

Содержание

Введение	5
1 Аналитический раздел	6
1.1 Цель и задачи работы	6
1.2 Обнаружение аномалий	6
1.2.1 Классификация методов обнаружений аномалий	7
1.2.2 Признаковое представление данных	7
1.3 Результат метода обнаружения аномалий	8
1.4 Виды аномалий	8
1.4.1 Нормализация данных	10
1.4.1.1 Основные методы нормализация данных	10
1.5 Неконтролируемые алгоритмы обнаружения аномалий	11
1.5.1 Вероятностный подход	11
1.5.2 Линейный подход	11
1.5.3 Метрические методы	12
1.6 Виды аномалий	13
1.7 Постановка задачи	14
1.8 Существующие подходы к созданию всячины	14
2 Конструкторский раздел	17
2.1 Оценка качества изображений	17
2.1.1 PSNR	17
2.1.2 SSIM	17
2.1.3 Алгоритм случайного распределения	18
2.2 Виды случайных распределений	18
2.2.1 Белый шум	18
2.2.2 Коричневый шум	19
2.2.3 Гауссовский шум	20
2.2.4 Фиолетовый шум	20
2.2.5 Розовый и синий шумы	21
2.3 Алгоритм Байера	22
2.4 Алгоритм Флойда-Стейнберга	22
2.5 Ложный алгоритм Флойда-Стейнберга	23
2.6 Алгоритм Джарвиса,Джунка и Нинка	23
2.7 Первый алгоритм Юлиомы	24
3 Технологический раздел	25
3.1 Выбор языка программирования	25
3.2 Выбор вспомогательных библиотек	25
3.2.1 Диаграмма классов	25

4	Исследовательский раздел	26
4.1	Время дизеринга раличных алгоритмов	26
4.2	Качество получаемого изображения	27
4.3	Размер получаемого изображения	27
	Заключение	29
	Список использованных источников	30

Глоссарий

Выборка/выборка данных — конечный набор прецедентов (объектов, случаев, событий, испытуемых, образцов, и т.п.), некоторым способом выбранных из множества всех возможных прецедентов, называемого генеральной совокупностью[1].

— Метка(ярлык) - порция данных, идентифицирующая набор данных, описывающая его определенные свойства и обычно хранимая в том же пространстве памяти, что и набор данных[2].

классификатор?

Теория распознавания образа — раздел информатики и смежных дисциплин, развивающий основы и методы классификации и идентификации предметов, явлений, процессов, сигналов, ситуаций и т.п. объектов, которые характеризуются конечным набором некоторых свойств и признаков. ddos-атака? Датасет - набор данных[3]

Введение

Задача поиска аномалий относится к одному из популярных способов машинного обучения - обучению без учителя. В настоящее время задачу поиска аномалий активно решают во многих областях жизнедеятельности:

- а) Защита информации и безопасность
- б) Социальная сфера и медицина
- в) Банковская и финансовая отрасль
- г) Распознавание и обработка текста, изображений, речи
- д) Другие сферы деятельности(например, мониторинг неисправностей механизмов)

Задачей поиска выбросов, как частный случай задачи поиска аномалий так же занимаются во вышеперечисленных отраслях.

Количество данных в мире удваивается примерно каждые два года. Поэтому актуальной задачей является разработка новых методов и усовершенствования старых методов поиска выбросов.

В данной работе предлагается новый метод, позволяющий найти аномалии в выборках данных.

1 Аналитический раздел

1.1 Цель и задачи работы

Целью данной работы является создание программного комплекта для обнаружения выбросов временных рядов в собираемых данных. Для достижения данной цели необходимо решить следующие задачи:

- ПЕРЕПИСАТЬ ИЗ ПРЕЗЕНТАЦИИ
- проанализировать предметную область и существующие методы обнаружения выбросов
- разработать метод обнаружения выбросов
- smth
- создать ПО, реализующего разработанный метод обнаружения выбросов
- провести вычислительный эксперименты с использованием разработанного метода

1.2 Обнаружение аномалий

В машинном обучении обнаружение "ненормальных" экземпляров в наборах данных всегда представляло большой интерес. Этот процесс широко известен как обнаружение аномалий или обнаружение выбросов. Вероятно, первое определение было дано Граббсом[4] в 1969 году: "Относительное наблюдение или выброс - это элемент выборки, который, заметно отличается от других членов выборки, в которых он встречается ". Хотя это определение по-прежнему актуально и сегодня, мотивация для обнаружения этих выбросов сейчас совсем другая. Тогда основная причина обнаружения заключалась в том, чтобы удалить выбросы из данных для обучения, поскольку алгоритмы распознавания были весьма чувствительны к выбросам в данных. Эта процедура также называется очищением данных. После разработки более надежных классификаторов интерес к обнаружению аномалий значительно снизился. Однако в 2000 году произошел поворотный момент, когда исследователи стали больше интересоваться самими аномалиями, поскольку они часто связаны с особенно интересными событиями или подозрительными данными. С тех пор было разработано много новых алгоритмов, которые оцениваются в этой статье. В этом контексте определение Граббса также было расширено, так что сегодня аномалии, как известно, имеют две важные характеристики:

- а) Аномалия отличается от нормы по своим особенностям
- б) Аномалия редко встречается в наборе данных по сравнению с "нормальными" данными

1.2.1 Классификация методов обнаружений аномалий

В отличие от хорошо известной системы классификации, где учебные данные используются для обучения классификатора, а результаты измерений данных оцениваются впоследствии, возможно множество вариантов, когда речь идет об обнаружении аномалий. Метод обнаружения аномалий, которая будет использоваться, зависит от ярлыков, доступных в наборе данных, и мы можем выделить три основных типа:

а) Обучение с учителем. Доступны полностью размеченные данные для обучения и для тестов. Обычный классификатор может быть обучен один раз и применяться впоследствии. Это похоже на традиционное распознавание образов, за исключением того, что классы обычно сильно не сбалансированы. Поэтому не все алгоритмы классификации идеально подходят для этой задачи. Для многих применений аномалии не известны заранее или могут возникать спонтанно в качестве новинок на этапе тестирования.

б) Обучение с частичным привлечением учителя. Обучение использует учебные и тестовые наборы данных. Данные обучения состоят только из нормальных данных без каких-либо аномалий. Основная идея заключается в том, что модель нормального класса изучается, а аномалии могут быть обнаружены впоследствии, отклоняясь от этой модели. Эта идея также известна как «одноклассовая» классификация [5].

