# Understanding Mode Collapse in Generative Adversarial Networks

Subhasish Bandyopadhyay

August 9, 2025

## 1 Introduction

Generative Adversarial Networks (GANs), introduced in 2014 by Ian Goodfellow and his team, have become a cornerstone of modern machine learning, enabling the generation of realistic data samples from complex distributions. GANs consist of two neural networks: a generator $G$ that produces synthetic data from random noise $z$, and a discriminator $D$ that distinguishes real data from generated samples. The adversarial training objective is defined as:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}}[\log D(x)] + \mathbb{E}_{z \sim p_z}[\log(1 - D(G(z)))] \tag{1}$$

At equilibrium, the generator's distribution $p_g$ matches the real data distribution $p_{\text{data}}$, and $D(x) = 0.5$ for all $x$. However, training instability, particularly mode collapse, poses significant challenges. Mode collapse occurs when the generator produces limited or repetitive samples, failing to capture the full diversity of the data distribution. This chapter explores the causes of mode collapse, focusing on divergence measures, mode-seeking versus mass-covering behaviors, and mitigation strategies.

## 2 Understanding Mode Collapse in GAN Training

Mode collapse is a critical issue where the generator fails to cover all modes of the data distribution, instead focusing on a subset of high-density regions. This results in repetitive or nearly identical outputs, limiting the model's utility for tasks requiring diversity, such as data augmentation or creative content generation.

The root cause lies in the adversarial dynamics. The generator minimizes the discriminator's ability to classify its outputs as fake, often by perfecting a few modes that are easier to model or align with the current discriminator's weaknesses. For example, in a dataset with multiple Gaussian distributions, the generator might focus on one, ignoring others, as the standard GAN loss does not explicitly enforce diversity.

Key factors contributing to mode collapse include:

- **Imperfect Optimization**: Alternating gradient descent on non-convex objectives leads to local minima, where weak discriminators allow the generator to exploit low-diversity solutions.

- **Finite Sample Effects**: Mini-batches provide limited views of the data, enabling the generator to overfit to batch-specific patterns.

- **Hyperparameter Sensitivity**: Learning rates, architectures, and noise distributions can exacerbate collapse.

- **Divergence Choice**: The implicit divergence minimized by GANs encourages mode-seeking behavior.

Mode collapse can be detected through visual inspection, clustering metrics, or diversity scores in latent space. Its impact is significant in applications like image generation (monotonous outputs) or reinforcement learning (suboptimal policies).

# 3 Forward and Reverse KL Divergence: Foundations and Properties

To understand mode collapse, we examine divergence measures that quantify differences between distributions. The Kullback-Leibler (KL) divergence is central, with distinct behaviors depending on its direction.

The forward KL divergence is:

$$\mathrm{KL}(p_{\mathrm{data}}\|p_g) = \int_{-\infty}^{\infty} p_{\mathrm{data}}(x) \log\left(\frac{p_{\mathrm{data}}(x)}{p_g(x)}\right) dx = \mathbb{E}_{x\sim p_{\mathrm{data}}}\left[\log\frac{p_{\mathrm{data}}(x)}{p_g(x)}\right] \tag{2}$$

Properties:

- Non-negative, zero only if $p_{\mathrm{data}} = p_g$ almost everywhere.
- Asymmetric: $\mathrm{KL}(p_{\mathrm{data}}\|p_g) \neq \mathrm{KL}(p_g\|p_{\mathrm{data}})$.
- Infinite if $p_{\mathrm{data}}(x) > 0$ and $p_g(x) = 0$, requiring $p_g$ to cover $p_{\mathrm{data}}$'s support.

Minimizing forward KL promotes coverage, as $p_g$ must assign probability where $p_{\mathrm{data}}$ has mass to avoid infinite penalties, often leading to conservative, averaged outputs.

The reverse KL divergence is:

$$\mathrm{KL}(p_g\|p_{\mathrm{data}}) = \int_{-\infty}^{\infty} p_g(x) \log\left(\frac{p_g(x)}{p_{\mathrm{data}}(x)}\right) dx = \mathbb{E}_{x\sim p_g}\left[\log\frac{p_g(x)}{p_{\mathrm{data}}(x)}\right] \tag{3}$$

Properties:

- Non-negative and asymmetric.
- Infinite if $p_g(x) > 0$ and $p_{\mathrm{data}}(x) = 0$, penalizing $p_g$ for extraneous mass.
- Finite even if $p_g$ ignores parts of $p_{\mathrm{data}}$'s support.

