

# Исследование устойчивости генеративно-состязательных сетей

М. Зубков<sup>а</sup>, А. Чубчева<sup>а</sup>, А. Филиппова<sup>а</sup>, Н. Сергеев<sup>а</sup>

<sup>а</sup>Московский физико-технический институт (национальный исследовательский университет), Москва, Россия

## Abstract

Генеративно-состязательные сети (GAN) — это разновидность генеративных моделей, способная генерировать такие сложные типы данных, как картинки, звук и видео, однако, такие модели часто демонстрируют нестабильное поведение во время обучения. Целью данного исследования является обзор и применение теоретически обоснованных методов стабилизации процесса обучения GAN в задаче генерации мультимодального распределения.

**Keywords:** GAN, Устойчивость, Оптимизация, WGAN, Спектральная нормализация

## 1. Введение

Обучение генеративно-состязательной сети (GAN) представляет собой «состязание» двух нейронных сетей: генератора и дискриминатора. Задача генератора заключается в создании объектов (например, изображения или видео), по правдоподобности конкурирующие с обучающей выборкой, взятой из некоторого распределения  $\mu_0$ . Задача дискриминатора состоит в оценке вероятности того, что образец семплирован из распределения  $\mu_0$ , а не создан генератором. Классическая схема обучения GAN представлена на рисунке 1.

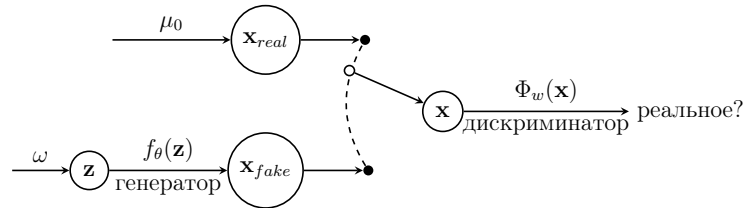


Рис. 1: Структура GAN

Для формулировки нашей задачи определим генератор формально:

**Определение 1.** Пусть  $X \subset \mathbb{R}^{n \times n}$  — множество исходных данных (например, картинок),  $Z \subset \mathbb{R}^{m \times m}$  — некоторое измеримое множество, а  $\omega$  — фиксированное распределение вероятности над множеством  $Z$ . Будем называть генератором параметризованную с параметром  $\theta \in \mathbb{R}^d$  нейронную сеть  $f_{\theta} : Z \rightarrow X$ . Распределением генератора будем называть множество образов вероятностных мер  $\{\mu_{\theta} = f_{\theta}(\omega), \text{ где } \theta \in \mathbb{R}^d\}$ .

Также приведем формально определение дискриминатора:

Email addresses: zubkov.md@phystech.edu (М. Зубков), chubcheva.ad@phystech.edu (А. Чубчева), filippova.am@phystech.edu (А. Филиппова), sergeev.ng@phystech.edu (Н. Сергеев)

15 **Определение 2.** Пусть  $Y \subset \mathbb{R}^{n \times n}$  – множество данных, содержащее образцы как из исход-  
 16 ной выборки, так и созданные генератором. Будем называть дискриминатором нейронную  
 17 сеть  $\Phi_w : Y \rightarrow [0, 1]$ , параметризованную с параметром  $w \in \mathbb{R}^l$ , сопоставляющую каждому  
 18 элементу  $y \in Y$  вероятность того, что этот  $y$  пришел из исходного распределения  $\mu_0$ .

19 Задача дискриминатора заключается в том, чтобы научиться отличать образцы, предло-  
 20 женные генератором  $Y_{\mu_\theta}$ , от образцов из обучающей выборки  $Y_{\mu_0}$ . Пусть  $\mathcal{J} : (\mu_\theta, \Phi_w) \rightarrow \mathbb{R}$  –  
 21 функция потерь, зависящая от дискриминатора  $\Phi_w$  и распределения  $\mu_\theta$ , созданного генера-  
 22 тором. Чем больше при фиксированном  $\mu_\theta$  значение функции  $\mathcal{J}$ , тем лучше дискриминатор  
 23  $\Phi_w$  отличает выборку из  $\mu_0$  от выборки из  $\mu_\theta$ . Тогда задача определения параметров дискри-  
 24 минатора формулируется следующим образом:

$$\max_{\Phi_w} \mathcal{J}(\mu_\theta, \Phi_w) = J(\mu_0, \mu_\theta) \quad (1)$$

25 Задача генератора заключается в создании таких объектов, чтобы фиксированный опти-  
 26 мальный дискриминатор как можно хуже отличал их от реальных объектов из распределения  
 27  $\mu_0$ . Таким образом, цель генератора в том, чтобы минимизировать меру отклонения сгенера-  
 28 рованной выборки от изначальной  $J(\mu_0, \mu_\theta)$ .

$$\min_{\mu_\theta} J(\mu_0, \mu_\theta) = \min_{\mu_\theta} \max_{\Phi_w} \mathcal{J}(\mu_\theta, \Phi_w) \quad (2)$$

29 Существует много различных вариантов определения функции  $\mathcal{J}$ , мы остановились на  
 30 функции потерь WGAN [1], ее преимущества будут описаны в разделе 2.

$$J(\mu_0, \mu_\theta) = \max_{\Phi_w} \mathcal{J}(\mu_\theta, \Phi_w) = \max_{\|\Phi_w\|_{Lip} \leq 1} \mathbb{E}_{x \sim \mu_0} [\Phi_w(x)] - \mathbb{E}_{x \sim \mu_\theta} [\Phi_w(x)], \quad (3)$$

31 где  $\|\Phi_w\|_{Lip} \leq 1$  обозначает, функция имеет константу Липшица не более единицы.