в) Обучение без учителя. Самый гибкий способ, который не требует каких-либо меток. Кроме того, нет различия между учебным и тестовым наборами данных. Идея заключается в том, что алгоритм обнаружения аномалии оценивает данные исключительно на основе внутренних свойств набора данных. Как правило, расстояния или плотности используются для оценки того, что является нормальным, а что является выбросом. В этой работе основное внимание будет уделено именно этому способу. Так же этот способ иногда называют "неконтролируемый способ обнаружения аномалий".

1.2.2 Признаковое представление данных

В дальнейшем будем исходить из предположения что данные имеют признаковое представление, т.е. каждый объект x представлен в виде вектора \mathbb{R}^d . В задаче обучения без учителя задача обнаружения аномалий формулируется следующим образом: в заданном множестве X для каждого элемента выдать 0, если этот объект относится. При этом правильных ответов не предоставляется.

В аналогичной задаче обучения с учителем на некоторой выборке входных данных X_{train} , называемой тренировочной выборкой, известен правильный ответ, т.е. для каждого элемента $x \in X_{train}$ представлены метки $y \in \{0,1\}$, характеризующие является ли объект аномалией или нет. Для выборки данных, для которой метки не

предоставлены, задача сводится к задаче бинарной классификации, а значит может решаться при помощи любых алгоритмов машинного обучения с учителем. Возможны и "вырожденные" случаи, когда все метки тренировочного набора данных одинаковы. В таком случае алгоритмы выдают неправильный результат.

к классу нормальных данных, и 1, если этот объект аномален.

1.3 Результат метода обнаружения аномалий

Существует два варианта выходных данных алгоритма обнаружения аномалии. Во-первых, метка может использоваться как результат, указывающий, является ли экземпляр аномалией или нет. Во-вторых, оценка или достоверность могут быть более информативным результатом, указывающим на степень аномалии. А алгоритмах метода обучения с учителем зачастую используются метки как выходные данные. С другой стороны, для алгоритмах с частичным привлечением учителя и без учителя обнаружения аномалий чаще встречаются оценки.

1.4 Виды аномалий

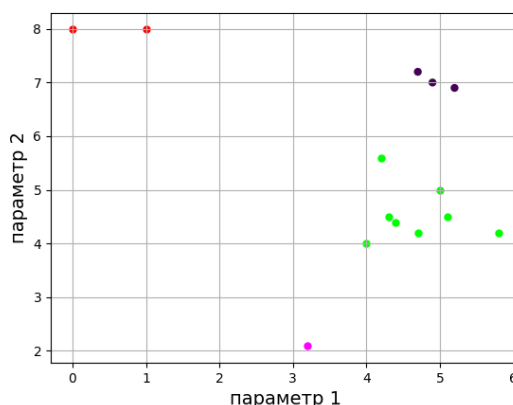


Рисунок 1.1 — Простой двумерный пример

Основная идея алгоритмов обнаружения аномалий заключается в обнаружении экземпляров данных в наборе данных, которые отклоняются от нормы. Однако на практике существует множество случаев, когда это основное предположение является неоднозначным. На рис. 2 показаны некоторые из этих случаев с использованием простого двумерного набора данных. Две аномалии могут быть легко идентифицированы визуально: красные точки сильно отличаются значениям параметров от областей плотной группировки точек и поэтому называются глобальными аномалиями. Когда мы смотрим на весь набор данных в целом, то фиолетовую точку можно отнести к тому же классу, что и зеленые точки. Однако, когда мы фокусируемся только на кластере зеленых точек и сравниваем его с фиолетовой, пренебрегая

всеми другими точками, то её можно рассматривать как аномалию. Поэтому фиолетовая точка называется локальной аномалией, так как она аномальна по сравнению с ее близкой окрестностью. В зависимости от цели исследования, нас могут интересовать местные аномалии или нет. Другой интересный вопрос заключается в том, следует ли рассматривать точки черного кластера как три аномалии или как (небольшой) кластер. Эти явления называются микрокластерами, а алгоритмы обнаружения аномалий должны присваивать оценку (вероятность того, что точка является аномалией) каждой точке этого кластера значения большие, чем точкам зеленого кластера, но меньше, чем красным точкам. Этот простой пример уже показывает, что аномалии не всегда очевидны, а оценка намного полезнее, чем назначение двоичных меток.

Обычно под аномалией принимают точки красные точки, так как их характеристики значительно отличаются от характеристик датасета, а так же их небольшое количество. Однако, такой принцип обнаружения иногда терпит неудачу. Например, при хакерских ddos-атаках, большая часть трафика - необычная, аномальная. В этом случае алгоритм обучения без учителя потерпит неудачу и не сможет выделить хакерскую атаку как аномальное поведение.

Задача обнаружения одиночных аномальных экземпляров в более крупном наборе данных (как это представлено до сих пор) называется обнаружением точечной аномалии [6]. Сегодня почти все доступные неконтролируемые алгоритмы обнаружения относятся к этому типу. Если аномальная ситуация представлена как множество многих случаев, это называется коллективной аномалией. Каждый из этих экземпляров не обязательно является точечной аномалией, но только определенная их комбинация определяет аномалию. Предыдущим приведенным примером возникновения нескольких специфических шаблонов доступа при обнаружении ddos-атак является такая коллективная аномалия. Третий вид - это контекстуальные аномалии, которые описывают эффект, что точка может рассматриваться как нормальная, но когда данный контекст учитывается, то точка оказывается аномалией. Самым распространенным контекстом является время. В качестве примера предположим, что мы измеряем температуру в диапазоне от -30°C до $+40^{\circ}\text{C}$ в течение года. Таким образом, температура 25°C кажется довольно нормальной, но когда мы учитываем контекстное время (например, месяц), такая высокая температура 25°C в течение зимы будет рассматриваться как аномалия.

Алгоритмы обнаружения точечных аномалий так же можно использовать для обнаружения контекстуальных и коллективных аномалий. Для этого нужно включить сам контекст как параметр. В вышеприведенном примере включение месяца как дополнительного параметра поможет обнаружить аномалию. Однако в более сложных сценариях может потребоваться одна или несколько новых параметров, чтобы преобразовать задачу определения контекстной аномалии в проблему обнаружения точечной аномалии. Преобразование поиска коллективной аномалии в поиск оди-

ночную может быть нетривиальной. Корреляция, агрегация и группировка используются для создания нового набора данных с другим представлением признаков[7]. Преобразование из задачи обнаружения коллективной аномалии в задачу обнаружения точечной аномалии требует глубоких знаний о наборе исходных данных и часто приводит к существенным искажениям при переводе данных в новый формат. Такое семантическое преобразование называется генерированием представления данных (англ. *data view generation*).

Таким образом можно сделать вывод, что многие задачи обнаружения аномалий требуют предварительной обработки данных перед передачей их на вход алгоритму. В противном случае можно получить формально верные, но фактически бесполезные результаты.

