



Министерство образования и науки Российской Федерации
Федеральное государственное бюджетное образовательное учреждение
высшего образования
«Московский государственный технический университет
имени Н.Э. Баумана
(национальный исследовательский университет)
(МГТУ им. Н.Э. Баумана)

ФАКУЛЬТЕТ «ФУНДАМЕНТАЛЬНЫЕ НАУКИ»

КАФЕДРА «ПРИКЛАДНАЯ МАТЕМАТИКА»

РАСЧЕТНО-ПОЯСНИТЕЛЬНАЯ ЗАПИСКА

К ВЫПУСКНОЙ КВАЛИФИКАЦИОННОЙ РАБОТЕ

НА ТЕМУ:

НОВЫЕ РАСЧЕТНЫЕ СХЕМЫ ВИХРЕВЫХ МЕТОДОВ
И ПРОГРАММНЫЙ КОМПЛЕКС
ДЛЯ МАТЕМАТИЧЕСКОГО МОДЕЛИРОВАНИЯ
ОБТЕКАНИЯ ПРОФИЛЕЙ

Студент ФН2-81
(Группа)

Рятина Е.П.
(И.О.Фамилия)

Руководитель ВКР

Марчевский И.К.
(И.О.Фамилия)

Консультант

(И.О.Фамилия)

Консультант

(И.О.Фамилия)

Нормоконтролер

Лукашин М.М.
(И.О.Фамилия)

2018 г.

Содержание

ВВЕДЕНИЕ	3
1 ЗАДАЧА ДЕТЕКЦИИ ИЗОБРАЖЕНИЙ	5
2 ПОДХОДЫ К ЗАДАЧЕ ДЕТЕКТИРОВАНИЯ	7
2.1 Метод скользящего окна	7
2.2 Метод предложения регионов	8
2.3 Детекция за один проход	9
3 НЕЙРОСЕТЕВЫЕ АРХИТЕКТУРЫ В ДЕТЕКТИРОВАНИИ	12
3.1 Архитектура Transformer	12
3.2 Архитектура YOLO	18
4 МОДЕЛИРОВАНИЕ НЕЙРОННОЙ СЕТИ	21
4.1 Описание используемого датасета	21
4.2 Предобработка датасета для сети ViT	22
4.3 Венгерский алгоритм	23
4.4 Функция потерь сети ViT и сопоставление истинных и предсказанных меток	26
4.5 Метрики сети ViT	27
4.6 Слои и структуры сети ViT	30
4.7 Реализация сети ViT	33
4.8 Предобработка датасета для сети YOLO	34
5 РЕЗУЛЬТАТЫ МОДЕЛИРОВАНИЯ	35
5.1 Сеть ViT	35
5.2 Сеть YOLO	37
ЗАКЛЮЧЕНИЕ	38
Список использованных источников	39

ВВЕДЕНИЕ

Задача детекции изображений, связанная с выделением и классификацией объектов на изображении, является важной областью компьютерного зрения.

Одним из ранних методов детекции является метод Виолы-Джонса, основанный на каскадных классификаторах, который был популярен в прошлом, но имеет некоторые недостатки, такие как невысокая точность обнаружения на изображениях в условиях изменения масштаба и поворота объектов.

Современные подходы для детекции объектов на изображениях основаны на использовании нейронных сетей.

R-CNN – это первый метод детекции объектов на изображении, основанный на нейронных сетях. Он использует сверточные нейронные сети для извлечения признаков из областей изображения, а затем классифицирует и локализует объекты. Однако сеть R-CNN работает медленно из-за необходимости обрабатывать каждую область изображения отдельно.

Fast R-CNN – улучшенная версия R-CNN, которая использует более эффективный способ обработки областей изображения с помощью RoI-pooling. Он также улучшает скорость и точность детекции объектов на изображении.

Faster R-CNN – это метод, основанный на более эффективном процессе выделения областей изображения и RoI-pooling, называемом Region Proposal Network (RPN). Он может автоматически генерировать регионы интереса, что увеличивает скорость детекции.

С другой стороны, YOLO (You Only Look Once) – это метод, который выделяет объекты и классифицирует их на изображении в один прогон. Он использует сеть, которая делит изображение на сетку ячеек, каждая из которых ответственна за предсказание объектов в этой ячейке.