32 Использование GAN позволяет с высокой точностью решать такие задачи, как восстано-  
 33 вление изображений [2], увеличение разрешения изображения [3], увеличение размеров выбор-  
 34 ки данных [4], преобразование текста в изображения [5]. Однако, обучение GAN нестабильно,  
 35 как показано в [1]. Данный вопрос будет подробнее изучен в следующем разделе.

## 36 2. Обзор литературы

37 Одной из наиболее важных проблем, возникающих при обучении GAN, является их неста-  
 38 бильность. Под нестабильностью понимается затухание градиентов функции потерь генера-  
 39 тора  $J(\mu_0, \mu_\theta)$  по параметрам  $\theta$ , переобучение дискриминатора и чувствительность к гипер-  
 40 параметрам, таким как число итераций обучения дискриминатора  $n_{cr}$  и размер шага в гради-  
 41 ентном методе оптимизации. Для решения перечисленных проблем были предложены новые  
 42 архитектуры для генератора и дискриминатора [6], целые функции [1] и регуляризации [7, 8].

43 Ключевой шаг в решении проблемы нестабильности был сделан в статье [1]. В этой работе  
 44 предложена функция потерь 2, которая порождена метрикой Вассерштайна [9], являющейся  
 45 наиболее слабой метрикой в пространстве распределений. Для того чтобы определить наибо-  
 46 лее слабую метрику, рассмотрим последовательность нейронных сетей-генераторов  $f_{\theta_n}$ , где  
 47  $\theta_n \in \mathbb{R}^d$  – параметризация на  $n$ -том шаге обучения генератора. Тогда метрика  $\rho_0$  являеся наи-  
 48 более слабой, если произвольная последовательность  $f_{\theta_n}$ , сходящаяся по некоторой метрике  
 49  $\rho$ , будет также сходиться и по данной метрике  $\rho_0$ .

50 Функция потерь WGAN обладает перечисленными свойствами только в том случае, когда  
 51 выполнена липшицевость и оптимальность дискриминатора на каждой итерации. Свойство  
 52 липшицевости было подробно изучено в [7, 10, 11]. Вопрос оптимальности дискриминатора  
 53 изучался в [12, 13, 14].

54 Также в статье [1] был предложен алгоритм 1 обучения GAN, который мы будем изучать  
 55 в данной работе.

---

**Algorithm 1** Используемые гиперпараметры:  $\alpha = 9 \cdot 10^{-5}$ ,  $m = 64$ .

---

**Require:**  $\alpha$ , шаг градиентного спуска.  $m$ , размер батча.  $n_{cr}$ , количество итераций обучения дискриминатора.  $w_0$ , начальные параметры дискриминатора.  $\theta_0$ , начальные параметры генератора.

```

1: while не выполнен стоп критерий на  $\theta$  do
2:   for  $t = 0, \dots, n_{cr}$  do
3:     Сформировать выборку  $\{\mathbf{x}^{(i)}\}_{i=1}^m \sim \mu_0$  батч из реальных данных.
4:     Сформировать выборку  $\{\mathbf{z}^{(i)}\}_{i=1}^m \sim \mu_\theta$  батч из данных генератора.
5:      $g_w \leftarrow \nabla_w \left[ \frac{1}{m} \sum_{i=1}^m \Phi_w(\mathbf{x}^{(i)}) - \frac{1}{m} \sum_{i=1}^m \Phi_w(f_\theta(\mathbf{z}^{(i)})) \right]$ 
6:      $w \leftarrow w + \alpha \cdot \text{Adam}(w, g_w)$ 
7:   end for
8:   Сформировать выборку  $\{\mathbf{z}^{(i)}\}_{i=1}^m \sim \mu_\theta$  батч из данных генератора.
9:    $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^m \Phi_w(f_\theta(\mathbf{z}^{(i)}))$ 
10:   $\theta \leftarrow \theta - \alpha \cdot \text{Adam}(\theta, g_\theta)$ 
11: end while

```

---

56 Следующий шаг в решении проблемы нестабильности был сделан в статье [15], где была  
 57 доказана теорема 2, утверждающая, что при выполнении некоторых условий  $\|J_k(\mu_\theta)\| \rightarrow 0$   
 58 при  $k \rightarrow \infty$ , где  $k$  соответствует числу итераций обучения генератора. Это гарантирует схо-  
 59 димость функции потерь генератора к стационарной точке при бесконечном числе итераций.  
 60 Данная теорема играет ключевую роль в нашей статье, сформулируем ее в следующей сек-  
 61 ции.

### 62 3. Постановка задачи

63 В данной статье мы будем изучать влияние параметра  $n_{cr}$  на обучение GAN. Как показано  
 64 в статье [1], в общем случае GAN не сходится. Для решения данной проблемы в статье [15]  
 65 рассматривается теорема, которая формулирует достаточное условие сходимости GAN:

66 **Теорема 1 (Bertsekas [16]).** *Предположим, что  $J : \mu_\theta \rightarrow \mathbb{R}$  является ограничена снизу и*  
 67  *$L$ -гладкой, то есть  $\|\nabla J\|_{Lip} \leq L$ . Пусть  $\theta_{k+1} = \theta_k - \alpha \nabla J(\mu_{\theta_k})$ . Тогда  $\|\nabla J(\mu_{\theta_k})\| \rightarrow 0$  при*  
 68  *$k \rightarrow \infty$ .*

69 Данная теорема гарантирует сходимость к стационарной точке при бесконечном числе  
 70 итераций обновления параметров генератора и дискриминатора. Однако, согласно статье [17]  
 71 функция  $J$  в общем случае не является  $L$ -гладкой. Поэтому для обоснования сходимости была  
 72 доказана следующая теорема, которая гарантирует  $L$ -гладкость функции  $J$ :

73 **Теорема 2 (Chu et al. [15]).** *Пусть  $J : \mu_\theta \rightarrow \mathbb{R}$  выпуклая функция. Зафиксируем  $\mu := \mu_0$*   
 74 *и рассмотрим оптимальный дискриминатор:  $\Phi_\mu : Y \rightarrow [0, 1]$ . Пусть он удовлетворяет*  
 75 *следующим условиям:*

(D1)  $x \mapsto \Phi_\mu(x)$  –  $\alpha$ -липшицева,

(D2)  $x \mapsto \nabla_x \Phi_\mu(x)$  –  $\beta_1$ -липшицева,

(D3)  $\mu \mapsto \nabla_x \Phi_\mu(x)$  –  $\beta_2$ -липшицева по метрике 1-Wasserstein [9].