1.4.1 Нормализация данных

Когда мы получили предварительно обработанный датасет для поиска точечной аномалии, то последним шагом перед передачей в алгоритм, является нормализация данных. Нормализация данных предназначена для устранения зависимости от выбора единицы измерения и заключается в преобразовании диапазонов значений всех атрибутов к стандартным интервалам $([0,1]$ или $[-1,1])$ [?]. Нормализация данных направлена на придание всем атрибутам одинакового "веса".

1.4.1.1 Основные методы нормализация данных

а) Min-max нормализация заключается в применении к диапазону значений атрибута x линейного преобразования, которое отображает $[\min(x), \max(x)]$ в $[A, B]$.

$$x'_i = \tau(x_i) = \frac{x_i - \min(x)}{\max(x) - \min(x)} * (B - A) + A \quad (1.1)$$

$$x \in [\min(x), \max(x)] \Rightarrow \tau(x) \Rightarrow [A, B] \quad (1.2)$$

Min-max нормализация сохраняет все зависимости и порядок оригинальных значений атрибута. Недостатком этого метода является то, что выбросы могут сжать основную массу значений к очень маленькому интервалу

б) Z-нормализация основывается на приведении распределения исходного атрибута x к центрированному распределению со стандартным отклонением, равным 1 [?]

$$x'_i = \tau(x_i) = \frac{x_i - \bar{x}}{\sigma_x} \quad (1.3)$$

$$M[x'] = 1 \quad (1.4)$$

$$D[\bar{x}] = 0 \quad (1.5)$$

Метод полезен когда в данных содержат выбросы.

в) Масштабирование заключается в изменении длины вектора значений атрибута путем умножения на константу [?] .

$$x'_i = \tau(x_i) = \lambda * x_i \quad (1.6)$$

Длина вектора x уменьшается при $|\lambda| < 1$ и увеличивается при $|\lambda| > 1$

1.5 Неконтролируемые алгоритмы обнаружения аномалий

1.5.1 Вероятностный подход

Основная идея генеративного подхода заключается в использование вероятностного смесового моделирования данных. Предлагается подобрать такую вероятностную модель, из которой было получены нормированные данные. Такие модели обычно называются генеративными моделями, где для каждой точки(элемента данных) можем посчитать генеративную вероятность(или вероятность правдоподобия).Т.е. задача сводится к нахождению плотности распределения $p(x)$. Аномалиями при этом считаются точки(элементы набора данных), имеющую низкое правдоподобие. В качестве показателя аномальности выступает функция p .

Для построения генеративной модели нужно решить следующую задачу:

$$\prod_{x \in X_{norm}} p(x, \theta) \rightarrow \max_{\theta} \quad (1.7)$$

где X_{norm} - нормальные данные представленного набора данных $p(x, \theta) | \theta \in \omega$ - семейство плотностей вероятностей, параметризованные θ ;

Этот метод редко используется на практике, так как тяжело проверить полученную генеративную модель на адекватность, сложно убедиться в правильном выборе семейства смесовых распределений. Это связано с тем, что низкое значение функции правдоподобия может означать как и аномальное значение, так и неудачно подобранную модель. Этот метод применяется с опорой на априорную информацию, в случае когда можно проверить полученную модель на адекватность.

1.5.2 Линейный подход

Основной идеей линейного подхода является построение некой модели, характеризующей нормальные данные. Точки, которые значительно отклоняются от этой модели, считаются аномалиями.

Предполагается, что нормальные данные находятся в подпространстве пространства атрибутов данных(размер подпространства атрибутов данных равен размерности данных). В свою очередь, задача линейного метода - найти низкоразмерное подпространства, такие что, выборка данных этого подпространства значительно отличается от остальных точек пространства данных.

Одним из возможных вариантов решения является использование линейной регрессии. Выбирается одна из наблюдаемых переменных набора данных и относительно неё решается задача линейной регрессии оставшихся атрибутов. Итоговым ответом будет является усредненное значения показателя аномалии по всем атрибутам.

Алгоритмы, основанные на линейном подходе, требуют наличия линейной зависимости атрибутов данных.

1.5.3 Метрические методы

Метрические методы хорошо подходят в случае когда данные не размечены. Сложность вычисления прямо как пропорциональна размерности данных m , как и их количеству n . При росте набора данных наблюдается экспоненциальный рост сложности вычислений. Однако, эти методы хорошо проявляют себя на ограниченных наборах данных[?]. Следовательно такие методы как k -ближайших соседей (так же известный как обучение на основе примеров и описанный позднее) с нотацией ассимптотического роста $O(n^2m)$ недопустимы для наборов данных с большой размерности, если их размерность не может быть уменьшена.

Существуют много различных вариации алгоритма k -ближайших соседей для обнаружения аномалий, но все они основаны на вычислении некой метрики "расстояния до соседей такой как Евклидово расстояние или расстояние Махаланобиса. Евклидово расстояние задается следующей формулой:

$$\sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (1.8)$$

и является просто расстоянием между двумя точка, когда как расстояние Махаланобиса, задаваемое следующей формулой

$$\sqrt{(x - \mu)^T C^{-1} (x - \mu)} \quad (1.9)$$

вычисляет расстояние от точки до центра тяжести (μ), определяемого формулой коррелированных атрибутов, заданных матрицей ковариации (C). Расстояние Махаланобиса рассчитывается значительно дольше по сравнению с евклидовым по сравнению с евклидовым расстоянием для для больших объемов данных, поскольку оно требует пройти через весь набор данных, чтобы идентифицировать корреляции атрибутов.

Рамасвани в 2000 году предложил оптимизацию метода k -ближайших соседей (с англ. k -Nearest Neighbour - k -NN) в виде составления ранжированного списка потенциальных выбросов. Точка p является выбросом если не более $n - 1$ других точек в наборе данных имеют более высокий D_m (расстояние до m соседей), где m задается.

