

# Fair machine learning

## Lecture 4

Changho Suh  
EE, KAIST

Aug. 4, 2021

# **A generative model: A hinge-loss-based GAN**

**Reading: TN4**

# Recap: Fair classifiers

---

1. Explore fairness measures in fair classifiers:  
DDP and DEO
2. Study an MI-based fair classifier which offers a good tradeoff while suffering from training instability.
3. Investigate another fair classifier based on KDE, which addresses the training instability issue.

# Outline of Lectures 4/5/6

---

Move onto a different context: **Fair generative model.**

1. Study a generative model which forms the basis of a fair generative model to be explored in depth:

## **A hinge-loss-based GAN**

2. Make a connection with a well-known statistical measure prevalent in information theory:

## **Total Variation Distance (TVD)**

3. Investigate a TVD-based fair generative model.
4. Discuss a couple of other relevant issues.

# Focus of Lecture 4

---

Move onto a different context: **Fair generative model.**

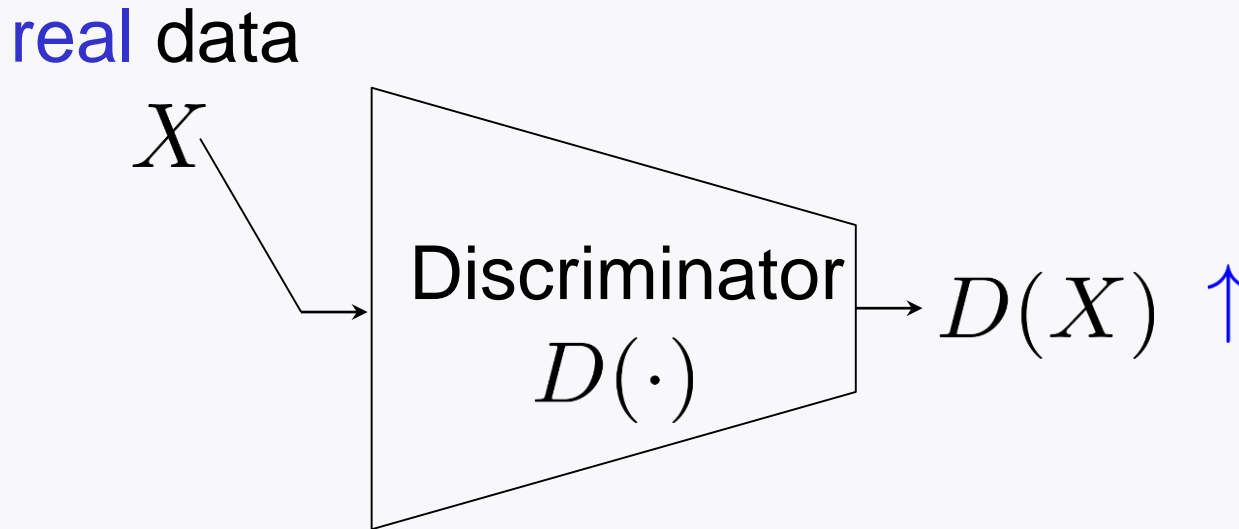
1. Study a generative model which forms the basis of a fair generative model to be explored in depth:

## **A hinge-loss-based GAN**

2. Make a connection with a well-known statistical measure prevalent in information theory:

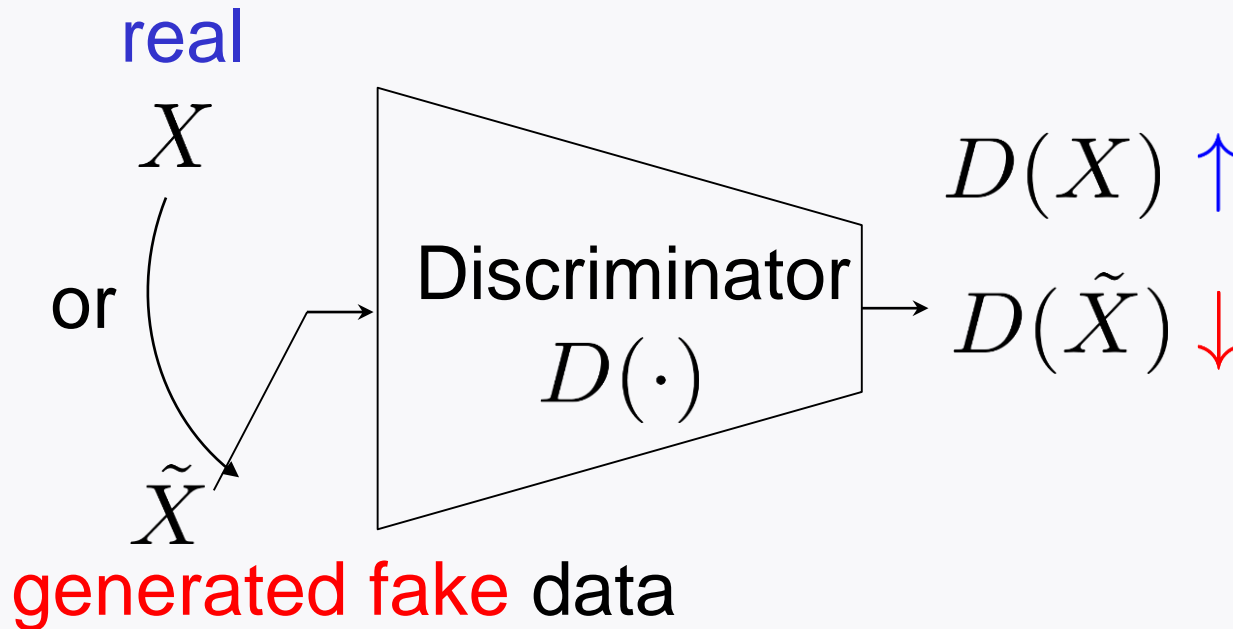
## **Total Variation Distance (TVD)**

3. Investigate a TVD-based fair generative model.
4. Discuss a couple of other relevant issues.



**Role:** Discriminate real from generated fake samples

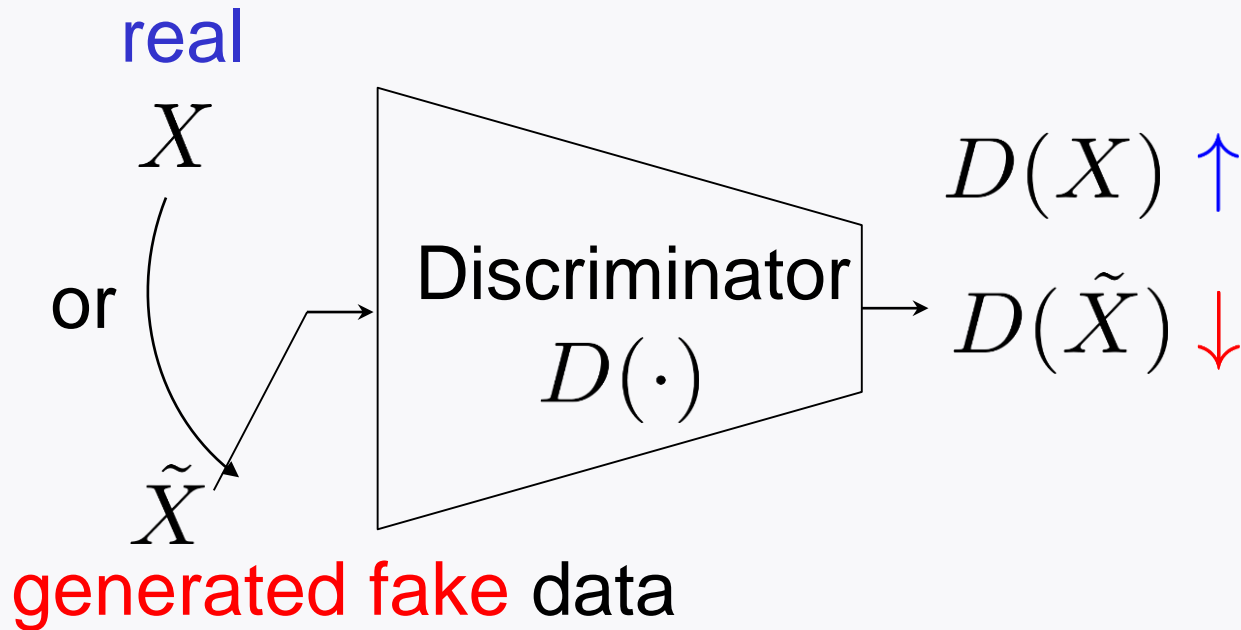
Intend to yield a **large**  $D(\cdot)$  if the input is **real** data;



**Role:** Discriminate real from generated fake samples

Intend to yield a **large**  $D(\cdot)$  if the input is **real** data;  
a **small**  $D(\cdot)$  for **generated** data.