Также, пусть семейство генераторов  $f_\theta(\omega)$  удовлетворяет условиям:

(G1)  $\theta \mapsto f_\theta(z)$   $A$ -липшицева в среднем для  $z \sim \omega$ , то есть,

$$\mathbb{E}_{z \sim \omega}[\|f_{\theta_1}(z) - f_{\theta_2}(z)\|_2] \leq A\|\theta_1 - \theta_2\|_2, \text{ и}$$

(G2)  $\theta \mapsto D_\theta f_\theta(z)$   $B$ -липшицева в среднем для  $z \sim \omega$ , то есть,

$$\mathbb{E}_{z \sim \omega}[\|D_{\theta_1} f_{\theta_1}(z) - D_{\theta_2} f_{\theta_2}(z)\|_2] \leq B\|\theta_1 - \theta_2\|_2.$$

Тогда  $\theta \mapsto J(\mu_\theta)$  является  $L$ -гладкой, где  $L = \alpha B + A^2(\beta_1 + \beta_2)$ .

Одним из ключевых условий теоремы 2 является оптимальность дискриминатора на каждой итерации обучения генератора. Изучение данного вопроса представлено в статье [12]. Автор утверждает, что при оптимальном дискриминаторе процесс обучения генератора сходится с вероятностью 1. С другой стороны, в статье [13] показано, что обучение дискриминатора до оптимальности ведет к затуханию градиента для генератора. В связи с этим противоречием мы решили изучить влияние  $n_{cr}$  на процесс обучения GAN.

#### 4. Методы

В статье [15] было доказано, что условия (D1)-(D3), (G1), (G2) теоремы 2 эквиваленты уже изученным техникам стабилизации GAN, представленным в таблице 2.

Условие	Методы решения
(D0)	Оптимальность дискриминатора
(D1)	Спектральная нормализация [10] липшицева регуляризация [11]
(D2)	Гладкие функции активации, спектральная нормализация [10]
(D3)	Состязательная атака [8] WGAN-GP [7]
(G1)	Нет метода
(G2)	Нет метода

Рис. 2: Методы

В нашем исследовании мы будем использовать функцию потерь (2) и алгоритм 1, предложенные в статье [1]. При этом из описанных в таблице 2 методов мы будем использовать гладкие функции активации и спектральную нормализацию, предложенную в статье [10]. В нашей задаче дискриминатор является многослойным перцептроном:

$$\Phi_w(x) = W_L(a_L(W_{L-1}(\dots(W_1x)\dots)),$$

где  $W_l \in \mathbb{R}^{(d_l+1) \times d_{l+1}}$  –  $l$ -тый линейный слой,  $a_l$  – функция активации после  $l$ -того линейного слоя,  $L$  – количество слоев дискриминатора. Такая архитектура представлена на рисунке 3(с) для случая  $L = 4$ .

Также в нашей модели дискриминатора мы будем использовать спектральную нормализацию (4). В статье [10] показывается, что:

$$\|\Phi_w\|_{Lip} \leq \prod_{i=1}^L \|W_i\|_2$$

$l$ -тый слой можно представить в виде матрицы  $W_l$ , для которой  $\|W_l\|_2 = \rho(W_l^T W_l)$ , где  $\rho(W_l) = \max\{|\sigma| : \text{где } \sigma - \text{сингулярное число матрицы } W_l^T W_l\}$ . Таким образом, спектральной нормализации будет соответствовать преобразование матриц линейных слоев вида:

$$W_l^{SN} = \frac{W_l}{\|W_l\|_2} \quad (4)$$

В экспериментах мы будем использовать как модель со спектральной нормализацией, так и без. Для краткости будем обозначать их SNGAN и noSNGAN соответственно. Архитектуры их дискриминаторов и генераторов представлены на рисунке 3.