1.6 Виды аномалий

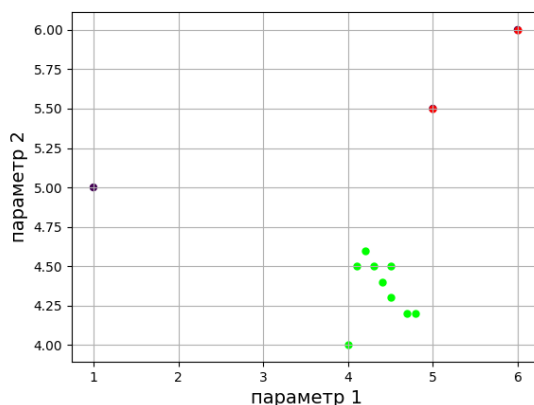


Рисунок 1.2 — Простой двумерный пример

В анализе данных есть два основных направления, которые занимаются поиском аномалий - это детектирование новизны и обнаружение выбросов. "Новый объект это так же объект, который отличается по своим свойствам от объектов выборки. Однако, в отличие от выброса, он его ещё нет в самой выборке и задача анализа сводится к его обнаружению при появлении. Например, если вы анализируете замеры уровня шума и отбрасываете слишком высокие или слишком низкие значения, то вы боретесь с выбросом. А если Вы создаёте алгоритм, который для каждого нового замера оценивает, насколько он похож на прошлые, и выбрасывает аномальные — вы "боретесь с новизной"[8]. Выбросы являются следствием:

- а) ошибок в данных
- б) неверно классифицированных объектов
- в) присутствием объектов других выборок
- г) намеренным искажением данных

На рисунке 1.1 можно увидеть желтые точки - выброс "слабом смысле". Они незначительно отклоняются от основных данных(зеленые точки). Красные же точки являются аномальными - выбросами "в сильном смысле они значительно отклоняются от основных данных. В данной работе будет изучаться вопрос нахождения "сильных выбросов"и критериев отличия сильного выброса от основных данных. В дальнейшем под словом "выброс"будет подразумеваться "сильный выброс а под аномалией - в выброс(выброс является частным случаем аномалии). Понятие аномалии зачастую интерпретируют по-разному в зависимости от характера данных. Обычно аномалией называют некоторое отклонение от нормы. Это определение нуждается в формальном уточнении.

1.7 Постановка задачи

Кстати, про картинки. Во-первых, для фигур следует использовать `[ht]`. Если и после этого картинки вставляются «не по ГОСТ», т.е. слишком далеко от места ссылки, — значит у вас в РПЗ **слишком мало текста!** Хотя и ужасный параметр `!ht` у окружения `figure` тоже никто не отменял, только при его использовании документ получается страшный, как в ворде, поэтому просьба так не делать по возможности.

1.8 Существующие подходы к созданию всячины

Известны следующие подходы...

- а) Перечисление с номерами.
- б) Номера первого уровня. Да, ГОСТ требует именно так — сначала буквы, на втором уровне — цифры. Чуть ниже будет вариант «нормальной» нумерации и советы по её изменению. Да, мне так нравится: на первом уровне выравнивание элементов как у обычных абзацев. Проверим теперь вложенные списки.
 - 1) Номера второго уровня.
 - 2) Номера второго уровня. Проверяем на длиииинной-предлиииииииинной строке, что получается.... Сойдёт.
- в) По мнению Лукьяненко, человеческий мозг старается подвести любую проблему к выбору из трех вариантов.
- г) Четвёртый (и последний) элемент списка.

Теперь мы покажем, как изменить нумерацию на «нормальную», если вам этого захочется. Пара команд в начале документа поможет нам.

- 1) Изменим нумерацию на более привычную...
- 2) ... нарушим этим гост.
 - а) Но, пожалуй, так лучше.

В заключение покажем произвольные маркеры в списках. Для них нужен пакет **enumerate**.

- 1. Маркер с арабской цифрой и с точкой.
- 2. Маркер с арабской цифрой и с точкой.
 - I. Римская цифра с точкой.
 - II. Римская цифра с точкой.

В отчётах могут быть и таблицы — см. табл. 1.1 и 1.2. Небольшая таблица делается при помощи **tabular** внутри **table** (последний полностью аналогичен **figure**, но добавляет другую подпись).

Таблица 1.1 — Пример короткой таблицы с коротким названием

Тело	F	V	E	$F + V - E - 2$
Тетраэдр	4	4	6	0
Куб	6	8	12	0
Октаэдр	8	6	12	0
Додекаэдр	20	12	30	0
Икосаэдр	12	20	30	0
Эйлер	666	9000	42	$+\infty$

Для больших таблиц следует использовать пакет **longtable**, позволяющий создавать таблицы на несколько страниц по ГОСТ.

Для того, чтобы длинный текст разбивался на много строк в пределах одной ячейки, надо в качестве ее формата задавать `p` и указывать явно ширину: в мм/дюймах (`110mm`), относительно ширины страницы (`0.22\textwidth`) и т.п.

Можно также использовать уменьшенный шрифт — но, пожалуйста, тогда уж во **всей** таблице сразу.

Таблица 1.2 — Пример длинной таблицы с длинным названием на много длинных-длинных строк

Вид шума	Громкость, дБ	Комментарий
Порог слышимости	0	
Шепот в тихой библиотеке	30	Конечно, это было до эпохи мобильных (внутри машины)
Обычный разговор	60-70	
Звонок телефона	80	
Уличный шум	85	
Гудок поезда	90	
Шум электрички	95	
Порог здоровой нормы	90-95	Длительное пребывание на более громком шуме может привести к ухудшению слуха
Мотоцикл	100	(модель бензокосилки) (Доом в целом вреден для здоровья)
Power Mower	107	
Бензопила	110	
Рок-концерт	115	

Продолжение на след. стр.

Продолжение таблицы 1.2

Порог боли	125	feel the pain
Клепальный молоток	125	(автор сам не знает, что это)
Порог опасности	140	Даже кратковременное пребывание на шуме большего уровня может привести к необратимым последствиям
Реактивный двигатель	140 180	Необратимое полное повреждение слуховых органов
Самый громкий возможный звук	194	Интересно, почему?..

2 Конструкторский раздел

2.1 Оценка качества изображений

2.1.1 PSNR

Пиковое отношение сигнала к шуму (англ. peak signal-to-noise ratio) - соотношение между максимумом возможного значения сигнала и мощностью шума, искажающего значения сигнала.[?]

$$PSNR = 20 \log_{10} \frac{MAX_i}{\sqrt{MSE}} \quad (2.1)$$

Где MAX_i - это максимальное значение, принимаемое пикселем изображения, MSE - среднеквадратичное отклонение. Для двух монохромных изображений I и K размера $m \times n$, одно из которых считается зашумленным приближением другого, вычисляется так:

$$MSE = \frac{1}{m * n} \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} |I(i,j) - K(i,j)|^2 \quad (2.2)$$

2.1.2 SSIM

Индекс структурного сходства (SSIM от англ. structure similarity) — метод измерения схожести между двумя изображениями путем полного сопоставления. SSIM-индекс является развитием традиционных методов, таких как PSNR (peak signal-to-noise ratio) и метод среднеквадратичной ошибки MSE, которые оказались несовместимы с физиологией человеческого восприятия.