# A reasonable interpretation on $D(\cdot)$

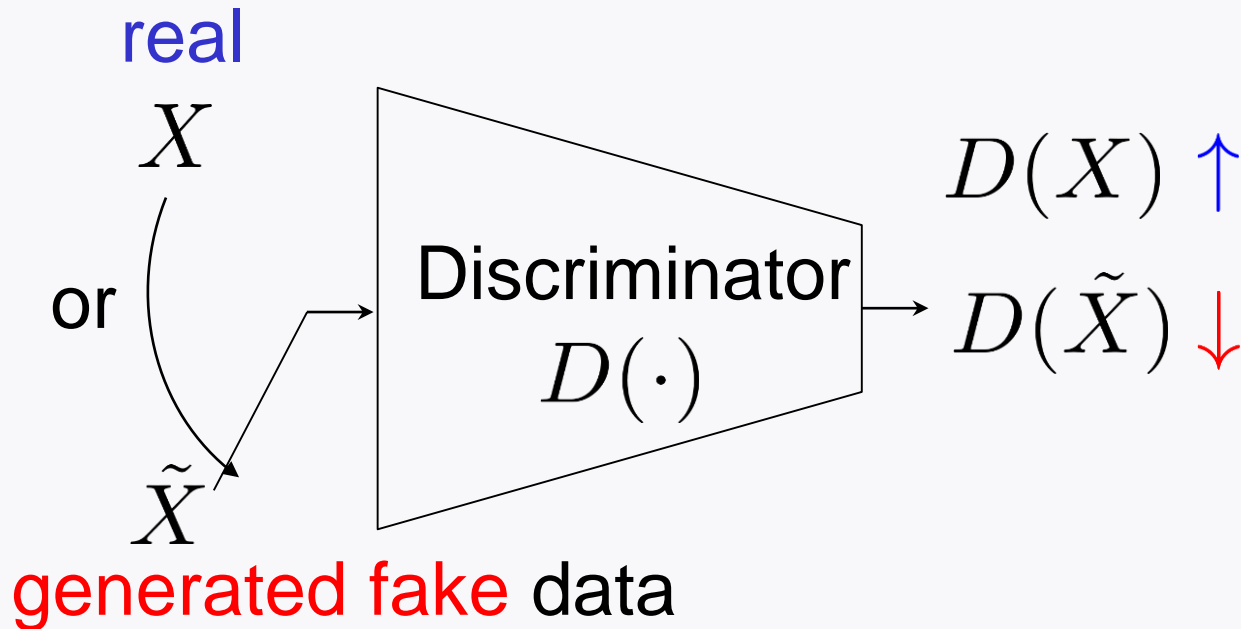


Probability of the input being **real**:

$$D(\cdot) = \mathbb{P}((\cdot) = \text{real})$$



# A reasonable interpretation on $D(\cdot)$

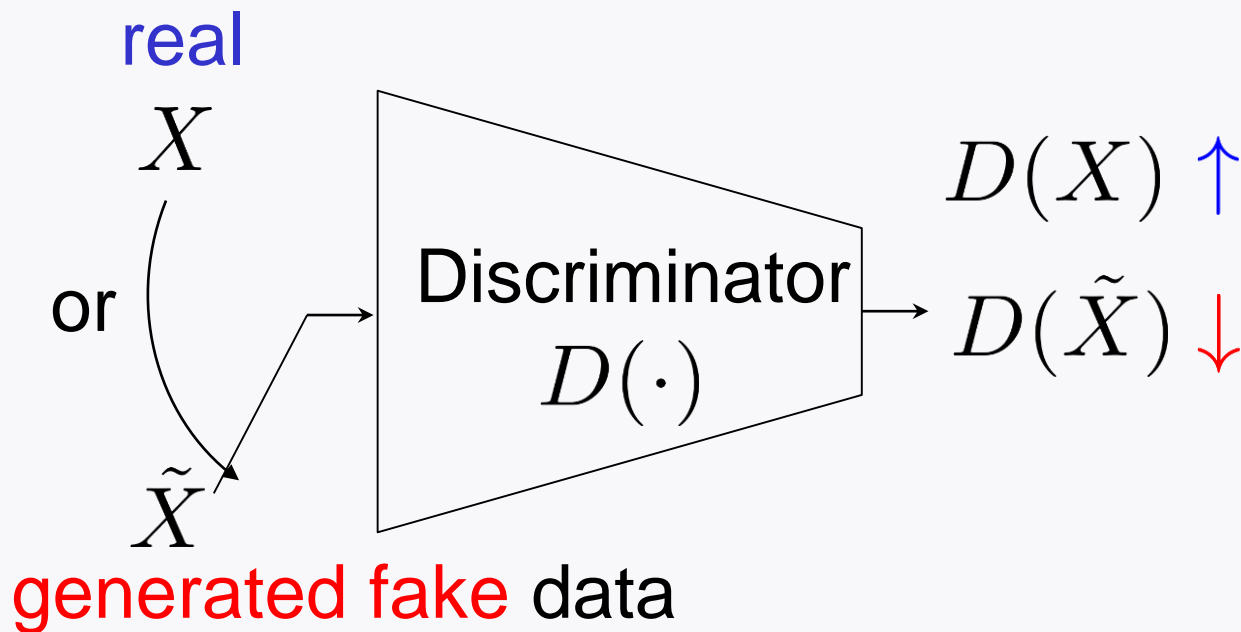


Probability of the input being **real**:

$$\uparrow D(X) = \mathbb{P}(X = \text{real}) = 1$$

$$\downarrow D(\tilde{X}) = \mathbb{P}(\tilde{X} = \text{real}) = 0$$

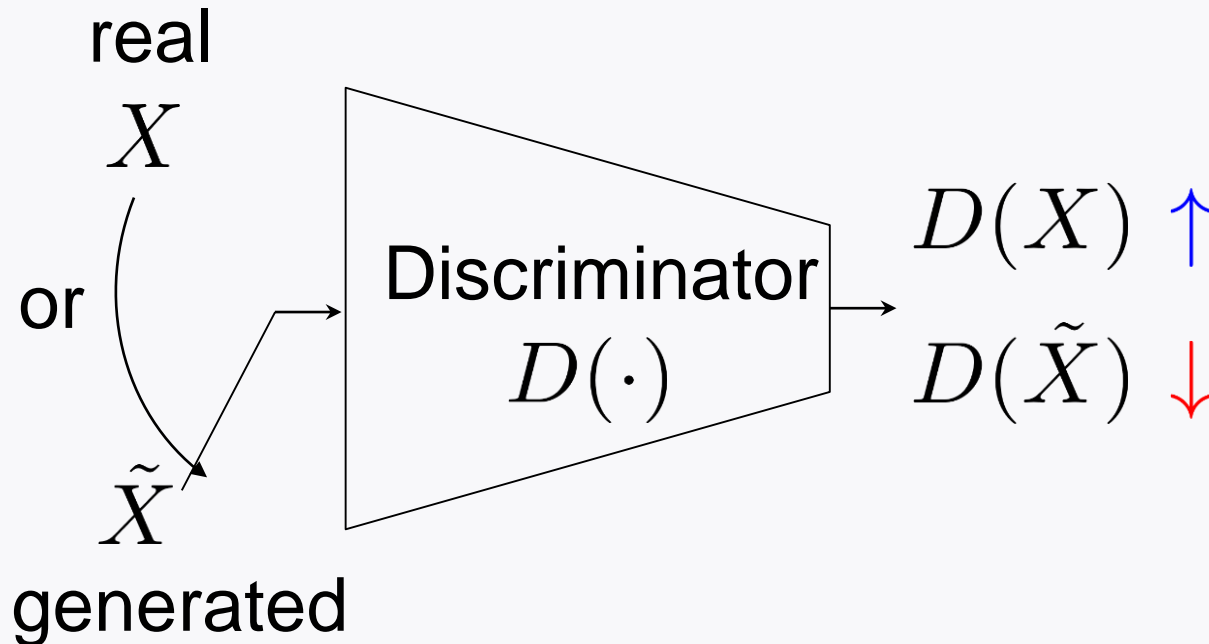
# Goal



Maximize:  $D(X)$  &  $1 - D(\tilde{X})$

A natural optimization:

$$\max_D \mathbb{E}_{\mathbb{P}_{\text{real}}} [D(X)] + \mathbb{E}_{\mathbb{P}_G} [1 - D(\tilde{X})]$$