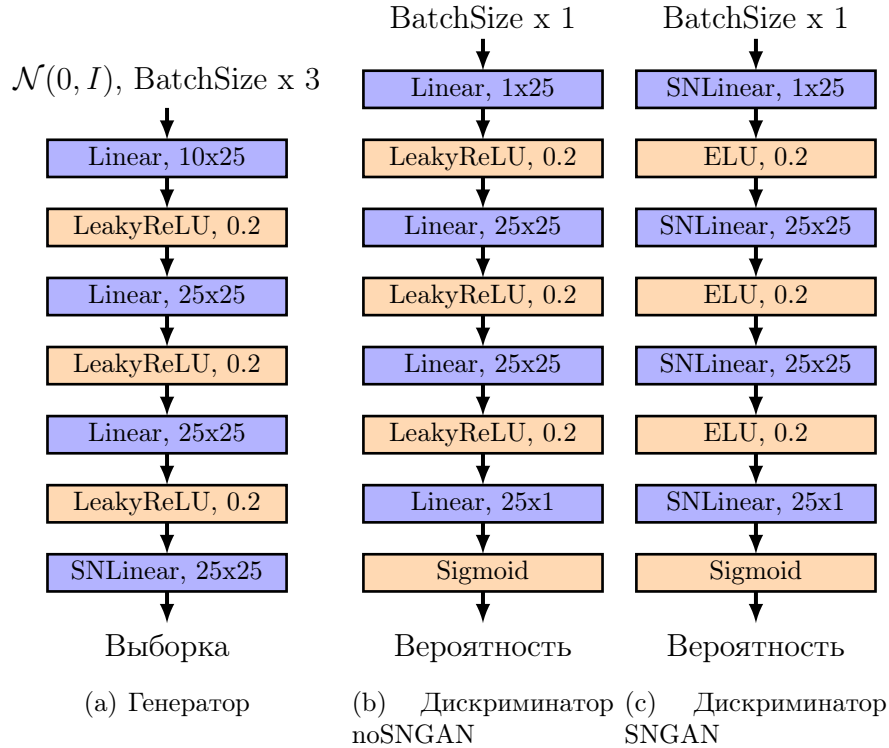


Рис. 3: Архитектуры нейронных сетей

В качестве исходных данных мы взяли выборки из линейной комбинации нормальных распределений. Пусть  $\xi \sim \mathcal{N}(\mu_1, \sigma_1)$ ,  $\eta \sim \mathcal{N}(\mu_2, \sigma_2)$ ,  $\zeta \sim \mathcal{N}(\mu_3, \sigma_3)$ . Задача заключается в том, чтобы научить генератор создавать выборку из целевого распределения  $\xi + \eta + \zeta \sim \mu_0$ . Пример такого распределения представлен на рисунке 4, где  $\mu_1 = 2, \sigma_1 = 0.4, \mu_2 = 0, \sigma_2 = 0.55, \mu_3 = 5, \sigma_3 = 0.25$ .

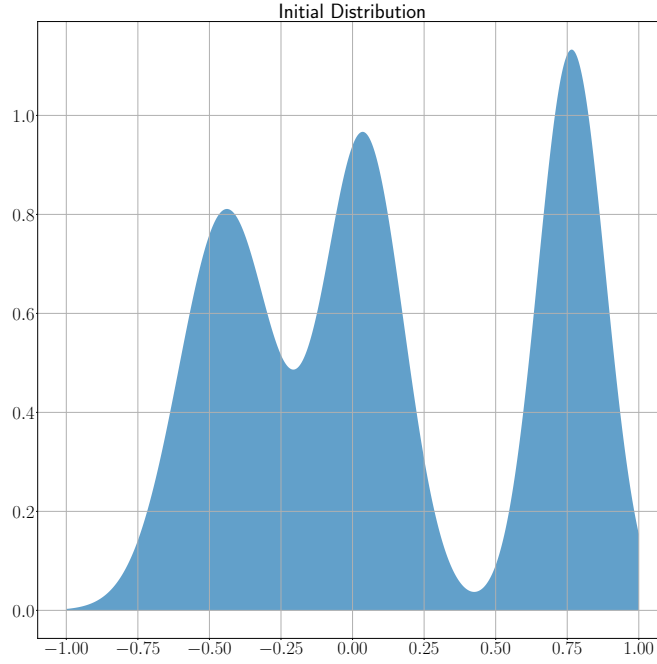


Рис. 4: Исходное мультимодальное распределение  $\mu_0$

105 Далее необходимо определиться с метрикой качества для объектов, созданных генерато-  
 106 ром. Для сравнения распределений двух выборок существует ряд статистических критериев  
 107 однородности. Мы воспользовались  $p$ -value критерия Колмогорова-Смирнова [18], так как  
 108 при верности нулевой гипотезы он зависит от  $\mu_0$  и  $\mu_\theta$ .

Поставим две гипотезы:

$$H_0 : \mu_\theta = \mu_0$$

$$H_1 : \mu_\theta \neq \mu_0$$

109  $p$ -value будем называть вероятность получить такое же или более экстремальное значение  
 110 некоторой статистики, при условии, что гипотеза  $H_0$  верна. В нашем случае такой статисти-  
 111 кой будет статистика Колмогорова-Смирнова 5. С помощью  $p$ -value мы будем оценивать,  
 112 насколько похожи исходное и сгенерированное распределения, причем чем больше значение  
 113  $p$ -value, тем лучше они совпадают. Главное свойство статистики Колмогорова-Смирнова ука-  
 114 зана в следующей теореме:

**Теорема 3.** Пусть  $(X_1, \dots, X_n)$  – выборка из распределения генератора  $\mu_\theta$ , а выборка  $(Y_1, \dots, Y_m)$  – выборка из целевого распределения  $\mu_0$ . Пусть  $\hat{F}_n, \hat{G}_m$  – эмпирические функции распределения выборок  $\{X_i\}_{i=1}^n$  и  $\{Y_i\}_{i=1}^m$  соответственно, что по определению значит:  $\hat{F}(x) = \sum_{i=1}^n \mathbb{I}(X_i \leq x)$ , где:

$$\mathbb{I}(X_i \leq x) = \begin{cases} 1, & \text{если } X_i \leq x \\ 0, & \text{иначе} \end{cases}$$

Пусть также:

$$D_{n,m} = \sup_x |\hat{F}_n(x) - \hat{G}_m(x)| \quad (5)$$

115 При верной гипотезе  $H_0 : \mu_0 = \mu_\theta$  статистика Колмогорова-Смирнова имеет распреде-  
 116 ление, не зависящее от распределений  $\mu_0, \mu_\theta$ .

