# QUALITY AND DIVERSITY IN GENERATIVE MODELS THROUGH THE LENS OF f-DIVERGENCES

#### Alexandre Vérine

Centre Sciences des Données DI, École Normale Supérieure PSL











# THE VARIOUS PERFORMANCES OF GENERATIVE MODELS MOTIVATION

As the generation becomes better, the evaluation becomes more challenging.



#### DALL·E 2 (2023)

#### Midjourney v5 (2023)



#### Prompt: A dog playing with a child.

# THE VARIOUS PERFORMANCES OF GENERATIVE MODELS MOTIVATION

As the generation becomes better, the evaluation becomes more challenging.



DALL·E 2 (2023)

### $\neq$

Midjourney v5 (2023)



Prompt: A dog playing with a child.

VERINE - SEMINAR LPSM - QUALITY AND DIVERSITY IN GENERATIVE MODELS

In this presentation, we discuss on evaluating, optimizing and improving quality and diversity of generative models:

1. Evaluating: How can we assess quality and diversity independently in Generative Models?

2. **Optimizing:** Can we optimize a specific trade-off between quality and diversity?

3. Improving: How can we improve the quality and diversity of a pre-trained generative models?

SECTION 1: EVALUATING

### **Evaluating:**

### How can we assess quality and diversity independtly in Generative Models?

#### GENERATIVE MODELS FRAMEWORK



• Assumption: There is an unknown *target distribution* P in  $\mathcal{X} \subset \mathbb{R}^d$ .



• Assumption: There is an unknown *target distribution* P in  $\mathcal{X} \subset \mathbb{R}^d$ .

#### GENERATIVE MODELS Framework



- Assumption: There is an unknown *target distribution* P in  $\mathcal{X} \subset \mathbb{R}^d$ .
- Goal: Learn a *parameterized distribution*  $\hat{P}$  that approximate *P*:

### GENERATIVE MODELS



- Assumption: There is an unknown *target distribution* P in  $\mathcal{X} \subset \mathbb{R}^d$ .
- Goal: Learn a *parameterized distribution*  $\hat{P}$  that approximate *P*:
  - 1. Consider a distribution Q in *a latent space*  $\mathcal{X} \subset \mathbb{R}^m$ , usually  $\mathcal{N}(0, I_m)$ .

### GENERATIVE MODELS



- Assumption: There is an unknown *target distribution* P in  $\mathcal{X} \subset \mathbb{R}^d$ .
- Goal: Learn a *parameterized distribution*  $\widehat{P}_G$  that approximate *P*:
  - 1. Consider a distribution Q in *a latent space*  $\mathcal{X} \subset \mathbb{R}^m$ , usually  $\mathcal{N}(0, I_m)$ .
  - 2. Take *a generator model G* represented by a neural network. Take  $\hat{P}_G = G \# Q$ .

#### GENERATIVE MODELS Framework

- Assumption: There is an unknown *target distribution* P in  $\mathcal{X} \subset \mathbb{R}^d$ .
- Goal: Learn a *parameterized distribution*  $\hat{P}_G$  that approximate *P*:
  - 1. Consider a distribution Q in *a latent space*  $\mathcal{X} \subset \mathbb{R}^m$ , usually  $\mathcal{N}(0, I_m)$ .
  - 2. Take *a generator model G* represented by a neural network. Take  $\hat{P}_G = G \# Q$ .
  - 3. Compute  $G^{\text{opt}}$  that minimize *a dissimilarity measure D* between *P* and  $\hat{P}_G$ :

$$G^{\text{opt}} = \operatorname*{argmin}_{G} D(P, \widehat{P}_G)$$

### GENERATIVE MODELS



- Assumption: There is an unknown *target distribution* P in  $\mathcal{X} \subset \mathbb{R}^d$ .
- Goal: Learn a *parameterized distribution*  $\hat{P}$  that approximate *P*:
  - 1. Consider a distribution Q in *a latent space*  $\mathcal{X} \subset \mathbb{R}^m$ , usually  $\mathcal{N}(0, I_m)$ .
  - 2. Take *a generator model G* represented by a neural network. Take  $\hat{P} = G \# Q$ .
  - 3. Compute  $G^{\text{opt}}$  that minimize *a dissimilarity measure D* between *P* and  $\widehat{P}$ :

$$G^{\text{opt}} = \operatorname*{argmin}_{G} D(P, \widehat{P})$$

#### GENERATIVE MODELS IN PRACTICE



#### GENERATIVE MODELS IN PRACTICE



### Low Diversity

#### GENERATIVE MODELS IN PRACTICE



#### PRECISION AND RECALL FOR GENERATIVE MODELS METRICS TO EVALUATE QUALITY AND DIVERSITY

To assess models, we use the notion of Precision and Recall, inspired from Information Retrieval:

### Quality



#### PRECISION AND RECALL FOR GENERATIVE MODELS METRICS TO EVALUATE QUALITY AND DIVERSITY

To assess models, we use the notion of Precision and Recall, inspired from Information Retrieval:



What proportion of generated samples are realistic?

#### PRECISION AND RECALL FOR GENERATIVE MODELS METRICS TO EVALUATE QUALITY AND DIVERSITY

To assess models, we use the notion of Precision and Recall, inspired from Information Retrieval:



What proportion of generated samples are realistic? What proportion of real samples can be generated?