Отличительной особенностью метода, в отличие от MSE и PSNR, является то, что он учитывает «восприятие ошибки» благодаря учёту структурного изменения информации. Идея заключается в том, что пиксели имеют сильную взаимосвязь, особенно когда они близки пространственно. Данные зависимости несут важную информацию о структуре объектов и о сцене в целом. Особенностью является, что SSIM всегда лежит в промежутке от -1 до 1, причем при его значении равном 1, означает, что мы имеем две одинаковые картинки. Общая формула имеет вид

$$SSIM(x,y) = \frac{(2\mu_x\mu_y + c_1)(\sigma_{xy} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x^2 + \sigma_y^2 + c_2)} \quad (2.3)$$

Тут μ_x среднее значение для первой картинки, μ_y для второй, σ_x среднеквадратичное отклонение для первой картинки, и соответственно σ_y для второй, σ_{xy} это уже ковариация. Она находится следующим образом:

$$\sigma_{xy} = \mu_{xy} - \mu_x\mu_y \quad (2.4)$$

c_1 и c_2 - поправочные коэффициенты, которые нужны вследствие малости знаменателя.

$$c_1 = (0,01 * d)^2 \quad (2.5)$$

$$c_2 = (0,03 * d)^2 \quad (2.6)$$

d - количество цветов, соответствующих данной битности изображения. Для подтверждения или опровержения вышеописанных гипотезы реализуются несколько алгоритмов случайного распределения, алгоритм Флойда-Стейнберга, модификации на основе алгоритма Флойда-Стейнберга. Результаты работы алгоритмов сравниваются по времени работы, по количеству затрачиваемой памяти, а так же по SSIM и PSNR.

2.1.3 Алгоритм случайного распределения

$P(x,y)$ - цвет конкретного пикселя

Листинг 2.1 — Алгоритм случайного распределения

```

1 for x in range(height):
2     for y in range(weight):
3         if P(x,y) > 127:
4             P(x,y) = 255
5         else:
6             P(x,y) = 0

```

2.2 Виды случайных распределений

2.2.1 Белый шум

Белый шумом называют сигнал с равномерной спектральной плотностью на всех частотах и дисперсией, равной бесконечности. Является стационарным случайным процессом. В качестве сигнала в задаче дизеринга рассматривается последовательность последовательность чисел, получаемых от генератора случайных чисел.

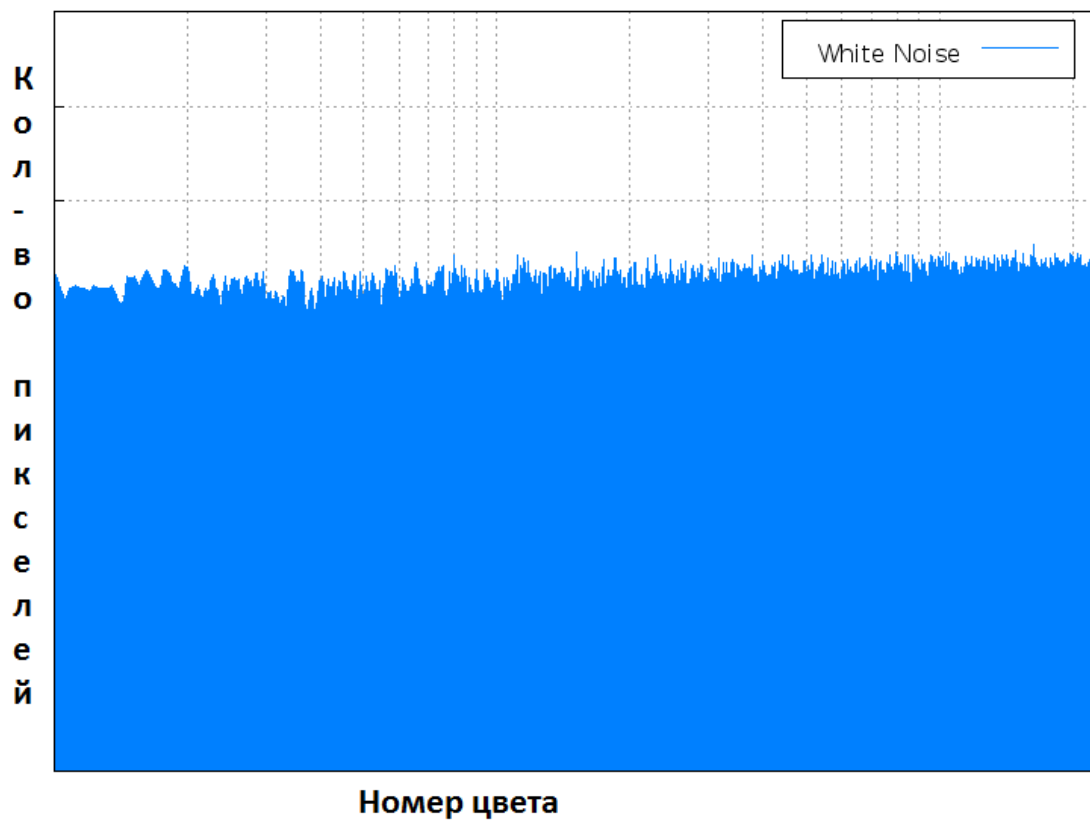


Рисунок 2.1 — Диаграмма белого шума

2.2.2 Коричневый шум

Спектральная плотность коричневого шума пропорциональна $1/f^2$, где f — частота. Это означает, что на низких частотах шум имеет больше энергии, чем на высоких. То есть пикселей темных цветов больше, чем пикселей светлого цвета. Применение фильтра коричневого шума в целом затемняет получаемое изображение.

Листинг 2.2 — Получение коричневого шума

```

1 def smoother(noise):
2     output = []
3     for i in range(len(noise) - 1):
4         output.append(0.5 * (noise[i] + noise[i+1]))
5 return output