SSD (Single Shot MultiBox Detector) – похож на YOLO, это ещё один метод одноступенчатого обнаружения объектов, который предсказывает класс и ограничивающие рамки в одну попытку. Однако он использует другой подход к разделению изображения на набор якорных коробок с заранее определенными соотношениями сторон, масштабами и позициями.

Трансформеры – это архитектура нейронных сетей, которая изначально

была разработана для обработки текстов, но также может быть использована для обработки изображений с помощью «механизма внимания». Использование трансформеров позволяет задействовать глобальный контекст изображения и агрегировать информацию из различных групп областей изображения. Примерами нейронных сетей, основанных на трансформерах, которые были успешно применены в задаче детекции объектов на изображении, являются DeTr и ViT.

DeTr – это одна из первых нейронных сетей, полностью основанных на трансформерах. Она показала превосходство над другими методами детекции объектов на изображении и демонстрирует высокую точность обнаружения.

ViT – это нейронная сеть, которая была разработана для обработки изображений в виде сетки пикселей, и основана на трансформере. Она также показывает хорошую точность при решении задачи детекции объектов на изображениях.

1 ЗАДАЧА ДЕТЕКЦИИ ИЗОБРАЖЕНИЙ

Популярная в области компьютерного зрения задача детекции заключается в нахождении на изображениях и видеозаписях различных объектов, а также локализации их местоположения. Преимущественно данная область ограничивается прямоугольником, допускающим определенную параметризацию – изменение угла, ширины, высоты и т.д. В соответствии с формой детектируемого объекта вместо прямоугольника может быть выбрана другая фигура, например, эллипс или окружность.

Формально задачу можно представить следующим образом:

- 1) пусть задан набор s исследуемых классов, которым принадлежат объекты;
- 2) на вход подаются цветные (черно-белые) изображения $I \in R^{N \times M \times C}$;
- 3) требуется получить выходной набор обнаруженных объектов

$$Obj = \{Obj_i | i = 1, \dots, p\},$$

где Obj_i – структура, которая описывает местоположение объекта (bounding box), а также класс $c_i \in [1, s]$, которому он принадлежит.

Достоверность детекции зависит от нескольких факторов. Так, чтобы улучшить точность локализации, необходимо добавить угол поворота объекта. Описание границы объекта или области может способствовать более точному описанию его положения.

Детектируемых объекты могут быть условно разделены на две группы. К первой категории относятся объекты, обнаружение которых может быть ограничено прямоугольником – люди, животные, материальные предметы и т.д. Вторую группу образуют «неисчисляемые» материалы, которые не имеют конкретной формы, размера, а их описание задано какой-либо текстурой – вода, небо, земля и т.д.

Ограничивающая рамка может быть задана с использованием четырех пространственных координат в двух видах:

- 1) классический формат – задание координат (x_{min}, y_{max}) левого верхнего и координат (x_{max}, y_{min}) правого нижнего углов;

2) центрированное задание – определение координат центра (c_x, c_y) прямогоугольника, а также высоты h и ширины w .

В настоящее время задача обнаружения объектов является довольно частой в сфере автотранспортного движения в связи с развивающимися технологиями беспилотного управления и контроля трафика. Большое число существующих методов требует изучения их применения и усовершенствования для реализации задач детектирования в режиме реального времени.

2 ПОДХОДЫ К ЗАДАЧЕ ДЕТЕКТИРОВАНИЯ

Наиболее популярным направлением в области распознавания и детектирования объектов являются искусственные нейронные сети. В тех случаях, когда заранее известны возможные типы объектов, хорошее качество решения задачи показывают алгоритмы, использующие сверточные нейронные сети (CNN). Использование подобного семейства сетей позволяет обнаруживать объекты разных классов и пространственных размеров.

Можно выделить два основных подхода применения CNN:

- 1) ряд задач может быть реализован посредством применения готовых или уже предобученных архитектур; они могут детектировать объекты, на которых производилось обучение;
- 2) для других более частных задач может быть применена новая база изображений с целью дообучения модели, для которой имеется предобученная архитектура, в соответствии с поставленной задачей; данных случай предполагает настройку весов, в основном, исключительно для последних полносвязных слоев.

Изучение и развитие детекторов (алгоритмов, позволяющих выделять объекты), а также хорошие показатели качества их работы привело к возникновению большого разнообразия нейросетевых архитектур.