#### Definition 2.1 (Support-Based Precision and Recall - Kynkäänniemi et al. [8].)

For any distributions  $P \in \mathcal{P}(\mathcal{X})$  and  $\widehat{P} \in \mathcal{P}(\mathcal{X})$ , we say that the distribution P has precision  $\overline{\alpha}$  at recall  $\overline{\beta}$  with respect to  $\widehat{P}$  if

$$\bar{\alpha} \coloneqq \widehat{P}(\operatorname{Supp}(P)) \quad and \quad \bar{\beta} \coloneqq P(\operatorname{Supp}(\widehat{P})). \tag{1}$$

Precision for finite support is the proportion of generated data that lies on the support of the real data:

$$\bar{\alpha} = \widehat{P}(\operatorname{Supp}(P)).$$



Precision for finite support is the proportion of generated data that lies on the support of the real data:

$$\bar{\alpha} = \widehat{P}(\operatorname{Supp}(P)).$$



Recall for finite support is the proportion of the support of the real data that is covered by the generated data:

$$\bar{\beta} = P(\operatorname{Supp}(\widehat{P})).$$



Recall for finite support is the proportion of the support of the real data that is covered by the generated data:

$$\bar{\beta} = P(\operatorname{Supp}(\widehat{P})).$$















# PRECISION AND RECALL FOR GENERATIVE MODELS IN PRACTICE

MNIST Dataset [14]



# PRECISION AND RECALL FOR GENERATIVE MODELS IN PRACTICE



Precision: 0.54 Recall: 0.91

Precision: 0.80 Recall: 0.70


# PRECISION AND RECALL FOR GENERATIVE MODELS FOR LLMS

On open-ended generation, the quality and diversity of LLMs can be evaluated using Precision and Recall for instance on Webtext: Bronnec et al. [3]



# $\begin{array}{l} Precision \text{ and } Recall \text{ for } Generative \text{ Models} \\ \hline \text{For } \text{LLMs} \end{array}$

We can also evaluate the quality and diversity of LLMs on Chatbot open-ended generation. We can for instance check the impact of In-Context examples on the quality and diversity of the generated text. For instance on Wikipedia Biographies generation:



#### PRECISION AND RECALL FOR GENERATIVE MODELS FOR INIFINITE SUPPORT



#### PRECISION AND RECALL FOR GENERATIVE MODELS



#### Both distributions have **perfect** Precision *and* Recall.

#### **Definition 2.2 (PR-Curve for Generative Models - Sajjadi et al. [11], Simon et al. [12])** Let $P, \hat{P} \in \mathcal{P}(\mathcal{X})$ be two distributions such that $P, \hat{P} \ll \mu$ . The PR-Curve is the set $PRD(P, \hat{P})$ defined as:

$$PRD(P,\widehat{P}) = \{(\alpha_{\lambda}, \beta_{\lambda}) \mid \lambda \in [0, \infty]\}$$
(2)

with:

$$\alpha_{\lambda} = \int_{\mathcal{X}} \min\left(\lambda p(\mathbf{x}), \hat{p}(\mathbf{x})\right) d\mu(\mathbf{x}) \quad and \quad \beta_{\lambda} = \int_{\mathcal{X}} \min\left(p(\mathbf{x}), \hat{p}(\mathbf{x})/\lambda\right) d\mu(\mathbf{x}). \tag{3}$$

For the Precision,  $\lambda p$  is compared to  $\hat{p}$  for different threshold  $\lambda \in [0, +\infty]$ :

$$\alpha_{\lambda} = \int_{\mathcal{X}} \min\left(\lambda p(\mathbf{x}), \hat{p}(\mathbf{x})\right) d\mu(\mathbf{x})$$
(4)

For the Precision,  $\lambda p$  is compared to  $\hat{p}$  for different threshold  $\lambda \in [0, +\infty]$ :

$$\alpha_{\lambda} = \int_{\mathcal{X}} \min\left(\lambda p(\mathbf{x}), \hat{p}(\mathbf{x})\right) d\mu(\mathbf{x})$$
(4)

For the Recall, *p* is compared to  $\hat{p}/\lambda$  for different threshold  $\lambda \in [0, +\infty]$ :

$$\beta_{\lambda} = \int_{\mathcal{X}} \min\left(p(\mathbf{x}), \hat{p}(\mathbf{x})/\lambda\right) d\mu(\mathbf{x})$$
(5)

For the Recall, *p* is compared to  $\hat{p}/\lambda$  for different threshold  $\lambda \in [0, +\infty]$ :

$$\beta_{\lambda} = \int_{\mathcal{X}} \min\left(p(\mathbf{x}), \hat{p}(\mathbf{x})/\lambda\right) d\mu(\mathbf{x})$$
(5)



Figure. Learning distribution with low recall and high precision.



Figure. Learning distribution with high recall and low precision.



Figure. Learning distribution with low recall and low precision.



Figure. Learning distribution with high recall and high precision.

# PR-CURVE AND SUPPORT-BASED PRECISION AND RECALL RELATION

The PR-Curve is a generalization of the Precision and Recall for finite support:

#### Theorem 2.3 (Support-based and PR-Curves - Siry et al. [13])

Let  $P, \hat{P} \in \mathcal{P}(\mathcal{X})$  be two distributions. Then, the support-based Precision and Recall  $(\bar{\alpha}, \bar{\beta})$  are related to the PR-Curve values  $PRD(P, \hat{P})$  for  $\lambda = 0$  and  $\lambda = \infty$ :

$$\bar{\alpha} = \max_{\lambda} \alpha_{\lambda} = \alpha_{\infty} \quad and \quad \bar{\beta} = \max_{\lambda} \beta_{\lambda} = \beta_0.$$
 (6)

# PR-CURVE AND SUPPORT-BASED PRECISION AND RECALL RELATION



# PR-CURVE AND SUPPORT-BASED PRECISION AND RECALL RELATION



# PR-CURVE FOR GENERATIVE MODELS IN PRACTICE



Precision: 0.54 Recall: 0.91

Precision: 0.80 Recall: 0.70

# PR-CURVE FOR GENERATIVE MODELS IN NLP



Figure. PR-Curve for distributions journal articles: AG News.

SECTION 2: OPTIMIZING

#### **Optimizing:** Can we optimize a specific trade-off between Precision and Recall?

# TUNING PRECISION AND RECALL IN GENERATIVE MODELS TRUNCATION

Hard Trunctation Karras et al. [6]

Soft Trunctation Kingma and Dhariwal [7]

#### SECTION 2: OPTIMIZING HARD TRUNCATION



**Figure.** From left to right:  $\psi = 0.0$ ,  $\psi = 0.3 \ \psi = 0.7 \ \psi = 1.0$ .



## SECTION 2: OPTIMIZING SOFT TRUNCATION



(a)  $\psi = 0.04$ 





(d)  $\psi = 2.0$ 

Figure. Soft-Truncation on BigGAN. Source:[2].