```

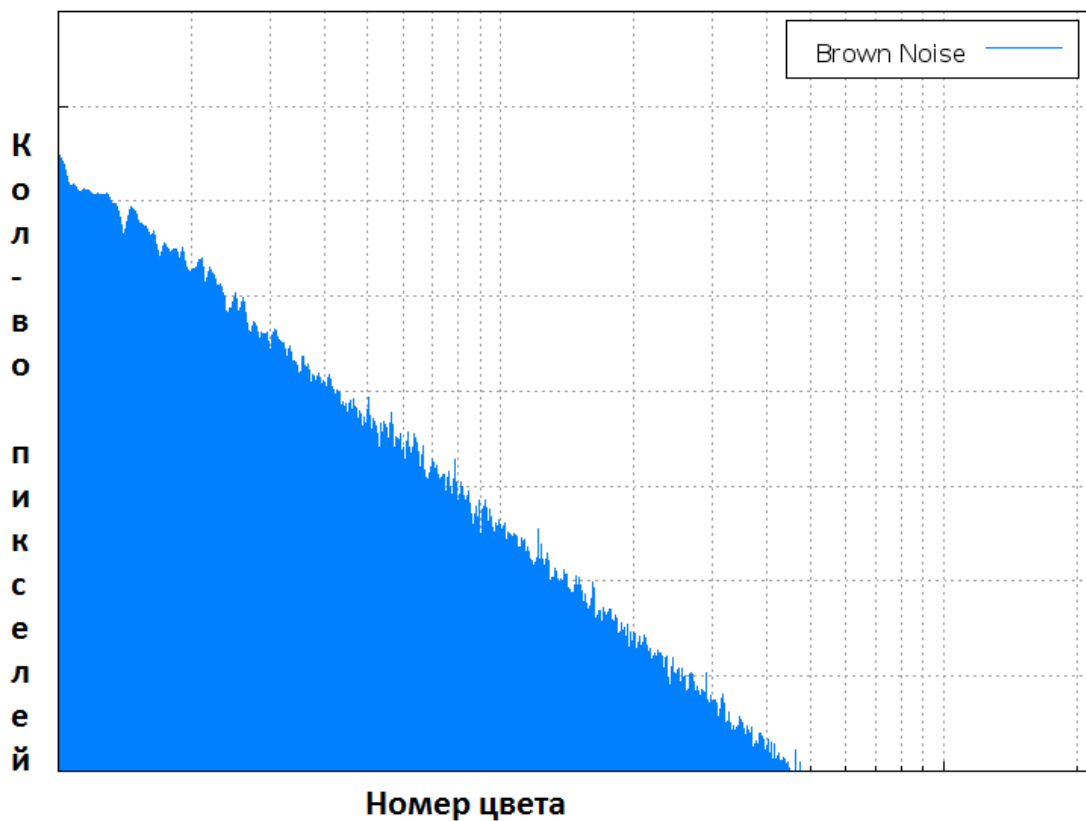


Рисунок 2.2 — Диаграмма красного шума

2.2.3 Гауссовский шум

Гауссовский шум - шум, имеющий функцию плотности вероятности (PDF), равную нормальному распределению, которое также известно как гауссово распределение.

$$p_g(z) = \frac{1}{\sigma\sqrt{2 * \pi}} e^{-\frac{(z-\mu)^2}{2\sigma^2}} \quad (2.7)$$

z - количество цветов, μ -среднее значение, σ - стандартное отклонение.

2.2.4 Фиолетовый шум

Фиолетовый шум представляет собой противоположность между коричневому шуму. Получение его аналогично получению коричневого шума. Применение фильтра фиолетового шума в целом засветляет получаемое изображение.

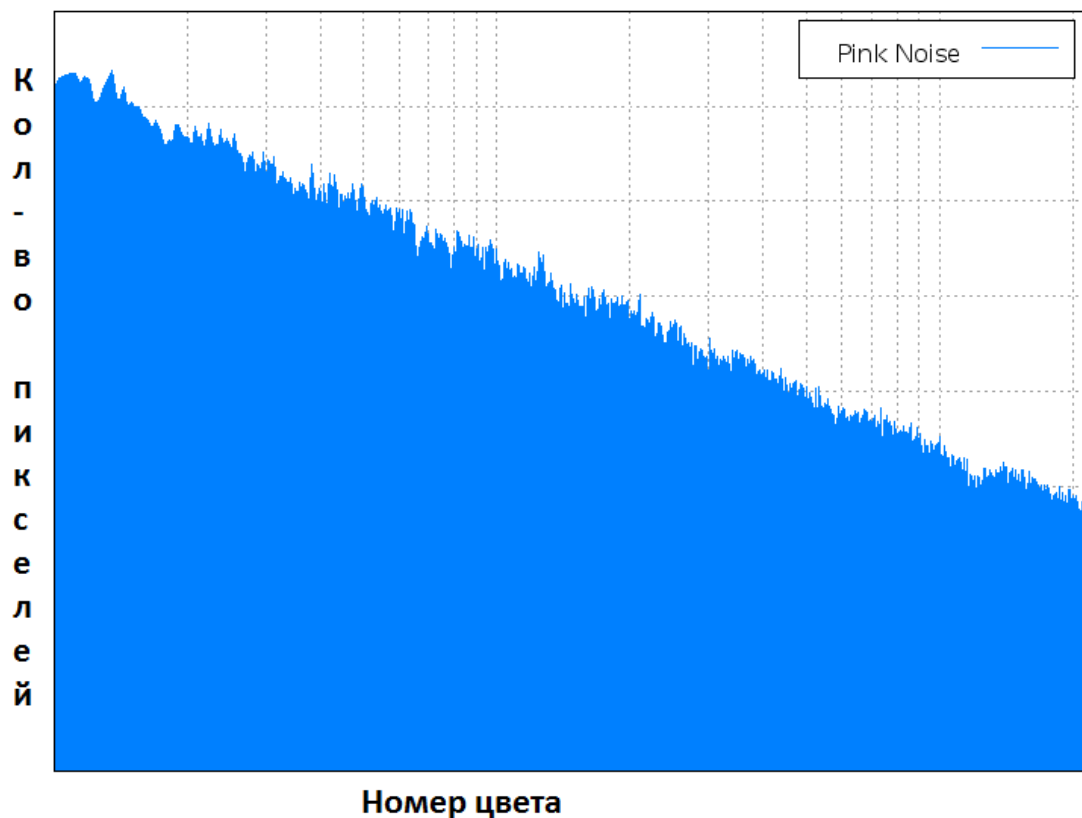


Рисунок 2.3 — Диаграмма розового шума

Листинг 2.3 — Получение розового шума

```

1 def rougher(noise):
2     output = []
3     for i in range(len(noise) - 1):
4         output.append(0.5 * (noise[i] - noise[i+1]))
5     return output

```

2.2.5 Розовый и синий шумы

Розовый и синий шумы представляют собой "промежуточные" шумы. Изображение с розовым шумом темнее изображения с белым шумом, но светлее изображения с коричневым шумом. Изображение с синим шумом светлее изображения с белым шумом, но темнее чем изображение с фиолетовым шумом. Их получение аналогично получению со коричневого и фиолетового шумов.