Minimizing reverse KL encourages the generator to focus on high-density regions of $p_{\mathrm{data}}$, potentially ignoring other modes, which aligns with mode-seeking behavior.

GANs theoretically minimize the Jensen-Shannon (JS) divergence:

$$\mathrm{JS}(p\|q) = \frac{1}{2}\mathrm{KL}(p\|m) + \frac{1}{2}\mathrm{KL}(q\|m), \quad \text{where} \quad m = \frac{p+q}{2} \tag{4}$$

However, their training dynamics often resemble reverse KL, contributing to mode collapse.

# 4 Mode Seeking and Mass Covering: Behavioral Intuitions

The choice of divergence influences whether the model exhibits mode-seeking or mass-covering behavior:

- **Mass Covering (Forward KL)**: The generator spreads probability mass to cover all of $p_{\mathrm{data}}$, avoiding penalties in regions where $p_{\mathrm{data}} > 0$ but $p_g \approx 0$. This can lead to diffuse outputs, as seen in Variational Autoencoders (VAEs), which produce blurry but comprehensive samples.

- **Mode Seeking (Reverse KL)**: The generator concentrates on high-density modes of $p_{\mathrm{data}}$, minimizing penalties over its own samples. This produces sharp but less diverse outputs, as the generator can ignore low-density regions.

In GANs, the generator's updates, driven by samples from $p_g$, align with reverse KL, explaining their tendency for mode collapse. For example, in a bimodal dataset, a mass-covering model might generate samples between modes, while a mode-seeking model focuses on one.

# 5 Enforcing Coverage: Detailed Explanation

The difference between forward and reverse KL divergences explains why GANs are prone to mode collapse. The forward KL divergence can be rewritten as:

$$\text{KL}(p_{\text{data}}\|p_g) = \int p_{\text{data}}(x)\log p_{\text{data}}(x)\,dx - \int p_{\text{data}}(x)\log p_g(x)\,dx \tag{5}$$

$$= H(p_{\text{data}}) + \mathbb{E}_{x\sim p_{\text{data}}}[-\log p_g(x)] \tag{6}$$

The first term, $H(p_{\text{data}})$, is constant, so minimizing forward KL maximizes $\mathbb{E}_{x\sim p_{\text{data}}}[\log p_g(x)]$. This requires $p_g(x)$ to be large where $p_{\text{data}}(x) > 0$, enforcing coverage.

For reverse KL:

$$\text{KL}(p_g\|p_{\text{data}}) = -H(p_g) - \mathbb{E}_{x\sim p_g}[\log p_{\text{data}}(x)] \tag{7}$$

The expectation is over $p_g$, so the loss does not penalize missing modes in $p_{\text{data}}$. The generator can achieve low loss by focusing on a subset of modes, leading to mode dropping. This is why GANs, with updates based on samples from $p_g$, prioritize high-quality samples over comprehensive coverage.

# 6 Mitigation Strategies and Advanced Techniques

To address mode collapse, several strategies have been developed:

1. **Architectural Modifications**: Minibatch discrimination detects low diversity within batches, while unrolled GANs improve generator guidance.

2. **Alternative Divergences**: Wasserstein GANs (WGANs) use the Wasserstein-1 distance:

$$W(p,q) = \inf_{\gamma\in\Pi(p,q)} \mathbb{E}_{(x,y)\sim\gamma}[\|x-y\|]$$

   with gradient penalties (WGAN-GP) for stability.

3. **Regularization**: Spectral normalization and label smoothing stabilize training.

4. **Diversity-Promoting Losses**: InfoGAN maximizes mutual information for diverse outputs.

5. **Training Protocols**: Progressive growing and two-timescale update rules (TTUR) balance training dynamics.

6. **Ensemble Methods**: Multiple generators or discriminators cover different modes.

Recent advancements include diffusion-GAN hybrids and self-supervised GANs using contrastive learning.

# 7 Case Studies and Applications

- **Computer Vision**: Early GANs like DCGAN collapsed on CIFAR-10, while BigGAN mitigated this with large batches and truncation tricks.

- **Natural Language Processing**: Text GANs collapse to common phrases; RelGAN uses relational memory to improve diversity.

- **Medical Imaging**: GANs for CT scans risk missing rare pathologies, mitigated by techniques like CycleGAN.

- **Climate Modeling**: GANs generating weather patterns may underestimate extremes due to collapse.

# 8  Future Perspectives

Advancements in GANs include integration with large language models for multimodal generation and theoretical work on convergence and divergence bounds. Ethical considerations, such as ensuring diverse representations, remain critical.