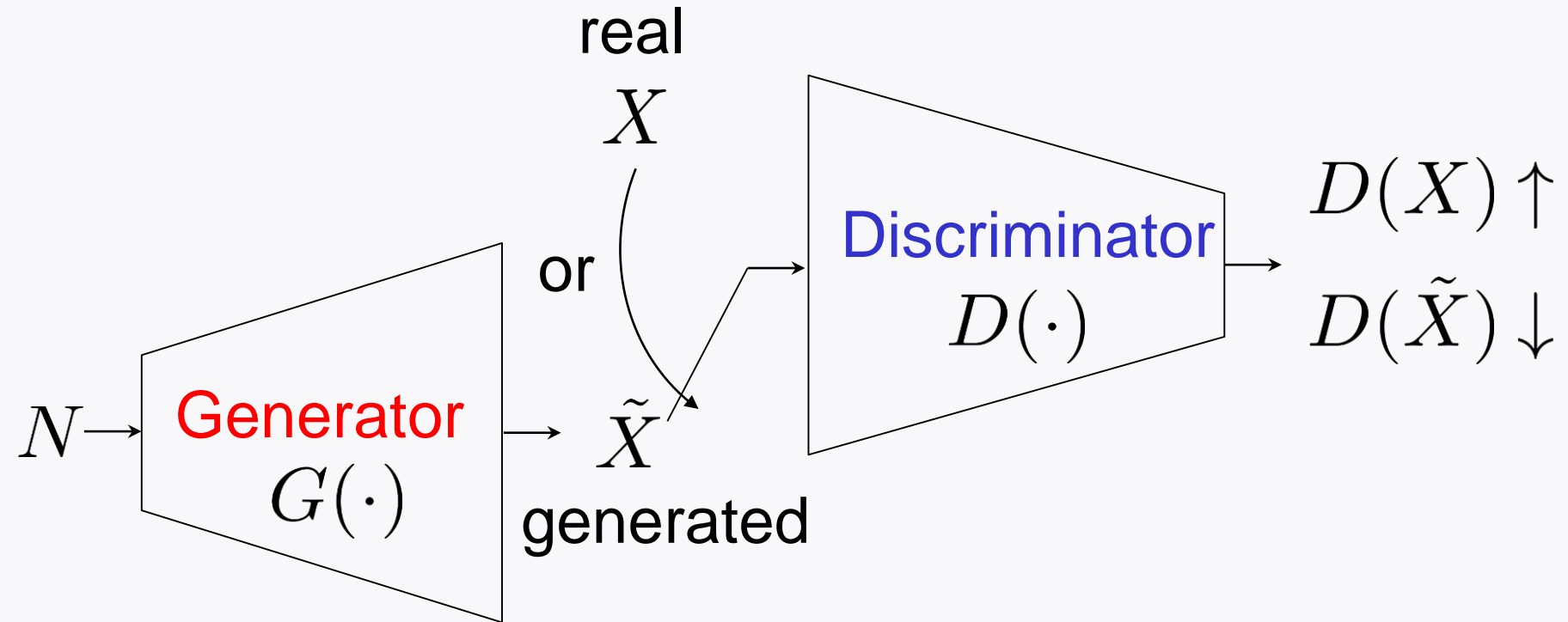
Goodfellow employed **log loss** instead:

$$\max_D \mathbb{E}_{\mathbb{P}_{\text{real}}} [\log D(X)] + \mathbb{E}_{\mathbb{P}_G} [\log(1 - D(\tilde{X}))]$$

Goodfellow mentioned: It is inspired by [Gutmann et al. AISTATS10] in his NeurIPS16's tutorial.

# A two-player game

Goodfellow et al. NeurIPS14



**Discriminator:**  $\max_D \mathbb{E}_{\mathbb{P}_{\text{real}}} [\log D(X)] + \mathbb{E}_{\mathbb{P}_G} [\log(1 - D(\tilde{X}))]$

**Generator:**  $\min_G \mathbb{E}_{\mathbb{P}_{\text{real}}} [\log D(X)] + \mathbb{E}_{\mathbb{P}_G} [\log(1 - D(\tilde{X}))]$