# TRAINING A GENERATIVE MODEL IN GENERAL

Traditionally, the goal is to minimize *a dissimilarity measure* between the target distribution P and the learned distribution  $\hat{P}$ :

$$\min_{G} D(P, \widehat{P}_{G}) \tag{6}$$



# TRAINING A GENERATIVE MODEL with f-divergences

Traditionally, the goal is to minimize *an* f-*divergence* between the target distribution P and the learned distribution  $\hat{P}$ :

$$\min_{G} \mathcal{D}_{f}(P \| \widehat{P}_{G}) \tag{6}$$



# *f*-DIVERGENCES DEFINITION

#### **Definition 3.1 (***f***-divergences)**

For any two probability distributions P and  $\hat{P}$  in  $\mathcal{P}(\mathcal{X})$  such that  $P, \hat{P} \ll \mu$ . Let p and  $\hat{p}$  be the Radon-Nikodym densities of P and  $\hat{P}$  with respect to  $\mu$ , respectively. Let f be any convex lower semi-continuous function  $f : [0, \infty] \rightarrow ] - \infty, +\infty$ ] such that f(1) = 0, the f-divergence between P and  $\hat{P}$  is

$$\mathcal{D}_{f}(P\|\widehat{P}) = \int_{\mathcal{X}} \widehat{p}(\mathbf{x}) f\left(\frac{p(\mathbf{x})}{\widehat{p}(\mathbf{x})}\right) d\mu(\mathbf{x}).$$
(7)

# *f*-DIVERGENCES DEFINITION

#### **Definition 3.1 (***f***-divergences)**

For any two probability distributions P and  $\hat{P}$  in  $\mathcal{P}(\mathcal{X})$  such that  $P, \hat{P} \ll \mu$ . Let p and  $\hat{p}$  be the Radon-Nikodym densities of P and  $\hat{P}$  with respect to  $\mu$ , respectively. Let f be any convex lower semi-continuous function  $f : [0, \infty] \rightarrow ] - \infty, +\infty$ ] such that f(1) = 0, the f-divergence between P and  $\hat{P}$  is

$$\mathcal{D}_{f}(P\|\widehat{P}) = \int_{\mathcal{X}} \widehat{p}(\mathbf{x}) f\left(\frac{p(\mathbf{x})}{\widehat{p}(\mathbf{x})}\right) d\mu(\mathbf{x}).$$
(7)

Usual divergences are *f*-divergences:

- ► Kullback-Leibler (KL),
- ▶ Reverse Kullback-Leibler (rKL),
- ► Jensen-Shannon (JS),
- ► Total Variation (TV),
- $\alpha$ -divergences.

# ESTIMATING *f*-DIVERGENCES DUAL VARIATIONAL FORM

*f*-divergences are *hardly tractable* in practice. However, they can be approximated by a dual approximation.

# ESTIMATING f-DIVERGENCES DUAL VARIATIONAL FORM

*f*-divergences are *hardly tractable* in practice. However, they can be approximated by a dual approximation.

- $f^*(t) = \sup_{u \in \mathbb{R}} \{tu f(u)\}$  be the Fenchel conjugate of f.
- $\mathcal{T}$  be the set of all measurable functions  $\mathcal{X} \to \mathbb{R}$ .

#### ESTIMATING *f*-DIVERGENCES DUAL VARIATIONAL FORM

*f*-divergences are *hardly tractable* in practice. However, they can be approximated by a dual approximation.

- $f^*(t) = \sup_{u \in \mathbb{R}} \{tu f(u)\}$  be the Fenchel conjugate of f.
- $\mathcal{T}$  be the set of all measurable functions  $\mathcal{X} \to \mathbb{R}$ .

#### Theorem 3.2 (Dual variational form of an *f*-divergence- Nguyen et al. [9])

Let  $P, \widehat{P} \in \mathcal{P}(\mathcal{X})$  two distributions such that P is absolutely continuous with respect to  $\widehat{P}$  and f a suitable generator function. The *f*-divergence between P and  $\widehat{P}$  admits a dual variational form:

$$\mathcal{D}_{f}(P \| \widehat{P}) = \sup_{T \in \mathcal{T}} \left( \mathbb{E}_{\boldsymbol{x} \sim P} \left[ T(\boldsymbol{x}) \right] - \mathbb{E}_{\boldsymbol{x} \sim \widehat{P}} \left[ f^{*}(T(\boldsymbol{x})) \right] \right).$$
(8)

We use  $T^{opt} \in \mathcal{T}$  to denote the function that achieves the supremum.

By doing so, we can rewrite the optimization problem as:

$$\min_{G} \max_{T} \underbrace{\mathbb{E}_{\boldsymbol{x} \sim P}\left[T(\boldsymbol{x})\right] - \mathbb{E}_{\boldsymbol{x} \sim \widehat{P}_{G}}\left[f^{*}(T(\boldsymbol{x}))\right]}_{\mathcal{D}_{f,T}^{\text{dual}}}$$

By doing so, we can rewrite the optimization problem as:

$$\min_{G} \max_{T} \underbrace{\mathbb{E}_{\boldsymbol{x} \sim P}\left[T(\boldsymbol{x})\right] - \mathbb{E}_{\boldsymbol{x} \sim \widehat{P}_{G}}\left[f^{*}(T(\boldsymbol{x}))\right]}_{\mathcal{D}_{f,T}^{\text{dual}}}$$

- ▶ The discriminator *T* is trained *to estimate* the divergence.
- ▶ The generator *G* is trained *to minimize* the divergence.

