

Fair machine learning

Lecture 5

Changho Suh
EE, KAIST

Aug. 4, 2021

A fair generative model using total variation distance

Reading: TN5

Recap: Hinge-loss-based GAN

Discriminator (**hinge** loss):

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

Generator (**linear** loss):

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

Made a connection w/ **total variation distance**:

$$\min_G \mathbb{E}_{\mathbb{P}_{\text{real}}} [D^*(X)] - \mathbb{E}_{\mathbb{P}_G} [D^*(\tilde{X})] = \min_G 2 \cdot \text{TV}(\mathbb{P}_{\text{real}}, \mathbb{P}_G)$$

Claimed: This connection gives insights into a fair generative model.

Outline of Lecture 5

Explore a **TVD**-based fair generative model.

1. Introduce a TVD-based regularization term that promotes **fair sample generation (class-balanced samples)**
2. Formulate a corresponding optimization.
3. Translate it into an **implementable** optimization that employs hinge loss.
4. Discuss experimental results.

A regularization term for fair sample gen.

Generator in hinge-loss-based GAN:

$$\min_G \text{TV}(\mathbb{P}_{\text{real}}, \mathbb{P}_G)$$

What we want for fair sample generation:

$$\mathbb{P}_G \approx \mathbb{P}_{\text{ref}} \leftarrow \text{a reference distribution respecting fair sample gen.}$$

An issue arises in satisfying this:

The given real data does not necessarily satisfy \mathbb{P}_{ref} .

Note: Interested scenarios: \mathbb{P}_{real} biased

A regularization term for fair sample gen.

Generator in hinge-loss-based GAN:

$$\min_G \text{TV}(\mathbb{P}_{\text{real}}, \mathbb{P}_G)$$

What we want for fair sample generation:

$\mathbb{P}_G \approx \mathbb{P}_{\text{ref}}$ ← a **reference** distribution
respecting fair sample gen.

A natural way to satisfy this:

Introduce a *new yet small* reference dataset
respecting \mathbb{P}_{ref} .

5-10% relative to original real data

A regularization term for fair sample gen.

Generator in hinge-loss-based GAN:

$$\min_G \text{TV}(\mathbb{P}_{\text{real}}, \mathbb{P}_G)$$

What we want for fair sample generation:

$\mathbb{P}_G \approx \mathbb{P}_{\text{ref}}$ ← a **reference** distribution
respecting fair sample gen.

A neutral regularization term:

$$\text{TV}(\mathbb{P}_{\text{ref}}, \mathbb{P}_G)$$

TVD-based optimization for a fair gen. model

[Um-Suh '21]:

$$\min_G (1 - \lambda) \cdot \text{TV}(\mathbb{P}_{\text{real}}, \mathbb{P}_G) + \lambda \cdot \text{TV}(\mathbb{P}_{\text{ref}}, \mathbb{P}_G)$$

Question:

How to solve the optimization?

Observation

[Um-Suh '21]:

$$\min_G (1 - \lambda) \cdot \text{TV}(\mathbb{P}_{\text{real}}, \mathbb{P}_G) + \lambda \cdot \text{TV}(\mathbb{P}_{\text{ref}}, \mathbb{P}_G)$$

Remember: $\text{TV}(\mathbb{P}_{\text{real}}, \mathbb{P}_G)$ was a consequence of evaluating Generator's objective at D^* , which was derived from:

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

Observation

[Um-Suh '21]:

$$\min_G (1 - \lambda) \cdot \text{TV}(\mathbb{P}_{\text{real}}, \mathbb{P}_G) + \lambda \cdot \text{TV}(\mathbb{P}_{\text{ref}}, \mathbb{P}_G)$$

Guess: $\text{TV}(\mathbb{P}_{\text{ref}}, \mathbb{P}_G)$ is a consequence of evaluating Generator's objective at another D_{ref}^* , which is derived from:

$$\max_{D_{\text{ref}}} \mathbb{E}_{\mathbb{P}_{\text{ref}}} [\min(0, -1 + D_{\text{ref}}(X_{\text{ref}}))] + \mathbb{E}_{\mathbb{P}_G} [\min(0, -1 - D_{\text{ref}}(\tilde{X}))]$$