## GAN:

$$D: \max_D \mathbb{E}_{\mathbb{P}_{\text{real}}} [\log D(X)] + \mathbb{E}_{\mathbb{P}_G} [\log(1 - D(\tilde{X}))]$$

$$G: \min_G \mathbb{E}_{\mathbb{P}_{\text{real}}} [\log D(X)] + \mathbb{E}_{\mathbb{P}_G} [\log(1 - D(\tilde{X}))]$$

Lim et al. employed different loss functions:

**hinge** loss for  $D$ ; **linear** loss for  $G$ .

## Hinge-loss-based GAN:

$$D: \max_D \mathbb{E}_{\mathbb{P}_{\text{real}}} [\min(0, -1 + D(X))] + \mathbb{E}_{\mathbb{P}_G} [\min(0, -1 - D(\tilde{X}))]$$

$$G: \min_G \mathbb{E}_{\mathbb{P}_{\text{real}}} [D(X)] - \mathbb{E}_{\mathbb{P}_G} [D(\tilde{X})]$$

# Hinge-loss-based GAN

Lim et al. arxiv17

$$D: \max_D \mathbb{E}_{\mathbb{P}_{\text{real}}} [\min(0, -1 + D(X))] + \mathbb{E}_{\mathbb{P}_G} [\min(0, -1 - D(\tilde{X}))]$$
$$G: \min_G \mathbb{E}_{\mathbb{P}_{\text{real}}} [D(X)] - \mathbb{E}_{\mathbb{P}_G} [D(\tilde{X})]$$

**Turns out:** Its advanced version offers the state-of-the-art performance in many of the real datasets.



Fake samples generated by its advanced version (CelebA)

$$D: \max_D \mathbb{E}_{\mathbb{P}_{\text{real}}} [\min(0, -1 + D(X))] + \mathbb{E}_{\mathbb{P}_G} [\min(0, -1 - D(\tilde{X}))]$$
$$G: \min_G \mathbb{E}_{\mathbb{P}_{\text{real}}} [D(X)] - \mathbb{E}_{\mathbb{P}_G} [D(\tilde{X})]$$

Looks like a heuristic choice on the loss functions.

But it turns out this yields an interesting connection w/:

## **Total Variation Distance (TVD)**

Also the connection gives insights into developing a fair generative model.

$$D: \max_D \mathbb{E}_{\mathbb{P}_{\text{real}}} [\min(0, -1 + D(X))] + \mathbb{E}_{\mathbb{P}_G} [\min(0, -1 - D(\tilde{X}))]$$
$$G: \min_G \mathbb{E}_{\mathbb{P}_{\text{real}}} [D(X)] - \mathbb{E}_{\mathbb{P}_G} [D(\tilde{X})]$$

One way to solve the optimization:

First solve  $D^*$

Then plug  $D^*$  into Generator's objective to solve  $G^*$



$$D: \max_D \mathbb{E}_{\mathbb{P}_{\text{real}}} [\min(0, -1 + D(X))] + \mathbb{E}_{\mathbb{P}_G} [\min(0, -1 - D(\tilde{X}))]$$
$$G: \min_G \mathbb{E}_{\mathbb{P}_{\text{real}}} [D(X)] - \mathbb{E}_{\mathbb{P}_G} [D(\tilde{X})]$$

**Theorem (Tan et al.):**

$$D^*(x) = \text{sign}(\mathbb{P}_{\text{real}}(x) - \mathbb{P}_G(x)) \quad \forall x$$

$$G^* = \arg \min_G \sum_{x \in \mathcal{X} \cup \tilde{\mathcal{X}}} |\mathbb{P}_{\text{real}}(x) - \mathbb{P}_G(x)|$$

set of real samples  $\nearrow$   $x \in \mathcal{X} \cup \tilde{\mathcal{X}}$   $\nwarrow$  set of generated samples

$$\mathbb{P}_{\text{real}}(x) = 0 \quad \text{if } x \in \tilde{\mathcal{X}} \setminus \mathcal{X}$$

$$\mathbb{P}_G(x) = 0 \quad \text{if } x \in \mathcal{X} \setminus \tilde{\mathcal{X}}$$

$$D: \max_D \mathbb{E}_{\mathbb{P}_{\text{real}}} [\min(0, -1 + D(X))] + \mathbb{E}_{\mathbb{P}_G} [\min(0, -1 - D(\tilde{X}))]$$
$$G: \min_G \mathbb{E}_{\mathbb{P}_{\text{real}}} [D(X)] - \mathbb{E}_{\mathbb{P}_G} [D(\tilde{X})]$$