By doing so, we can rewrite the optimization problem as:

$$\min_{G} \max_{T} \underbrace{\mathbb{E}_{\boldsymbol{x} \sim P}\left[T(\boldsymbol{x})\right] - \mathbb{E}_{\boldsymbol{x} \sim \widehat{P}_{G}}\left[f^{*}(T(\boldsymbol{x}))\right]}_{\mathcal{D}_{f,T}^{\text{dual}}}$$

- ▶ The discriminator *T* is trained *to estimate* the divergence.
- ▶ The generator *G* is trained *to minimize* the divergence.

By doing so, we can rewrite the optimization problem as:

$$\min_{G} \max_{T} \mathbb{E}_{\boldsymbol{x} \sim P} \left[ \log \left( D(\boldsymbol{x}) \right) \right] - \mathbb{E}_{\boldsymbol{x} \sim \widehat{P}_{G}} \left[ f^{*} \left( \log(D(\boldsymbol{x})) \right) \right]$$

- ▶ The discriminator *T* is trained *to estimate* the divergence.
- ▶ The generator *G* is trained *to minimize* the divergence.
- With  $T(\mathbf{x}) = \log(D(\mathbf{x}))$  with  $D(\mathbf{x}) \in [0, 1]$ .

By doing so, we can rewrite the optimization problem as:

$$\min_{G} \max_{T} \mathbb{E}_{\boldsymbol{x} \sim P} \left[ \log \left( D(\boldsymbol{x}) \right) \right] + \mathbb{E}_{\boldsymbol{x} \sim \widehat{P}_{G}} \left[ \log \left( 1 - D(\boldsymbol{x}) \right) \right]$$
(9)

- ▶ The discriminator *T* is trained *to estimate* the divergence.
- ▶ The generator *G* is trained *to minimize* the divergence.
- With  $T(\mathbf{x}) = \log(D(\mathbf{x}))$  with  $D(\mathbf{x}) \in [0, 1]$ .
- ►  $f^*(t) = f^*_{IS}(t) = -\log(1 \exp(t))$  for the Jensen-Shannon divergence.

We recover the original GAN framework.

By doing so, we can rewrite the optimization problem as:

$$\min_{G} \max_{T} \underbrace{\mathbb{E}_{\boldsymbol{x} \sim P}\left[T(\boldsymbol{x})\right] - \mathbb{E}_{\boldsymbol{x} \sim \widehat{P}_{G}}\left[f^{*}(T(\boldsymbol{x}))\right]}_{\mathcal{D}_{f,T}^{\text{dual}}}$$

- ▶ The discriminator *T* is trained *to estimate* the divergence.
- ▶ The generator *G* is trained *to minimize* the divergence.
- Generative Adversarial Networks [4] for *the Jensen-Shannon divergence*.
- Extended to other *f*-divergences by Nowozin et al. [10].
- Extend to other generative models such as Normalizing Flows by Grover et al. [5].

Effect of the f-divergence on the learned distribution

All *f*-divergences are not equal:

$$\mathcal{D}_{f}(P \| \widehat{P}) = \mathbb{E}_{\mathbf{x} \sim \widehat{P}} \left[ f \left( \frac{p(\mathbf{x})}{\widehat{p}(\mathbf{x})} \right) \right]$$


Effect of the f-divergence on the learned distribution

All *f*-divergences are not equal:



Effect of the f-divergence on the learned distribution

All *f*-divergences are not equal:



## Examples of f-divergence minimization

## Examples of f-divergence minimization



VERINE - SEMINAR LPSM - QUALITY AND DIVERSITY IN GENERATIVE MODELS

### SECTION 2: OPTIMIZING CONTRIBUTIONS

### Can we optimize a specific trade-off between Precision and Recall?

- ▶ What is the relation between the Precision-Recall curve and *f*-divergences?
- ► Can we optimize a point on the Precision-Recall curve using *f*-divergences?

# SECTION 2: OPTIMIZING CONTRIBUTIONS

### Can we optimize a specific trade-off between Precision and Recall?

- ▶ What is the relation between the Precision-Recall curve and *f*-divergences?
- ► Can we optimize a point on the Precision-Recall curve using *f*-divergences?

#### List of contributions:

- ▶ We show that the PR-Divergence is an *f*-divergence.
- ▶ We show that any *f*-divergence can be written as a weighted average PR-Divergences.
- We propose an algorithm to optimize the PR-Divergence.

# PRECISION-RECALL DIVERGENCE DEFINITION

#### **Definition 3.3 (PR-Divergence generator function** $f_{\lambda}$ **)**

*Given a trade-off parameter*  $\lambda \in [0, +\infty]$ *, we define the generator function*  $f_{\lambda} : [0, +\infty] \rightarrow ] -\infty, +\infty]$  *given by* 

$$f_{\lambda}(u) = \begin{cases} \max(\lambda u, 1) - \max(\lambda, 1) & \text{for } \lambda \in [0, +\infty[, \\ \mathbb{1}_{\{u=0\}} & \text{for } \lambda = +\infty. \end{cases}$$
(10)



# PRECISION-RECALL DIVERGENCE PROPERTIES

#### **Proposition 3.4 (PR-Divergence)**

For any distributions  $P, \hat{P} \in \mathcal{P}(\mathcal{X})$  such that  $P, \hat{P} \ll \mu$ , then for any  $\lambda \in [0, +\infty]$  the PR-Divergence defined as

$$\mathcal{D}_{\lambda-\mathrm{PR}}(P\|\widehat{P}) = \int_{\mathcal{X}} \widehat{p}(\mathbf{x}) f_{\lambda}\left(\frac{p(\mathbf{x})}{\widehat{p}(\mathbf{x})}\right) d\mu(\mathbf{x})$$
(11)

belongs to the class of *f*-divergences.