Equivalence

$$\min_G (1 - \lambda) \cdot \text{TV}(\mathbb{P}_{\text{real}}, \mathbb{P}_G) + \lambda \cdot \text{TV}(\mathbb{P}_{\text{ref}}, \mathbb{P}_G)$$

Turns out: Equivalent to

$$\begin{aligned} & \max_D \mathbb{E}_{\mathbb{P}_{\text{real}}} [\min(0, -1 + D(X))] + \mathbb{E}_{\mathbb{P}_G} [\min(0, -1 - D(\tilde{X}))] \\ & \max_{D_{\text{ref}}} \mathbb{E}_{\mathbb{P}_{\text{ref}}} [\min(0, -1 + D_{\text{ref}}(X_{\text{ref}}))] + \mathbb{E}_{\mathbb{P}_G} [\min(0, -1 - D_{\text{ref}}(\tilde{X}))] \\ & \min_G (1 - \lambda) \left\{ \mathbb{E}_{\mathbb{P}_{\text{real}}} [D(X)] - \mathbb{E}_{\mathbb{P}_G} [D(\tilde{X})] \right\} + \lambda \left\{ \mathbb{E}_{\mathbb{P}_{\text{ref}}} [D_{\text{ref}}(X_{\text{ref}})] - \mathbb{E}_{\mathbb{P}_G} [D_{\text{ref}}(\tilde{X})] \right\} \end{aligned}$$

To prove this, need to show:

$$\begin{aligned} & \min_G (1 - \lambda) \left\{ \mathbb{E}_{\mathbb{P}_{\text{real}}} [D^*(X)] - \mathbb{E}_{\mathbb{P}_G} [D^*(\tilde{X})] \right\} + \lambda \left\{ \mathbb{E}_{\mathbb{P}_{\text{ref}}} [D_{\text{ref}}^*(X_{\text{ref}})] - \mathbb{E}_{\mathbb{P}_G} [D_{\text{ref}}^*(\tilde{X})] \right\} \\ & = \min_G 2(1 - \lambda) \text{TV}(\mathbb{P}_{\text{real}}, \mathbb{P}_G) + 2\lambda \text{TV}(\mathbb{P}_{\text{ref}}, \mathbb{P}_G) \end{aligned}$$

Proof of equivalence

$$\begin{aligned} & \max_D \mathbb{E}_{\mathbb{P}_{\text{real}}} [\min(0, -1 + D(X))] + \mathbb{E}_{\mathbb{P}_G} [\min(0, -1 - D(\tilde{X}))] \\ & \max_{D_{\text{ref}}} \mathbb{E}_{\mathbb{P}_{\text{ref}}} [\min(0, -1 + D_{\text{ref}}(X_{\text{ref}}))] + \mathbb{E}_{\mathbb{P}_G} [\min(0, -1 - D_{\text{ref}}(\tilde{X}))] \\ & \min_G (1 - \lambda) \left\{ \mathbb{E}_{\mathbb{P}_{\text{real}}} [D(X)] - \mathbb{E}_{\mathbb{P}_G} [D(\tilde{X})] \right\} + \lambda \left\{ \mathbb{E}_{\mathbb{P}_{\text{ref}}} [D_{\text{ref}}(X_{\text{ref}})] - \mathbb{E}_{\mathbb{P}_G} [D_{\text{ref}}(\tilde{X})] \right\} \end{aligned}$$

Using the technique based on the lemma introduced earlier, one can show:

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

$$D_{\text{ref}}^*(x) = \text{sign}(\mathbb{P}_{\text{ref}}(x) - \mathbb{P}_G(x)) \quad \forall x \in \mathcal{X}_{\text{ref}} \cup \tilde{\mathcal{X}}$$