## 117 5. Эксперименты

118 Исходный код для воспроизведения результатов можно найти в репозитории [https://](https://github.com/maximzubkov/opt-project)  
 119 [github.com/maximzubkov/opt-project](https://github.com/maximzubkov/opt-project)

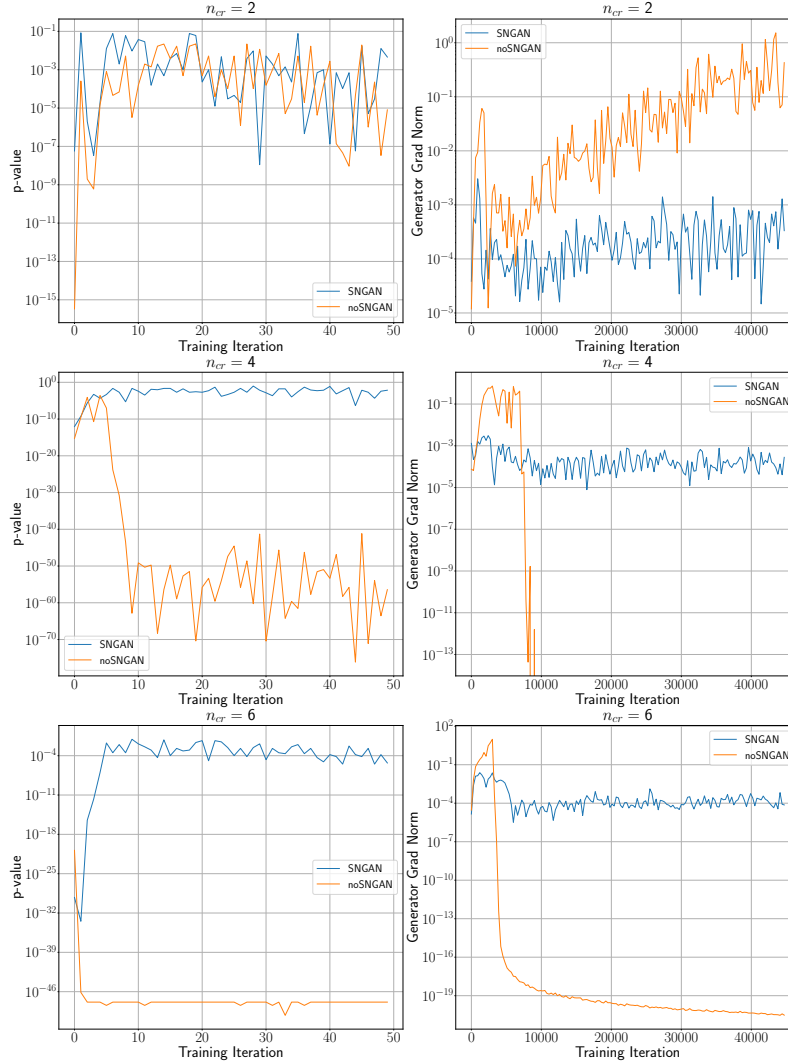


Рис. 5: Сравнение  $p$ -value и нормы градиента генератора для noSNGAN и SNGAN при различных значениях  $n_{cr}$

120 Следуя подходу, предложенному в теореме 2, мы взяли постоянный шаг спуска равный  
 121  $9 \cdot 10^{-5}$ . Процесс вычисления шага представлен в статье [15]. Были получены следующие  
 122 результаты:

- 123 1. При  $n_{cr} > 2$  noSNGAN показывает плохие результаты, что можно видеть на рисун-  
 124 ках 6(b), 7(b), 8(b). Такой вывод можно сделать по тому, что  $p$ -value noSNGAN стреми-  
 125 тельно убывает к нулю. Также по рисунку 5 можно видеть, что при  $n_{cr} > 2$  происходит  
 126 затухание градиента. SNGAN сходится при любых значениях  $n_{cr}$ , сравнение сходимости  
 127 по  $p$ -value представлено на рисунке 5.
- 128 2. Зависимость  $p$ -value генератора от числа  $n_{cr}$  представлена на рисунке 9. Можно видеть,  
 129  $p$ -value SNGAN почти не зависит от выбранного значения  $n_{cr}$ . Однако, по рисунку 8(a)  
 130 можно заметить, что при  $n_{cr} = 6$  мы столкнулись с проблемой коллапса мод. Можно

предположить, что при больших значениях  $n_{cr}$  генератору требуется больше эпох, чтобы обучиться, так как на рисунке 8(а) видно, что ассигасу достаточно сильно колеблется.

3. Кроме того мы оценили насколько влияет использование спектральной нормализации на время обучения, результаты приведены в таблице 1. Как и ожидалось, SNGAN требует больше времени на итерацию, при этом при увеличении  $n_{cr}$  обоим GAN требуется меньше времени, так как необходимо делать меньше обновлений параметров генератора.

SN	$n_{cr}$	Время одной эпохи, сек.
—	2	$2.92 \pm 0.04$
—	4	$2.77 \pm 0.04$
—	6	$2.62 \pm 0.03$
+	2	$4.59 \pm 0.04$
+	4	$4.26 \pm 0.05$
+	6	$4.12 \pm 0.02$

Таблица 1: Размер обучающей выборки для каждого эксперимента составлял  $5 \cdot 10^4$ . Шаг градиентного спуска  $9 \cdot 10^{-5}$ . Количество эпох обучения – 50.