#### PRECISION-RECALL DIVERGENCE LINKING THE PR-DIVERGENCE TO THE PR-CURVE

#### Theorem 3.5 (PR-Curves as a function of $\mathcal{D}_{\lambda-PR}$ )

Given  $P, \hat{P} \in \mathcal{P}(\mathcal{X})$  such that  $P, \hat{P} \ll \mu$  and  $\lambda \in [0, +\infty]$ , the *PR-Curve*  $\partial PRD$  is related to the *PR-Divergence*  $\mathcal{D}_{\lambda-PR}(P || \hat{P})$  as follows.

$$\alpha_{\lambda}(P\|\widehat{P}) = \min(1,\lambda) - \mathcal{D}_{\lambda-\mathrm{PR}}(P\|\widehat{P}).$$
  
$$\beta_{\lambda}(P\|\widehat{P}) = \min(1,\lambda) - \mathcal{D}_{\lambda-\mathrm{PR}}(\widehat{P}\|P).$$

#### PRECISION-RECALL DIVERGENCE LINKING THE PR-DIVERGENCE TO THE PR-CURVE

#### Theorem 3.5 (PR-Curves as a function of $\mathcal{D}_{\lambda-PR}$ )

Given  $P, \hat{P} \in \mathcal{P}(\mathcal{X})$  such that  $P, \hat{P} \ll \mu$  and  $\lambda \in [0, +\infty]$ , the *PR-Curve*  $\partial PRD$  is related to the *PR-Divergence*  $\mathcal{D}_{\lambda-PR}(P || \hat{P})$  as follows.

$$\begin{aligned} \alpha_{\lambda}(P\|\widehat{P}) &= \min(1,\lambda) - \mathcal{D}_{\lambda\text{-PR}}(P\|\widehat{P}).\\ \beta_{\lambda}(P\|\widehat{P}) &= \min(1,\lambda) - \mathcal{D}_{\lambda\text{-PR}}(\widehat{P}\|P). \end{aligned}$$

A direct consequence of Theorem 3.5:

 $\operatorname*{argmin}_{\widehat{P} \in \mathcal{P}(\mathcal{X})} \mathcal{D}_{\lambda-\Pr}(P \| \widehat{P}) = \operatorname*{argmax}_{\widehat{P} \in \mathcal{P}(\mathcal{X})} \alpha_{\lambda}(P \| \widehat{P}).$ 

#### EXPLAINING QUALITY / DIVERSITY CONNECTION BETWEEN PR-DIVERGENCE AND *f*-DIVERGENCES

# Theorem 3.6 (*f*-divergences as a weighted average of PR-Divergences)

*For any*  $P, \hat{P} \in \mathcal{P}(\mathcal{X})$  *supported on all*  $\mathcal{X}$  *and any*  $\lambda \in [0, +\infty]$ *, then:* 

$$\mathcal{D}_{f}(P\|\widehat{P}) = \int_{0}^{\infty} rac{1}{\lambda^{3}} f''\left(rac{1}{\lambda}
ight) \mathcal{D}_{\lambda ext{-PR}}(P\|\widehat{P}) \mathrm{d}\lambda,$$



# OPTIMIZING THE PR-DIVERGENCE EXAMPLES



VERINE - SEMINAR LPSM - QUALITY AND DIVERSITY IN GENERATIVE MODELS

# OPTIMIZING THE PR-DIVERGENCE EXAMPLES

### OPTIMIZING THE PR-DIVERGENCE IN PRACTICE

If we train a model to minimize the PR-Divergence, we can use the dual variational form:

$$\min_{G} \max_{T} \mathbb{E}_{\boldsymbol{x} \sim P} \left[ T(\boldsymbol{x}) \right] - \mathbb{E}_{\boldsymbol{x} \sim \widehat{P}_{G}} \left[ f_{\lambda}^{*}(T(\boldsymbol{x})) \right].$$
(12)

### OPTIMIZING THE PR-DIVERGENCE IN PRACTICE

If we train a model to minimize the PR-Divergence, we can use the dual variational form:



The naive approach *fails* to optimize the PR-Divergence.

### OPTIMIZING THE PR-DIVERGENCE IN PRACTICE

If we train a model to minimize the PR-Divergence, we can use the dual variational form:



We propose *a new approach* to optimize the PR-Divergence.

We choose an auxiliary function *g* to train  $T_g$  is trained to estimate the *f*-divergence  $\mathcal{D}_g$ :

$$\max_{T} \mathbb{E}_{\boldsymbol{x} \sim P} \left[ T(\boldsymbol{x}) \right] - \mathbb{E}_{\boldsymbol{x} \sim \widehat{P}} \left[ g^*(T(\boldsymbol{x})) \right]$$
(13)

We choose an auxiliary function *g* to train  $T_g$  is trained to estimate the *f*-divergence  $\mathcal{D}_g$ :

$$\max_{T} \mathbb{E}_{\boldsymbol{x} \sim P} \left[ T(\boldsymbol{x}) \right] - \mathbb{E}_{\boldsymbol{x} \sim \widehat{P}} \left[ g^*(T(\boldsymbol{x})) \right]$$
(13)

At optimality, we have:

$$\nabla g^*(T_g^{\text{opt}}(\mathbf{x})) = \frac{p(\mathbf{x})}{\widehat{p}(\mathbf{x})}.$$
(14)

We choose an auxiliary function *g* to train  $T_g$  is trained to estimate the *f*-divergence  $\mathcal{D}_g$ :

$$\max_{T} \mathbb{E}_{\boldsymbol{x} \sim P} \left[ T(\boldsymbol{x}) \right] - \mathbb{E}_{\boldsymbol{x} \sim \widehat{P}} \left[ g^{*}(T(\boldsymbol{x})) \right]$$
(13)

At optimality, we have:

$$\nabla g^*(T_g^{\text{opt}}(\mathbf{x})) = \frac{p(\mathbf{x})}{\widehat{p}(\mathbf{x})}.$$
(14)