Proof of equivalence

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

$$D_{\text{ref}}^*(x) = \text{sign}(\mathbb{P}_{\text{ref}}(x) - \mathbb{P}_G(x)) \quad \forall x \in \mathcal{X}_{\text{ref}} \cup \tilde{\mathcal{X}}$$

$$\min_G (1 - \lambda) \left\{ \mathbb{E}_{\mathbb{P}_{\text{real}}} [D(X)] - \mathbb{E}_{\mathbb{P}_G} [D(\tilde{X})] \right\} + \lambda \left\{ \mathbb{E}_{\mathbb{P}_{\text{ref}}} [D_{\text{ref}}(X_{\text{ref}})] - \mathbb{E}_{\mathbb{P}_G} [D_{\text{ref}}(\tilde{X})] \right\}$$

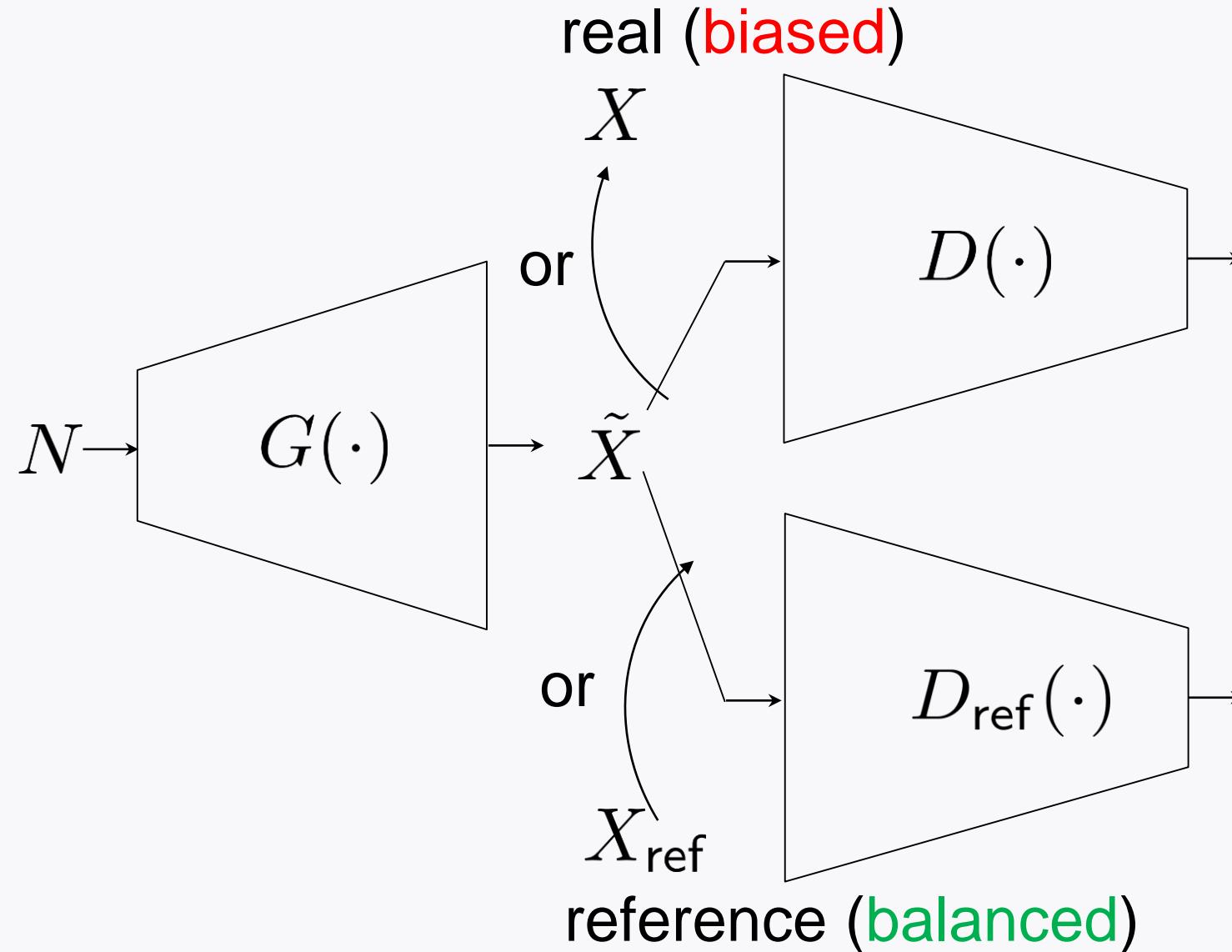
$$(1 - \lambda) \left\{ \mathbb{E}_{\mathbb{P}_{\text{real}}} [D^*(X)] - \mathbb{E}_{\mathbb{P}_G} [D^*(\tilde{X})] \right\} + \lambda \left\{ \mathbb{E}_{\mathbb{P}_{\text{ref}}} [D_{\text{ref}}^*(X_{\text{ref}})] - \mathbb{E}_{\mathbb{P}_G} [D_{\text{ref}}^*(\tilde{X})] \right\}$$