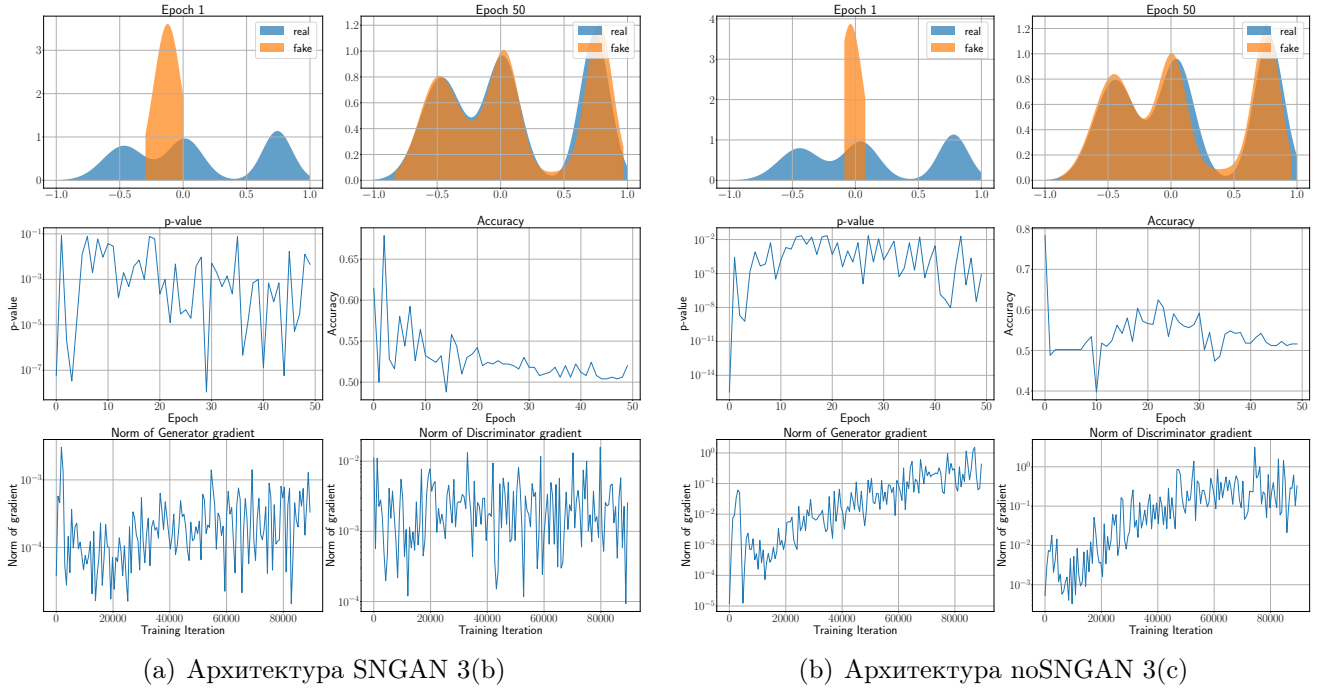


Рис. 6: Результаты экспериментов при  $n_{cr} = 2$



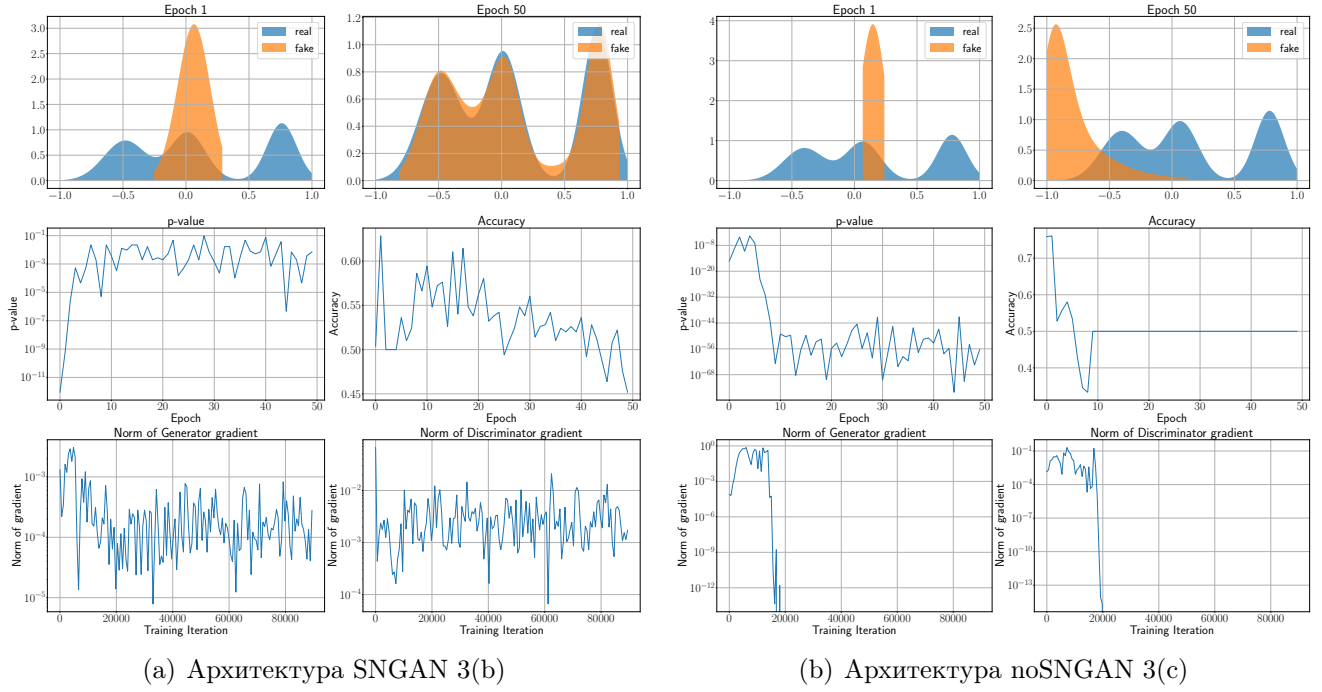


Рис. 7: Результаты экспериментов при  $n_{cr} = 4$

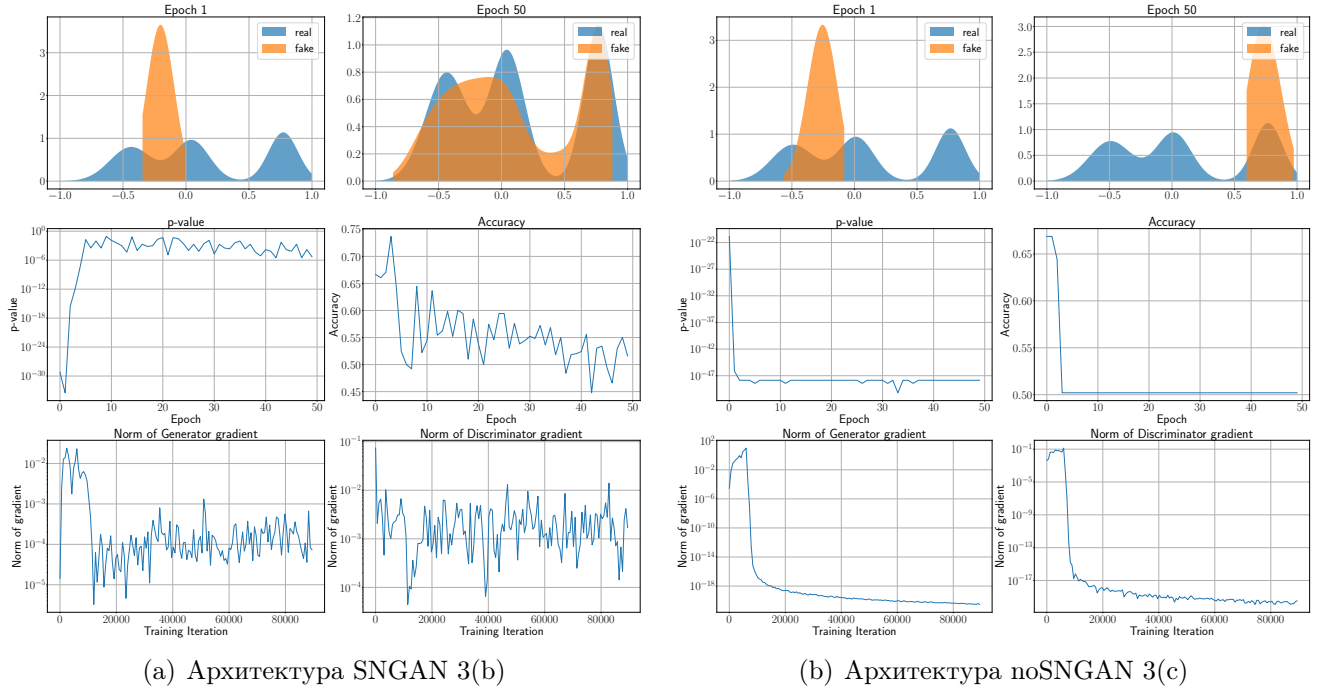


Рис. 8: Результаты экспериментов при  $n_{cr} = 6$

## 137 6. Выводы

138 По полученным результатам были сделаны следующие заключения:

- 139 1. Исследуемый параметр  $n_{cr}$  оказывает сильное влияние на обучение noSNGAN и SNGAN.
- 140 Спектральная нормализация ограничивает норму градиента, а в её отсутствие гради-

енты «взрываются», что приводит к переобучению дискриминатора. Как следствие, наблюдается затухание градиента, и остановка процесса обучения. В тех экспериментах, где