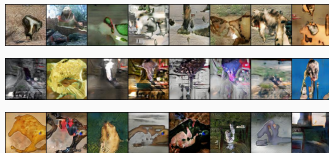
Minimizing FID



GAN



GAN_{FID}

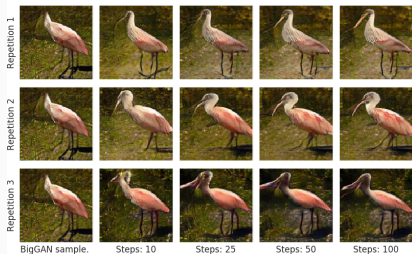


Can FID Loss Improve Generated Images?

What will happen if we just optimize for FID?

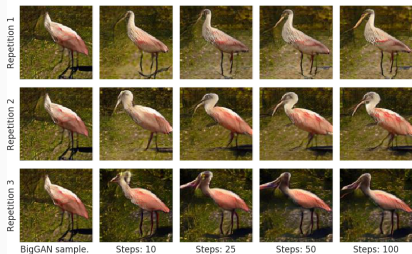
Can FID Loss Improve Generated Images?

What will happen if we just optimize for FID?



Can FID Loss Improve Generated Images?

What will happen if we just optimize for FID?



$$\underbrace{\|\mu_r - \mu_f\|_2^2}_{\text{mean difference}} + \underbrace{\text{Tr} [\Sigma_r] + \text{Tr} [\Sigma_f] - 2\text{Tr} \left[\sqrt{\Sigma_r \Sigma_f} \right]}_{\text{covariance difference}} \quad (10)$$