$$= (1 - \lambda) \sum_{x \in \mathcal{X} \cup \tilde{\mathcal{X}}} (\mathbb{P}_{\text{real}}(x) - \mathbb{P}_G(x)) D^*(x) + \lambda \sum_{x \in \mathcal{X}_{\text{ref}} \cup \tilde{\mathcal{X}}} (\mathbb{P}_{\text{ref}}(x) - \mathbb{P}_G(x)) D_{\text{ref}}^*(x)$$

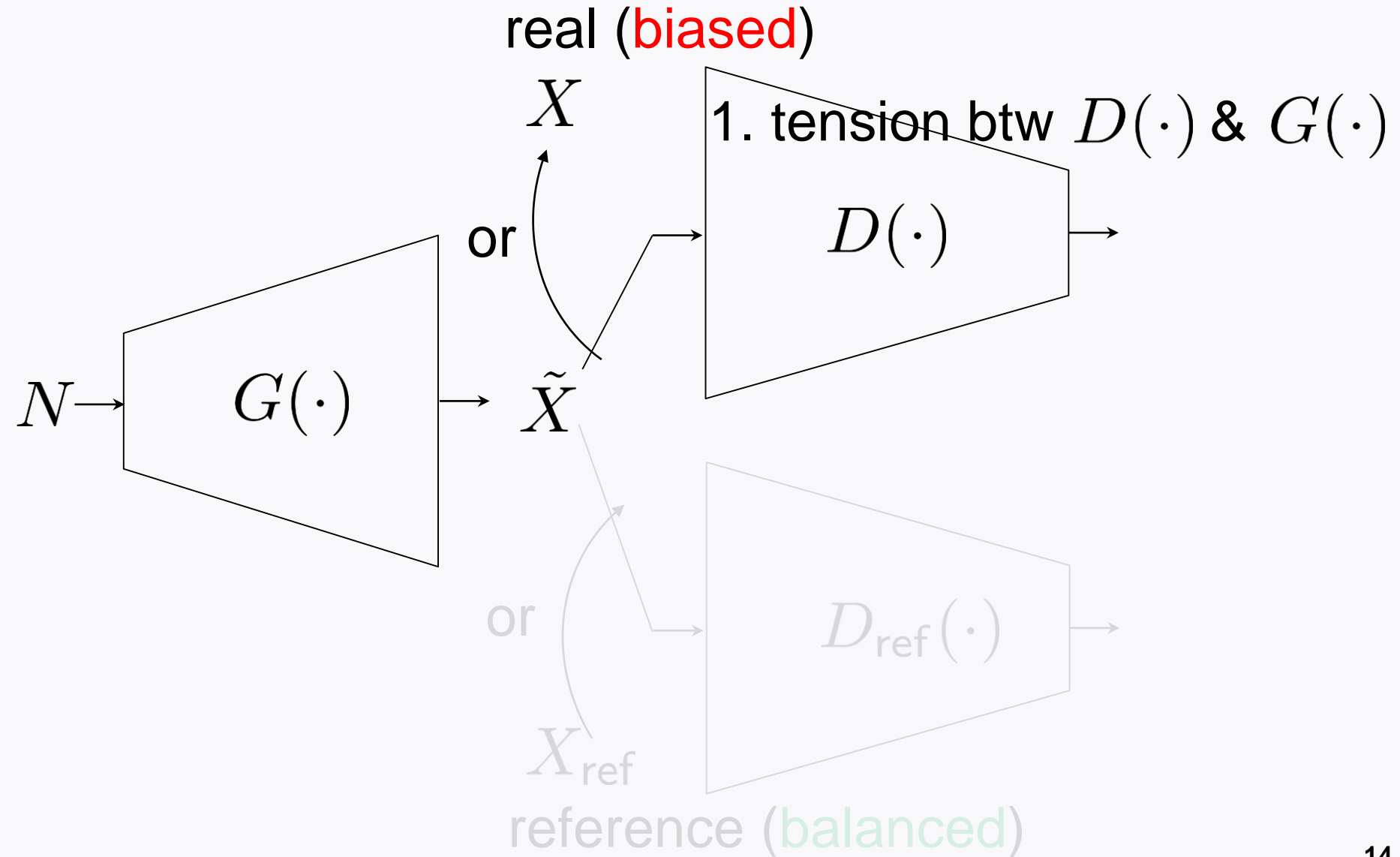
$$= (1 - \lambda) \sum_{x \in \mathcal{X} \cup \tilde{\mathcal{X}}} |\mathbb{P}_{\text{real}}(x) - \mathbb{P}_G(x)| + \lambda \sum_{x \in \mathcal{X}_{\text{ref}} \cup \tilde{\mathcal{X}}} |\mathbb{P}_{\text{ref}}(x) - \mathbb{P}_G(x)|$$