задействована спектральная нормализация, переобучение дискриминатора контролируется, в результате чего GAN генерирует более правдоподобное распределение.

2. В ходе исследования мы выяснили, что увеличение  $n_{cr}$  влечет за собой уменьшение времени на одну эпоху обучения, хотя с другой стороны требуется большее число эпох для обучения генератора.

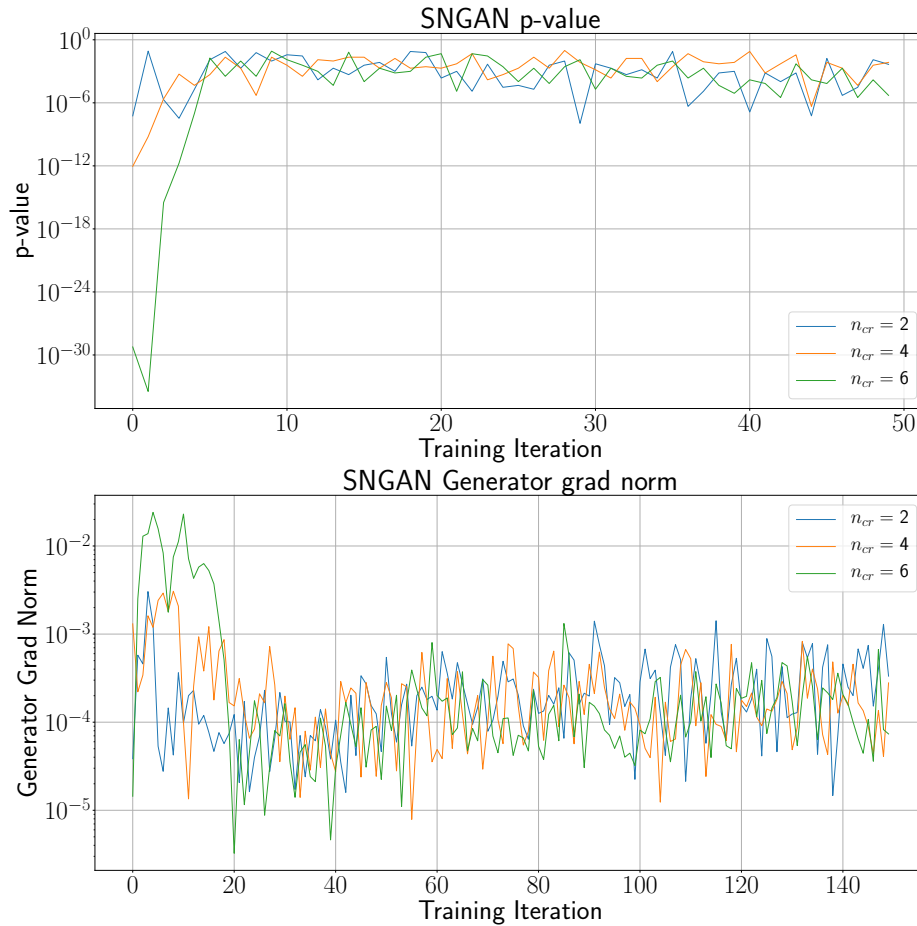


Рис. 9: Сравнение p-value и нормы градиента генератора при различных значениях  $n_{cr}$  для SNGAN

## 7. Результаты и будущая работа

В данной статье мы изучили влияние параметра  $n_{cr}$  на процесс обучения генеративно-состязательной нейронной сети, а также убедились в эффективности техник, предложенных теоремой 2.