Any *f*-divergence can be computed using the primal estimation as follows using  $T_g$ :

$$\mathcal{D}_{f}(P\|\widehat{P}) = \mathbb{E}_{\mathbf{x}\sim\widehat{P}}\left[f\left(\frac{p(\mathbf{x})}{\widehat{p}(\mathbf{x})}\right)\right] = \mathbb{E}_{\mathbf{x}\sim\widehat{P}}\left[f\left(\nabla g^{*}(T_{g}^{\text{opt}}(\mathbf{x}))\right)\right].$$
(15)

#### CONVERGENCE OF THE PROPOSED APPROACH BOUNDING THE ESTIMATION ERROR

#### Theorem 3.7 (Bound on the estimation of an *f*-divergence using an auxiliary *g*-divergence)

Let  $f, g : \mathbb{R}^+ \to \mathbb{R}$  be such that g is  $\mu$ -strongly convex, f is  $\sigma$ -Lipschitz, and  $\mathcal{D}_f$ ,  $\mathcal{D}_g$  be f-divergences. For any discriminator  $T : \mathcal{X} \to \operatorname{dom}(g^*)$ , let  $r(\mathbf{x}) = \nabla g^*(T(\mathbf{x}))$ . Then:

$$\mathcal{D}_{g}(P\|\widehat{P}) - \mathcal{D}_{g,T}^{\text{dual}} \le \epsilon \implies \left| \mathcal{D}_{f}(P\|\widehat{P}) - \mathcal{D}_{f,T}^{\text{primal}}(P\|\widehat{P}) \right| \le \sigma \sqrt{\frac{2\epsilon}{\mu}}.$$
(16)

#### CONVERGENCE OF THE PROPOSED APPROACH BOUNDING THE ESTIMATION ERROR

#### Theorem 3.7 (Bound on the estimation of an *f*-divergence using an auxiliary *g*-divergence)

Let  $f, g : \mathbb{R}^+ \to \mathbb{R}$  be such that g is  $\mu$ -strongly convex, f is  $\sigma$ -Lipschitz, and  $\mathcal{D}_f$ ,  $\mathcal{D}_g$  be f-divergences. For any discriminator  $T : \mathcal{X} \to \operatorname{dom}(g^*)$ , let  $r(\mathbf{x}) = \nabla g^*(T(\mathbf{x}))$ . Then:

$$\mathcal{D}_{g}(P\|\widehat{P}) - \mathcal{D}_{g,T}^{\text{dual}} \leq \epsilon \Longrightarrow \left| \mathcal{D}_{f}(P\|\widehat{P}) - \mathcal{D}_{f,T}^{\text{primal}}(P\|\widehat{P}) \right| \leq \sigma \sqrt{\frac{2\epsilon}{\mu}}.$$
(16)  
The smaller the error on *T* is, the smaller the error on estimating  $\mathcal{D}_{f}$  using *T* is.

# CONVERGENCE OF THE PROPOSED APPROACH EXAMPLES



# OPTIMIZING THE PR-DIVERGENCE WITH OUR APPROACH IN PRACTICE



# OPTIMIZING THE PR-DIVERGENCE WITH OUR APPROACH TRAINING GANS

Model		CIFA	R-10 32	× 32	Cele	CelebA $64 \times 64$				
		FID	Р	R	FID	Р	R			
Baseline	Big-	13.37	86.51	65.66	9.16	78.41	51.42			
GAN										
$\lambda = 0.05$		13.29	81.10	70.63	-	-	-			
$\lambda = 0.1$		11.62	81.78	74.58	-	-	-			
$\lambda = 0.2$		13.36	84.85	65.13	8.79	83.37	44.07			
$\lambda = 0.5$		14.50	83.27	68.23	6.03	77.60	55.98			
$\lambda = 1.0$		14.03	83.04	69.35	13.07	81.70	36.85			
$\lambda = 2.0$		16.94	84.93	59.79	14.23	82.98	32.87			
$\lambda = 5.0$		32.54	83.39	56.94	22.45	83.96	25.81			
$\lambda = 10.0$		39.69	84.11	39.29	-	-	-			
$\lambda = 20.0$		67.03	90.03	21.81	-	-	-			



When  $\lambda$  increases,  $\begin{cases}
Precision \uparrow \\
Recall \downarrow
\end{cases}$ 

 $\lambda = 0.1$ 

 $\lambda = 10$ 

#### **OPTIMIZING THE PR-DIVERGENCE WITH OUR APPROACH** FINE-TUNING GANS

Model	ImageNet 128 × 128			FFHQ 256 × 256		
	FID	Р	R	FID	Р	R
Baseline BigGAN	9.83	28.04	41.21	41.41	65.57	10.17
Soft $\psi = 0.7$	11.39	23.04	31.13	56.43	76.59	4.87
Soft $\psi = 0.5$	15.49	20.20	19.83	82.05	84.48	1.58
Hard $\psi = 2.0$	9.69	25.83	39.89	43.32	68.84	8.66
Hard $\psi = 1.0$	12.12	21.86	35.42	56.19	76.44	4.76
Hard $\psi = 0.5$	15.21	21.13	29.55	71.32	80.99	4.84
$\lambda = 0.2$	9.92	26.69	42.04	35.66	78.70	9.45
$\lambda = 0.5$	10.82	26.83	42.38	35.24	78.41	9.66
$\lambda = 1.0$	20.42	29.72	28.21	35.91	78.95	8.32
$\lambda = 2.0$	20.21	30.27	30.49	36.33	81.10	8.69
$\lambda = 5.0$	20.76	30.87	28.38	38.16	84.31	8.52



**SECTION 3: IMPROVING** 

# **Improving:** How can we improve the quality and diversity of a pre-trained generative models?

### SAMPLING FROM A GENERATIVE MODEL General Setting

To sample a point from the learned distribution  $\hat{P}$ :

- Sample  $z \sim Q$ .
- Compute x = G(z).



#### SAMPLING FROM A GENERATIVE MODEL General Setting

To sample a point from the learned distribution  $\hat{P}$ :

- Sample  $z \sim Q$ .
- Compute x = G(z).