$$= 2(1 - \lambda) \text{TV}(\mathbb{P}_{\text{real}}, \mathbb{P}_G) + 2\lambda \text{TV}(\mathbb{P}_{\text{ref}}, \mathbb{P}_G)$$

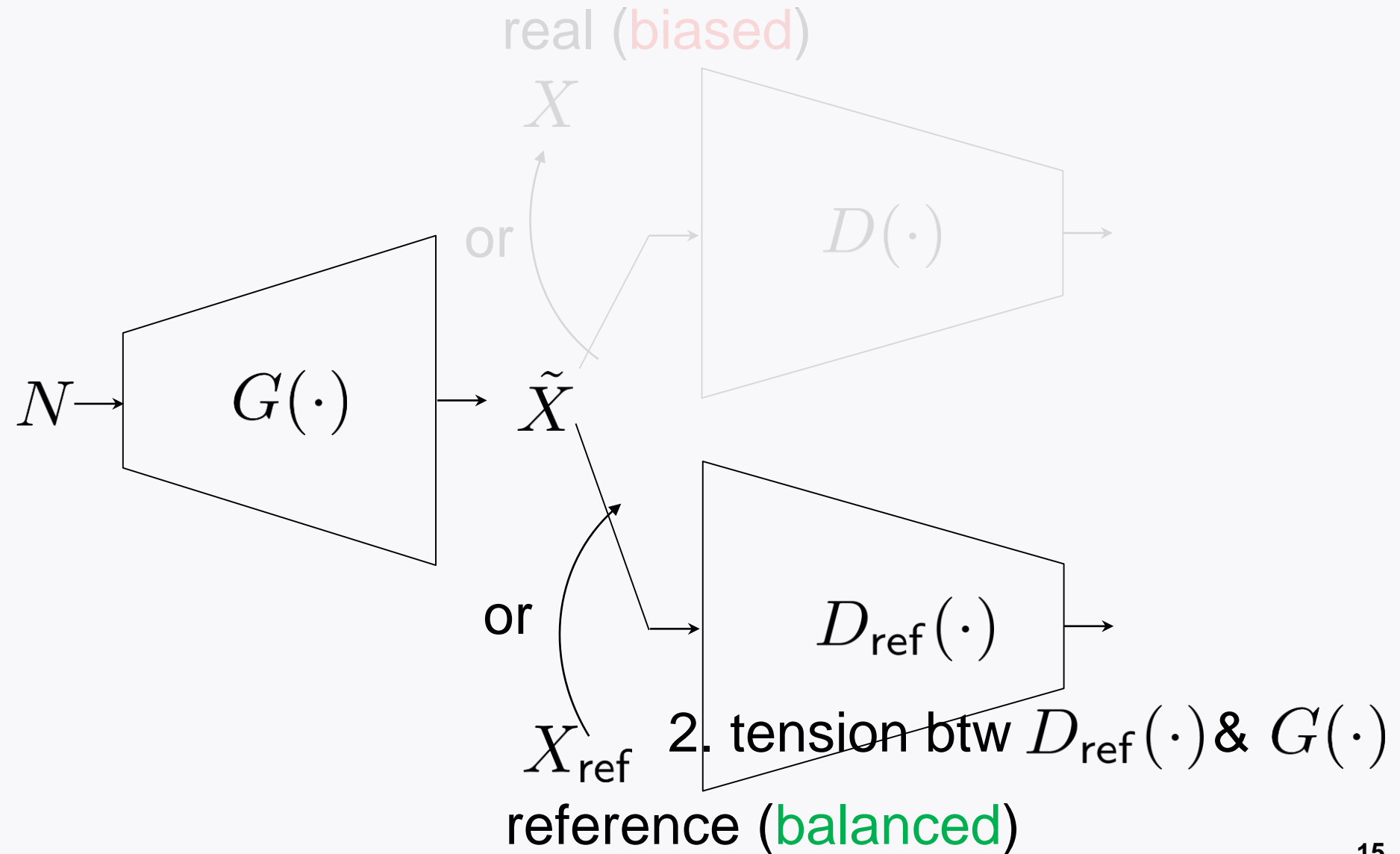
Architecture



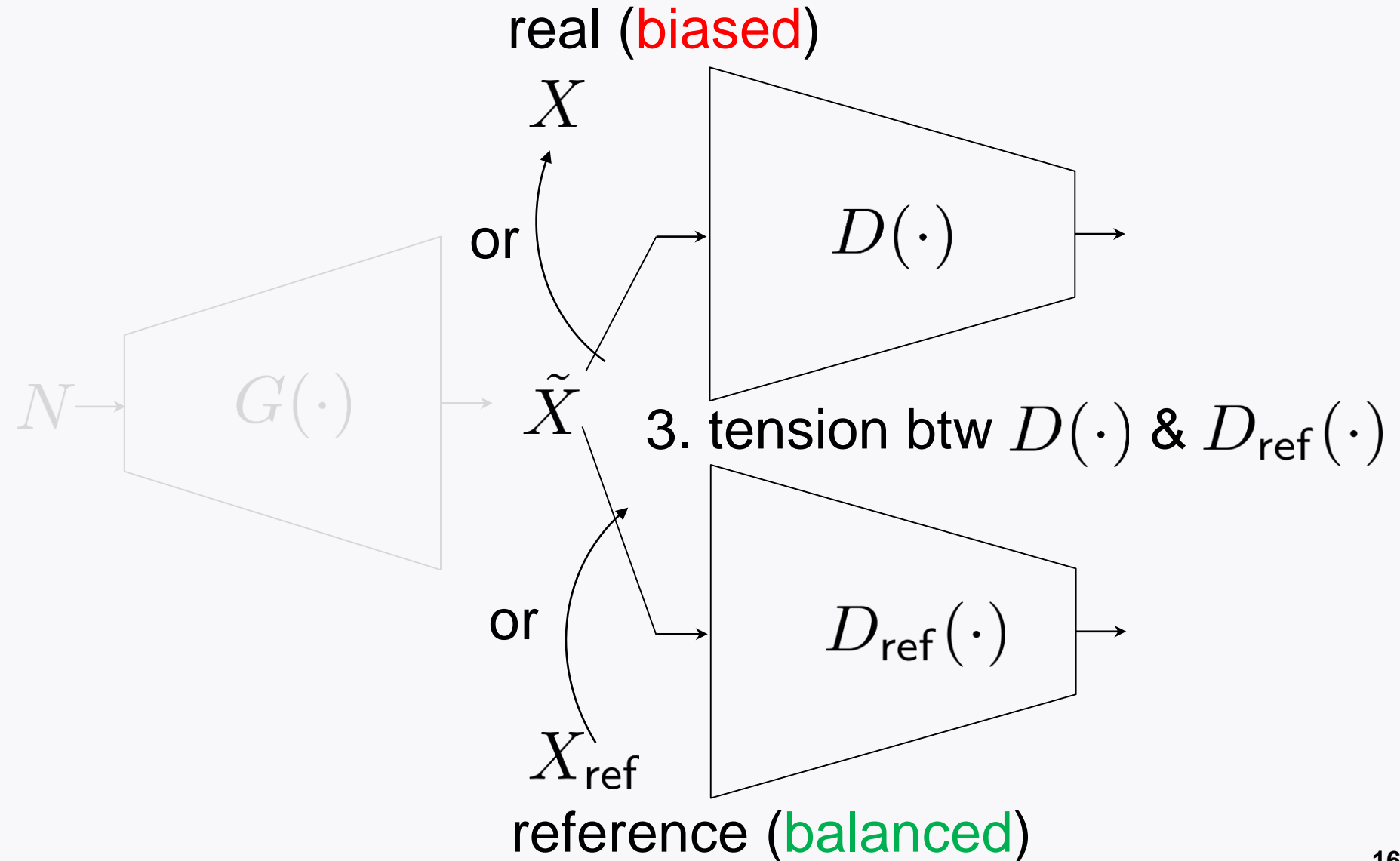
Three-way battle



Three-way battle



Three-way battle



Experiments

A benchmark real dataset: **CelebA**



Female: **male** \approx 90:10

$$m_{\text{real}} = 67,507$$

$$m_{\text{ref}} \approx 0.1m_{\text{real}}$$

$$\mathbb{P}_{\text{ref}}(z) \sim \text{Uniform}$$

Performance measure

1. A measure for the quality of generated samples

Fréchet Inception Distance (FID):