В будущем планируется применить полученные результаты к обучению генеративных моделей на более сложных данных, таких как изображения, с использованием сетей, имеющих большее число нейронов и слоев.

Вопрос об оценке качества работы генератора остается открытым: для выборки из одномерных данных мы использовали статистический критерий, но такой подход не работает в

157 задачах, где данные имеют более сложную структуру. Как правило в задаче генерации изображений, используются такие метрики как FID и Inception Score [19], однако данные метрики 158  
159 несовершенны, как показано в [20].

160 Для глубоких и сложных архитектур генератора и дискриминатора сложно оценить параметры  $A, B, \alpha, \beta_1, \beta_2$  из теоремы 2. Поэтому возникла идея воспользоваться ансамблем простых нейронных сетей. Как утверждается в статье [21], residual-соединения образуют ансамбль из простых моделей, именно поэтому сети, снабженные residual-блоками обучаются 162  
163 легче и быстрее остальных. Данную идею можно применить и в нашей задаче: как было показано ранее, GAN с глубокими сетями генератора и дискриминатора сложнее обучить, однако 164  
165 сравнительно неглубокие архитектуры, как 3, обучаются быстрее и проще. Следовательно 166  
167 использование residual-сетей оправданно.

## 168 Список литературы

- 169 [1] M. Arjovsky, S. Chintala, L. Bottou, Wasserstein gan, arXiv preprint arXiv:1701.07875 (2017).
- 170 [2] S. A. Hussein, T. Tirer, R. Giryes, Image-adaptive gan based reconstruction, arXiv preprint 171  
arXiv:1906.05284 (2019).
- 172 [3] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, 173  
J. Totz, Z. Wang, et al., Photo-realistic single image super-resolution using a generative 174  
adversarial network, in: Proceedings of the IEEE conference on computer vision and pattern 175  
recognition, 2017, pp. 4681–4690.
- 176 [4] F. H. K. d. S. Tanaka, C. Aranha, Data augmentation using gans, arXiv preprint 177  
arXiv:1904.09135 (2019).
- 178 [5] H. Zhang, T. Xu, H. Li, S. Zhang, X. Wang, X. Huang, D. Metaxas, Stackgan: Text 179  
to photo-realistic image synthesis with stacked generative adversarial networks, 2016. 180  
arXiv:1612.03242.
- 181 [6] L. Metz, B. Poole, D. Pfau, J. Sohl-Dickstein, Unrolled generative adversarial networks, arXiv 182  
preprint arXiv:1611.02163 (2016).
- 183 [7] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, A. C. Courville, Improved training of 184  
wasserstein gans, in: Advances in neural information processing systems, 2017, pp. 5767–5777.
- 185 [8] B. Zhou, P. Krähenbühl, Don't let your discriminator be fooled (2018).
- 186 [9] A. Andoni, P. Indyk, R. Krauthgamer, Earth mover distance over high-dimensional spaces., 187  
in: SODA, volume 8, 2008, pp. 343–352.
- 188 [10] T. Miyato, T. Kataoka, M. Koyama, Y. Yoshida, Spectral normalization for generative 189  
adversarial networks, arXiv preprint arXiv:1802.05957 (2018).
- 190 [11] Y. Yoshida, T. Miyato, Spectral norm regularization for improving the generalizability of 191  
deep learning, arXiv preprint arXiv:1705.10941 (2017).
- 192 [12] J. Li, A. Madry, J. Peebles, L. Schmidt, On the limitations of first-order approximation in 193  
gan dynamics, arXiv preprint arXiv:1706.09884 (2017).

- 194 [13] M. Wiatrak, S. V. Albrecht, Stabilizing generative adversarial network training: A survey,  
195 arXiv preprint arXiv:1910.00927 (2019).
- 196 [14] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, X. Chen, Improved  
197 techniques for training gans, in: Advances in neural information processing systems, 2016,  
198 pp. 2234–2242.
- 199 [15] C. Chu, K. Minami, K. Fukumizu, Smoothness and stability in gans, arXiv preprint  
200 arXiv:2002.04185 (2020).
- 201 [16] D. Bertsekas, Nonlinear Programming, Athena scientific optimization and computation series,  
202 Athena Scientific, 2016. URL: <https://books.google.ru/books?id=Tw0ujgEACAAJ>.
- 203 [17] L. Mescheder, S. Nowozin, A. Geiger, The numerics of gans, in: I. Guyon, U. V. Luxburg,  
204 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett (Eds.), Advances in Neural  
205 Information Processing Systems 30, Curran Associates, Inc., 2017, pp. 1825–1835. URL: <http://papers.nips.cc/paper/6779-the-numerics-of-gans.pdf>.
- 207 [18] S. Anulova, N. Krylov, R. Liptser, A. Shiryaev, A. Y. Veretennikov, Probability Theory III:  
208 Stochastic Calculus, volume 45, Springer Science & Business Media, 2013.
- 209 [19] M. J. Chong, D. Forsyth, Effectively unbiased fid and inception score and where to find them,  
210 arXiv preprint arXiv:1911.07023 (2019).
- 211 [20] M. Lucic, K. Kurach, M. Michalski, S. Gelly, O. Bousquet, Are gans created equal? a large-  
212 scale study, in: Advances in neural information processing systems, 2018, pp. 700–709.
- 213 [21] A. Veit, M. J. Wilber, S. Belongie, Residual networks behave like ensembles of relatively  
214 shallow networks, in: Advances in neural information processing systems, 2016, pp. 550–558.