 $P \neq \widehat{P}$ 



### SAMPLING FROM A GENERATIVE MODEL General Setting

To sample a point from the learned distribution  $\hat{P}$ :

- Sample  $z \sim Q$ .
- Compute x = G(z).

We have an estimation of  $\frac{p(x)}{\overline{p}(x)}$ .



#### SAMPLING FROM A GENERATIVE MODEL Rejection Sampling

To sample a point from the refined distribution  $\tilde{P}$ :

- Sample  $z \sim Q$ .
- Compute x = G(z).
- Accept *x* with probability a(x).

Using  $\frac{p(\mathbf{x})}{\hat{p}(\mathbf{x})}$  in  $a(\mathbf{x})$  allows sampling from *P*.



### SAMPLING FROM A GENERATIVE MODEL Rejection Sampling

To sample a point from the refined distribution  $\tilde{P}$ :

- Sample  $z \sim Q$ .
- Compute x = G(z).
- Accept *x* with probability a(x).

It defines a new distribution  $\tilde{P}$ .



### SAMPLING FROM A GENERATIVE MODEL Rejection Sampling

To sample a point from the refined distribution  $\tilde{P}$ :

- Sample  $z \sim Q$ .
- Compute x = G(z).
- Accept *x* with probability a(x).



 $\mathbb{E}_{\widehat{P}}\left[a(\boldsymbol{x})\right].$ 



### SAMPLING FROM A GENERATIVE MODEL

**REJECTION SAMPLING IN HIGH DIMENSION** 

## BUDGETED REJECTION SAMPLING

TUNING THE ACCEPTANCE RATE

#### Definition 4.1 (Discriminator Rejection Sampling (DRS) - Azadi et al. [1])

*Let*  $\gamma \in \mathbb{R}$ *, the acceptance probability is:* 

$$a_{\mathrm{DRS}}(\mathbf{x}) = rac{r(\mathbf{x})}{r(\mathbf{x})\left(1 - e^{\gamma}\right) + Me^{\gamma}}.$$

*If*  $\gamma < 0$ *, then the acceptance rate increases.* 

## BUDGETED REJECTION SAMPLING

TUNING THE ACCEPTANCE RATE

#### Definition 4.1 (Discriminator Rejection Sampling (DRS) - Azadi et al. [1])

*Let*  $\gamma \in \mathbb{R}$ *, the acceptance probability is:* 

$$a_{\mathrm{DRS}}(\mathbf{x}) = rac{r(\mathbf{x})}{r(\mathbf{x})\left(1 - e^{\gamma}\right) + Me^{\gamma}}.$$

*If*  $\gamma < 0$ *, then the acceptance rate increases.*
### SECTION 3: IMPROVING CONTRIBUTIONS

## How can we improve the quality and diversity of a pre-trained generative models?

- ► How can we apply Rejection Sampling with a limited budget?
- How does it improve Precision and Recall?

### SECTION 3: IMPROVING CONTRIBUTIONS

### How can we improve the quality and diversity of a pre-trained generative models?

- ► How can we apply Rejection Sampling with a limited budget?
- ► How does it improve Precision and Recall?

#### List of contributions:

- We introduce the Optimal Budgeted Rejection Sampling (OBRS).
- We show how OBRS improves the Precision the learned distribution.
- ▶ We show that training a generative model with OBRS improves the Recall.

Traditionally, the goal is:

 $\min_{G} \quad \mathcal{D}_{f}(P \| \widehat{P}_{G})$ 



With a given  $\hat{P}_G$ , our goal is:

 $\min_{a} \quad \mathcal{D}_{f}(P \| \widetilde{P}_{a})$ 



With a given  $\hat{P}_G$ , our goal is:

$$\min_{a} \quad \mathcal{D}_{f}(P \| \widetilde{P}_{a}) \\
\text{s.t.} \begin{cases} \mathbb{E}_{\widehat{P}} \left[ a(\mathbf{x}) \right] \geq 1/K, \\ \forall \mathbf{x} \in \mathcal{X}, \ 0 \leq a(\mathbf{x}) \leq 1. \end{cases}$$
(17)



#### **Theorem 4.2 (Optimal Acceptance Function)**

*For a sampling budget*  $K \ge 1$  *and finite* X*, the solution is,* 

$$a_{\text{OBRS}}(\mathbf{x}) = \min\left(\frac{p(\mathbf{x})}{\widehat{p}(\mathbf{x})}\frac{c_K}{M}, 1\right),\tag{18}$$

where  $c_K \geq 1$  is such that  $\mathbb{E}_{\boldsymbol{x} \sim \widehat{p}}[a_{\text{OBRS}}(\boldsymbol{x})] = 1/K$ .

#### **Theorem 4.2 (Optimal Acceptance Function)**

*For a sampling budget*  $K \ge 1$  *and finite* X*, the solution is,* 

$$a_{\text{OBRS}}(\boldsymbol{x}) = \min\left(\frac{p(\boldsymbol{x})}{\widehat{p}(\boldsymbol{x})}\frac{c_K}{M}, 1\right),\tag{18}$$

where  $c_K \geq 1$  is such that  $\mathbb{E}_{\mathbf{x} \sim \hat{p}}[a_{\text{OBRS}}(\mathbf{x})] = 1/K$ .

# **Theorem 4.2 (Optimal Acceptance Function)** For a sampling budget $K \ge 1$ and finite X, the solution is, $a_{\text{OBRS}}(\mathbf{x}) = \min$ (18)where $c_K \geq 1$ is such that $\mathbb{E}_{x \sim \hat{p}}[a_{OBRS}(x)] = 1/K$ . Does not depend on f. Same acceptance function to improve Precision and Recall.