2.3 Алгоритм Байера

Листинг 2.4 — Алгоритм Флойда-Стейнберга

```
1 For each pixel, Input, :  
2   Factor = ThresholdMatrix[xcoordinate % X][ycoordinate % Y];  
3   Attempt = Input + Factor * Threshold  
4   Color = FindClosestColorFrom(Palette, Attempt)
```

2.4 Алгоритм Флойда-Стейнберга

Рассмотрим более детально алгоритм Флойда-Стейнберга. $P(x,y)$ - цвет пикселя в точке x,y $I(x,y)$ - предполагаемый цвет пикселя с учетом ошибки (действительное число) Следует отметить, что сумма "долей" рассеивания ошибки

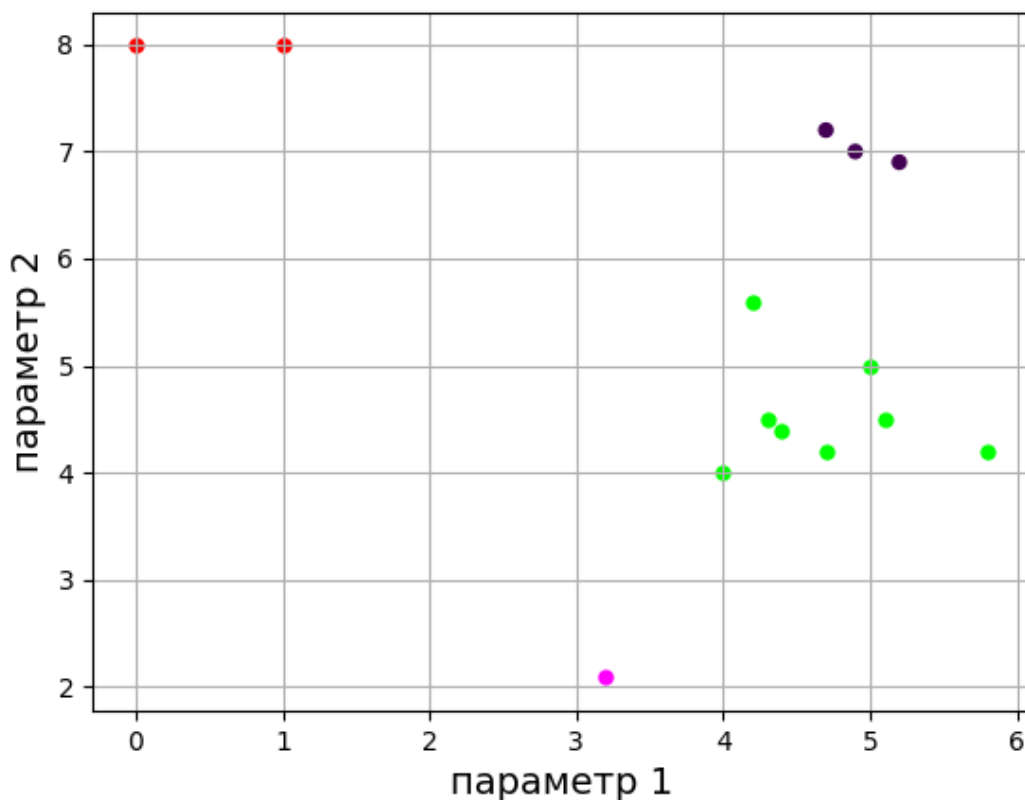


Рисунок 2.4 — Схема рассеивания ошибки

и для других алгоритмов упорядоченного дithering будет равна единице.

Листинг 2.5 — Алгоритм Флойда-Стейнберга

```
1 for x in range(width):  
2   for y in range(height):  
3      $P(x,y) = \text{trunc}(I(x,y) + 0.5)$ 
```

```

4      e = I(x,y) - P(x,y)
5      I(x,y-1) += 7/16*e
6      I(x+1, y-1) += 3/16*e
7      I(x+1, y) += 5/16*e
8      I(x+1, y+1) += 1/16*e

```

2.5 Ложный алгоритм Флойда-Стейнберга

Отличительной особенностью ложного алгоритма Флойда-Стейнберга является по пикселям изображения справа-налево, а не слева-направо как в большинстве других алгоритмов.

Листинг 2.6 — Ложный алгоритм Флойда-Стейнберга

```

1  for x in range(width):
2      for y in range(height):
3          P(x,y) = trunc(I(x,y)+0.5)
4          e = I(x,y) - P(x,y)
5          P(x,y) = trunc(I(x,y)+0.5)
6          e = I(x,y) - P(x,y)
7          I(x,y-1) += 3/8*e
8          I(x-1, y-1) += 2/8*e
9          I(x-1, y) += 3/8*e

```

2.6 Алгоритм Джарвиса,Джунка и Нинка

Листинг 2.7 — Алгоритм Джарвиса,Джунка и Нинка

```

1  for x in range(width):
2      for y in range(height):
3          P(x,y) = trunc(I(x,y)+0.5)
4          e = I(x,y) - P(x,y)
5          P(x,y) = trunc(I(x,y)+0.5)
6          e = I(x,y) - P(x,y)
7          I(x,y-1) += 3/48*e
8          I(x+1,y+1) += 5/48*e
9          I(x,y+1) += 7/48*e
10         I(x-1,y+1) += 5/48*e
11         I(x-2,y+1) += 3/48*e
12         I(x-1,y+1) += 5/48*e
13         I(x+2,y+2) += 1/48*e
14         I(x+1,y+2) += 3/48*e
15         I(x, y+2) += 5/48*e

```

```

16      I(x-1,y+2) += 3/48*e
17      I(x-2,y+2) += 1/48*e
18      I(x+1,y)   += 7/48*e
19      I(x+2,y)   += 5/48*e

```

2.7 Первый алгоритм Юлиомы

Листинг 2.8 — Первый алгоритм Юлиомы

```

1 For each pixel, Input, in the original picture:
2   Factor = ThresholdMatrix[xcoordinate % X][ycoordinate % Y];
3   Make a Plan, based on Input and the Palette.
4
5   If Factor < Plan.ratio,
6     Draw pixel using Plan.color2
7   else,
8     Draw pixel using Plan.color1

```

Листинг 2.9 — План нахождения цвета пикселя

```

1 SmallestPenalty = 10^99
2 For each unique combination of two colors from the palette, Color1
   and Color2:
3   For each possible Ratio, 0 .. (X*Y-1):
4     Mixed = Color1 + Ratio * (Color2 - Color1) / (X*Y)
5     Penalty = Evaluate the difference of Input and Mixed.
6
7     If Penalty < SmallestPenalty,
8       SmallestPenalty = Penalty
9     Plan = { Color1, Color2, Ratio / (X*Y) }1