**Theorem (Tan et al.):**

$$D^*(x) = \text{sign}(\mathbb{P}_{\text{real}}(x) - \mathbb{P}_G(x)) \quad \forall x$$

$$G^* = \arg \min_G \sum_{x \in \mathcal{X} \cup \tilde{\mathcal{X}}} |\mathbb{P}_{\text{real}}(x) - \mathbb{P}_G(x)|$$

$$\mathbb{P}_{\text{real}}(x) = 0 \quad \text{if } x \in \tilde{\mathcal{X}} \setminus \mathcal{X}$$

$$\mathbb{P}_G(x) = 0 \quad \text{if } x \in \mathcal{X} \setminus \tilde{\mathcal{X}}$$

$$=: 2 \cdot \text{TV}(\mathbb{P}_{\text{real}}, \mathbb{P}_G)$$

**Proof:**  $D^*(x) = \text{sign}(\mathbb{P}_{\text{real}}(x) - \mathbb{P}_G(x)) \quad \forall x$

$$D: \max_D \underbrace{\mathbb{E}_{\mathbb{P}_{\text{real}}} [\min(0, -1 + D(X))] + \mathbb{E}_{\mathbb{P}_G} [\min(0, -1 - D(\tilde{X}))]}_{\sum_{x \in \mathcal{X} \cup \tilde{\mathcal{X}}} \frac{\mathbb{P}_{\text{real}}(x) \min(0, -1 + D(x)) + \mathbb{P}_G(x) \min(0, -1 - D(x))}{2}}$$

Using Lemma, we get:

$$D^*(x) = \text{sign}(\mathbb{P}_{\text{real}}(x) - \mathbb{P}_G(x)) \quad \forall x$$

**Lemma:**

$$f(t) = \alpha \cdot \min\{0, -1 + t\} + \beta \cdot \min\{0, -1 - t\}$$

For given  $\alpha, \beta \geq 0$ :  $t^* = \arg \max_t f(t) = \text{sign}(\alpha - \beta)$

**Proof:**  $G^* = \arg \min_G \text{TV}(\mathbb{P}_{\text{real}}, \mathbb{P}_G)$

$$D^*(x) = \text{sign}(\mathbb{P}_{\text{real}}(x) - \mathbb{P}_G(x)) \quad \forall x$$

$$\mathbf{G}: \min_G \mathbb{E}_{\mathbb{P}_{\text{real}}} [D(X)] - \mathbb{E}_{\mathbb{P}_G} [D(\tilde{X})]$$

$$\begin{aligned} & \mathbb{E}_{\mathbb{P}_{\text{real}}} [D^*(X)] - \mathbb{E}_{\mathbb{P}_G} [D^*(\tilde{X})] \\ &= \sum_{x \in \mathcal{X} \cup \tilde{\mathcal{X}}} \mathbb{P}_{\text{real}}(x) D^*(x) - \sum_{x \in \mathcal{X} \cup \tilde{\mathcal{X}}} \mathbb{P}_G(x) D^*(x) \\ &= \sum_{x \in \mathcal{X} \cup \tilde{\mathcal{X}}} (\mathbb{P}_{\text{real}}(x) - \mathbb{P}_G(x)) D^*(x) \\ &= \sum_{x \in \mathcal{X} \cup \tilde{\mathcal{X}}} |\mathbb{P}_{\text{real}}(x) - \mathbb{P}_G(x)| = 2 \cdot \text{TV}(\mathbb{P}_{\text{real}}, \mathbb{P}_G) \end{aligned}$$

# Look ahead

---

Investigate a TVD-based fair generative model.

# Reference

---

- [1] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. *Advances in Neural Information Processing Systems 27 (NeurIPS)*, 2014.
- [2] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. *AISTATS* 2010.
- [3] J. H. Lim and J. C. Ye. Geometric gan. *arXiv:1705.02894*, 2017.
- [4] A. Brock, J. Donahue, and K. Simonyan. Large scale GAN training for high fidelity natural image synthesis. *International Conference on Learning Representations*, 2019.
- [5] Z. Tan, Y. Song, and Z. Ou. Calibrated adversarial algorithms for generative modelling. *Stat*, 2019.
- [6] L. Tierney. Introduction to general state-space markov chain theory. *Markov chain Monte Carlo in practice*, 1996.