$$W^2(\mathcal{N}(\mu_{\text{real}}, \Sigma_{\text{real}}), \mathcal{N}(\mu_G, \Sigma_G))$$

↖
2nd order Wasserstein distance

2. A measure for fair sample generation

Fairness Discrepancy (FD):

$$\sqrt{\sum_{z \in \mathcal{Z}} (\mathbb{P}_G(z) - \mathbb{P}_{\text{ref}}(z))^2}$$

FID vs FD tradeoff

| | FID (quality) | FD (fairness) |
|---|----------------------|----------------------|
| <i>Non-fair</i> model (Hinge-loss-based GAN) | 8.76 ± 0.196 | 0.539 ± 0.002 |
| TVD-based <i>fair</i> gen. model | 14.13 ± 0.343 | 0.0431 ± 0.0097 |

Generated samples

TVD-based fair generative model:



Female:male = 54:46

Look ahead

Discuss a couple of other relevant issues.

Reference

- [1] S. Um and C. Suh. A fair generative model using total variation distance. submitted to *Advances in Neural Information Processing Systems 34 (NeurIPS)*, 2021
- [2] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. *In Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- [3] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. *In Advances in Neural Information Processing Systems 30 (NeurIPS)*, 2017.
- [4] K. Choi, A. Grover, T. Singh, R. Shu, and S. Ermon. Fair generative modeling via weak supervision. *In Proceedings of the 37th International Conference on Machine Learning (ICML)*, 2020.