```

Листинг 2.10 — Вспомогательная матрица первого Алгоритма Юлиомы

```

1 [0/64, 48/64, 12/64, 60/64, 3/64, 51/64, 15/64, 63/64,
2 32/64, 16/64, 44/64, 28/64, 35/64, 19/64, 47/64, 31/64,
3 8/64, 56/64, 4/64, 52/64, 11/64, 59/64, 7/64, 55/64,
4 40/64, 24/64, 36/64, 20/64, 43/64, 27/64, 39/64, 23/64,
5 2/64, 50/64, 14/64, 62/64, 1/64, 49/64, 13/64, 61/64,
6 34/64, 18/64, 46/64, 30/64, 33/64, 17/64, 54/64, 29/64,
7 10/64, 58/64, 6/64, 54/64, 9/64, 57/64, 5/64, 53/64,
8 42/64, 26/64, 38/64, 22/64, 41/64, 25/64, 37/64, 21/64]

```


3 Технологический раздел

3.1 Выбор языка программирования

Для реализации данных алгоритмов был выбран язык C++. Данный язык был обоснован следующими причинами: Причины:

- 1) Компилируемый язык со статической типизацией.
- 2) Сочетание высокоуровневых и низкоуровневых средств.
- 3) Реализация ООП.
- 4) Наличие удобной стандартной библиотеки шаблонов

3.2 Выбор вспомогательных библиотек

Для реализации программы была выбрана библиотека Qt.

- 1) Широкие возможности работы с изображениями, в том числе и попиксельно
- 2) Наличии более функциональных аналогов стандартной библиотеки шаблонов в том числе для разнообразных структур данных

Так же были использованы библиотеки ImageMagick(для конвертации изображения в ограниченную цветовую палитру), OpenMP(многопоточность), OpenGL(работа с шейдерами)

3.2.1 Диаграмма классов

Для реализации различных алгоритмов была разработана следующая структура классов. Был создан абстрактный класс дизеринга Dithering с интерфейсом, наследуемом в дочерних классах. Так же была введена система менеджеров: DitherManager, MetricsManager, DataManager и MainManager.

4 Исследовательский раздел

4.1 Время дизеринга различных алгоритмов

Рассмотрим время работы различных алгоритмов для различных размеров изображения.

	Размер, пиксели	Время, мкс
White noise	133x90	862
Blue noise	133x90	930
Brown noise	133x90	934
Violet noise	133x90	937
Pink noise noise	133x90	930
Floyd-SD	133x90	1200
F. Floyd-SDe	133x90	1093
JJN	133x90	1909
White noise	458x458	15735
Blue noise	458x458	19374
Brown noise	458x458	19432
Violet noise	458x458	18787
Pink noise noise	458x458	18129
Floyd-SD	458x458	27173
F. Floyd-SDe	458x458	26424
JJN	458x458	47201
White noise	458x458	194376
Blue noise	458x458	200577
Brown noise	458x458	208400
Violet noise	458x458	251294
Pink noise noise	458x458	258775
Floyd-SD	458x458	251294
F. Floyd-SDe	458x458	387104
JJN	458x458	857481

Из рассмотрения вынесены алгоритм Юлиомы в вследствие того, что он значительно медленней других алгоритмов(2732568 мкс для изображения 113x90) в и алгоритм Байера, реализованный при помощи шейдеров, вследствие того, что он не укладывается в рамки требуемой палитры (при этом он работает очень быстро 64 мс для изображении 640x480).

4.2 Качество получаемого изображения

	PSNR	SSIM
White noise	33.2894	0.914778
Blue noise	36.1756	0.971626
Brown noise	33.32370	0.915767
Violet noise	37.63480	0.984574
Pink noise	36.4484	0.974718
Floyd-SD	37.0553	0.979173
F. Floyd-SDe	36.8401	0.976452
JJN	37.30740	0.981688
Yliouma	36.2359	0.967796
Without dithering	37.6348	0.984574

Несмотря на то, что некоторые сложные алгоритмы дизеринга диффузии ошибок обещают получения хорошего качества изображений, некоторые алгоритмы случайного дизеринга на конкретных изображениях дают лучший результат. Для того чтобы получить наилучший результат дизеринга, следует проанализировать результаты дизеринга нескольких изображений и выбрать среди них наилучшее. Так же следует отметить некоторую необъективность метрик: результат метрик не всегда совпадает с человеческим восприятием картинки.

4.3 Размер получаемого изображения

	Разрешение, пикс	Размер, кб	Исх. раз., кб
White noise	900x675	186	2373 bmp, 1779 png
Blue noise	900x675	135	
Brown noise	900x6750	186	
Violet noise	900x675	98	
Pink noise	900x675	1158	
Floyd-SD	900x675	1273	
F. Floyd-SDe	900x675	143	
JJN	900x675	117	
White noise	3984x32355	3431	50344 bmp, 37758 png
Blue noise	3984x3235	2570	
Brown noise	3984x3235	3432	
Violet noise	3984x3235	1950	
Pink noise	3984x3235	2406	
Floyd-SD	3984x32355	3605	
F. Floyd-SDe	3984x3235	4269	
JJN	3984x3235	3716	

Из вышеприведенной таблицы, можно заметить, размер изображения после дизеринга значительно уменьшается, достигается выигрыш в размере изображения до 15 раз, в зависимости от исходного контейнера изображения и выбранного способа дизеринга.

Заключение

В данной работе были реализованы различные алгоритмы дизеринга, было произведено сравнение и анализ этих алгоритмов. Программа позволяет получить изображение схожего визуального качества при значительном уменьшении размера. Был получен вывод о том, что для различных целей следует использовать различные алгоритмы дизеринга, универсального алгоритма дизеринга не существует. Программа не привязана к какой-то конкретной операционной системе и может быть скомпилирована и запущена на всех популярных ОС.

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. *machinelearning.ru*/. Выборка. <https://goo.gl/7gjJ6p>.
2. определения, ГОСТ 20886-85: Организация данных в системах обработки данных. Термины и. <http://www.gostrf.com/normadata/1/4294832/4294832686.pdf>.
3. Википедия. https://en.wikipedia.org/wiki/Data_set.
4. F.E., Grubbs. Procedures for Detecting Outlying Observations in Samples. Technometrics / Grubbs F.E. — 1969.
5. Moya M.M., Hush D.R. Network Constraints and Multi-objective Optimization for One-class Classification. Neural Networks / Hush D.R. Moya M.M. — 1996.
6. Chandola V Banerjee A, Kumar V. Anomaly Detection: A Survey / Kumar V. Chandola V, Banerjee A. — ACM Computing, 2009.
7. Goldstein M, Uchida S. Behavior Analysis Using Unsupervised Anomaly Detection / Uchida S. Goldstein M. — The 10th Joint Workshop on Machine Perception and Robotics, 2014.
8. Дьяконов, Александр. Поиск аномалий (Anomaly Detection) / Александр Дьяконов. — 2017. <https://goo.gl/Z43Ne9>.