#### **Proposition 4.3 (Precision and Recall Improvement)**

# IMPROVING PRECISION AND RECALL IN PRACTICE



1/K	FID	Р	R
0.25	1.57	78.48	86.73
0.50	1.58	78.23	86.05
0.75	1.77	77.94	86.54
1	1.97	77.91	86.62

## Diffusion Model on CIFAR-10

## GAN on CelebA

# TRAINING GENERATIVE MODELS WITH OBRS OBJECTIVE

With a given  $\hat{P}_G$ , our goal is:

 $\min_{a} \quad \mathcal{D}_{f}(P \| \widetilde{P}_{a})$ 



# TRAINING GENERATIVE MODELS WITH OBRS OBJECTIVE

We can train *G* to optimize:

 $\min_{G} \min_{a} \quad \mathcal{D}_{f}(P \| \widetilde{P}_{a,G})$ 



# TRAINING GENERATIVE MODELS WITH OBRS OBJECTIVE

We can train *G* to optimize:

 $\min_{G} \min_{a} \quad \mathcal{D}_{f}(P \| \widetilde{P}_{a,G})$ 



# TRAINING WITH OBRS



# TRAINING WITH OBRS

Dataset	Method	FID	Р	R
CelebA	Hinge Loss	9.33	80.23	57.78
64  imes 64	Tw/OBRS	3.74	74.40	65.15
ImageNet	Hinge Loss	12.18	27.75	34.33
$128 \times 128$	Tw/OBRS	11.65	26.84	46.16

Training with OBRS increases the Recall.

- Quality and diversity are two important aspects of generative models.
- Generally at odds with each other, but can be balanced.
- ► The trade-off must be optimized.

- Quality and diversity are two important aspects of generative models.
- Generally at odds with each other, but can be balanced.
- ► The trade-off must be optimized.

### Optimizing the trade-off

- Connecting PR-Curves and *f*-divergences
  - Focusing on the AUC
  - Building a symmetric PR-Divergence
- Minimizing the PR-Divergence
  - Apply the method to Diffusion Models
  - Investigate other training strategies

- Quality and diversity are two important aspects of generative models.
- Generally at odds with each other, but can be balanced.
- ► The trade-off must be optimized.

### Optimizing the trade-off

- Connecting PR-Curves and *f*-divergences
  - Focusing on the AUC
  - Building a symmetric PR-Divergence
- Minimizing the PR-Divergence
  - Apply the method to Diffusion Models
  - Investigate other training strategies

### Improving the quality and diversity

- Optimal Budgeted Rejection Sampling
  - Investigate the density ratio estimation error
  - Apply the method during denoising in Diffusion Models
- Training to improve diversity
  - Build formal proofs

### BROADER PERSPECTIVES FAIRNESS IN GENERATIVE MODELS

There are no metrics to distinguish between *fair* and *unfair* generations in terms of Precision and Recall.



### CONCLUSION

## Questions

### **REFERENCES I**

- [1] Azadi, S., Olsson, C., Darrell, T., Goodfellow, I., and Odena, A. (2019). Discriminator Rejection Sampling. arXiv:1810.06758 [cs, stat].
- [2] Brock, A., Donahue, J., and Simonyan, K. (2019). Large Scale GAN Training for High Fidelity Natural Image Synthesis. arXiv:1809.11096 [cs, stat].
- [3] Bronnec, F. L., Verine, A., Negrevergne, B., Chevaleyre, Y., and Allauzen, A. (2024). Exploring Precision and Recall to assess the quality and diversity of LLMs. *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*. arXiv:2402.10693 [cs].
- [4] Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative Adversarial Networks. In 27th Conference on Neural Information Processing Systems (NeurIPS 2014). arXiv: 1406.2661.
- [5] Grover, A., Dhar, M., and Ermon, S. (2018). Flow-GAN: Combining Maximum Likelihood and Adversarial Learning in Generative Models. arXiv:1705.08868 [cs, stat].
- [6] Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., and Aila, T. (2020). Analyzing and Improving the Image Quality of StyleGAN. arXiv:1912.04958 [cs, eess, stat].
- [7] Kingma, D. P. and Dhariwal, P. (2018). Glow: Generative Flow with Invertible 1x1 Convolutions. In 32nd Conference on Neural Information Processing Systems (NeurIPS 2018), Montréal, Canada., volume 31.

### **REFERENCES II**

- [8] Kynkäänniemi, T., Karras, T., Laine, S., Lehtinen, J., and Aila, T. (2019). Improved Precision and Recall Metric for Assessing Generative Models. In 33rd Conference on Neural Information Processing Systems (NeurIPS 2019), Vancouver, Canada. arXiv: 1904.06991.
- [9] Nguyen, X., Wainwright, M. J., and Jordan, M. I. (2009). On surrogate loss functions and \$f\$-divergences. *The Annals of Statistics*, 37(2). arXiv:math/0510521.
- [10] Nowozin, S., Cseke, B., and Tomioka, R. (2016). f-GAN: Training Generative Neural Samplers using Variational Divergence Minimization. arXiv:1606.00709 [cs, stat].
- [11] Sajjadi, M. S. M., Bachem, O., Lucic, M., Bousquet, O., and Gelly, S. (2018). Assessing Generative Models via Precision and Recall. In 32nd Conference on Neural Information Processing Systems (NeurIPS 2018), Montréal, Canada. arXiv: 1806.00035.
- [12] Simon, L., Webster, R., and Rabin, J. (2019). Revisiting precision recall definition for generative modeling. In *Proceedings of the 36th International Conference on Machine Learning*, pages 5799–5808. PMLR. ISSN: 2640-3498.
- [13] Siry, R., Webster, R., Simon, L., and Rabin, J. (2023). On the Theoretical Equivalence of Several Trade-Off Curves Assessing Statistical Proximity. *Journal of Machine Learning Research*, 24.
- [14] Yann LeCun, Corinna Cortes, and Burges, C. (2010). MNIST handwritten digit database. ATT Labs, 2.

#### **Temporary page!**

LATEX was unable to guess the total number of pages correctly. As there was some unprocessed data that should have been added to the final page this extra page has been added to receive it. If you rerun the document (without altering it) this surplus page will go away, because LATEX now knows how many pages to expect for this document.