2.1 Метод скользящего окна

Первым, однако наиболее затратным с точки зрения вычислений, считается метод скользящего окна. При таком подходе происходит пошаговый проход по всему изображению, во время которого происходит извлечение определенного региона, а также классификация объекта, попавшего в него.

Можно предположить, что простым вариантом послужит построение модели обнаружения объектов на основе модели, реализующей классификацию. Имея классификатор, который показывает хорошие результаты, можно обнаружать объекты, сдвигая «окно» по изображению и определяя, относится ли часть, ограниченная окном, к желаемому классу. Однако такой подход ограничивается двумя проблемами:

1) какой размер окна считать оптимальным, чтобы объект точно содержался в нем; даже однотипные объекты (к примеру, теннисный мяч, футбольный мяч, относящиеся к классу мяч) могут иметь сильно отличающиеся пространственные размеры;

2) соотношение сторон ограничивающего прямоугольника – сильно отличающиеся по форме объекты могут иметь разные пропорции сторон.

Решение описанных проблем требовало бы перебора различных скользящих окон, отличающихся размерами и формой. Это приводит к значительным вычислительным затратам.

2.2 Метод предложения регионов

Другим подходом в задаче детектирования является метод «предложений регионов» (region proposals), сочетающий в себе два алгоритма. Первый из них создает множество предположений – это регионы, в которых может находиться объект. Второй алгоритм производит классификацию объекта.

Данный метод предполагает сначала нахождение некоторых областей на изображении, а затем объединение их в группу, согласно определенной иерархии. Такие области впоследствии объединяются согласно различным цветовым пространствам, а также факторам схожести. В результате имеется несколько предложений регионов, в которых может быть обнаружен объект за счет слияния нескольких небольших регионов.

Главным направлением развития глубоких нейронных сетей стало обнаружение объектов в реальном времени. Здесь можно отметить сеть R-CNN (Region Based). Данный подход характеризуется разпознаванием не по всему входному изображению, а в локализованной области. Основные шаги в реализации детектирования с использованием R-CNN моделей [1]:

1) генерация областей интереса для входного изображения – предложение регионов, которые предположительно могут содержать объект;

2) формирования карты признаков – происходит масштабирование сгенерированных на предыдущем этапе областей до размеров, соответствующих архитектуре нейронной сети CNN, эти данные поступают на вход CNN;

3) классификация – каждый объект определяется к своему классу посредством применения полученного вектора признаков с использованием метода опорных векторов.

Прежде чем осуществить классификацию для каждого предложенного региона используется алгоритм подавления немаксимумов. Его суть заключается в том, что исключительно локальные максимумы могут быть рассмотрены в качестве контура объекта. Применение такого алгоритма обусловлено необходимостью избавления от дублирующих областей на входном изображении.

Оценка качества произведенной классификации производится с помощью часто применяемой в задачах компьютерного зрения метрики Intersection over Union (IoU). Данное значение равняется отношению площади пересечения прямоугольника (области интереса), который был получен в ходе обнаружения, и прямоугольника из разметки к площади их объединения:

$$IoU = \frac{B_1 \cup B_2}{B_1 \cup B_2}. \quad (2.1)$$

Среди минусов использования R-CNN сетей можно отметить их относительно медленную работу, а также то, что в ходе алгоритма рассматриваются лишь отдельные области, а не изображение в целом.

2.3 Детекция за один проход

Еще одну группу методов детекции составляют такие, в которых локализация и классификация осуществляются одновременно, то есть алгоритм проходит один раз. Данный подход получил название «Single Shot Detection» (SSD, «детекция за один проход»).

SSD-алгоритм представляет собой комбинацию двух компонентов – базовой модели и «головы». Первая модель, базовая, является сетью для классификации (средством для извлечения признаков на изображении), которая была предварительно обучена. Преимущественно ею является сеть наподобие ResNet, обученная на ImageNet. При этом из данной сети удален последний полносвязный слой. В результате имеется глубокая нейронная сеть, умеющая находить семантическое значение для изображения, поступающего на вход,

для которого сохранена пространственная структура, несмотря на более низкое разрешение.

Головная часть SSD задается одним или несколькими слоями, добавленными к базовой модели. Данные выхода могут быть интерпретированы как ограничивающие прямоугольники и классы объектов в пространственном расположении активаций конечных слоев. Архитектура сверточной нейронной сети с SSD-детектором схематично представлена на рис. 1.

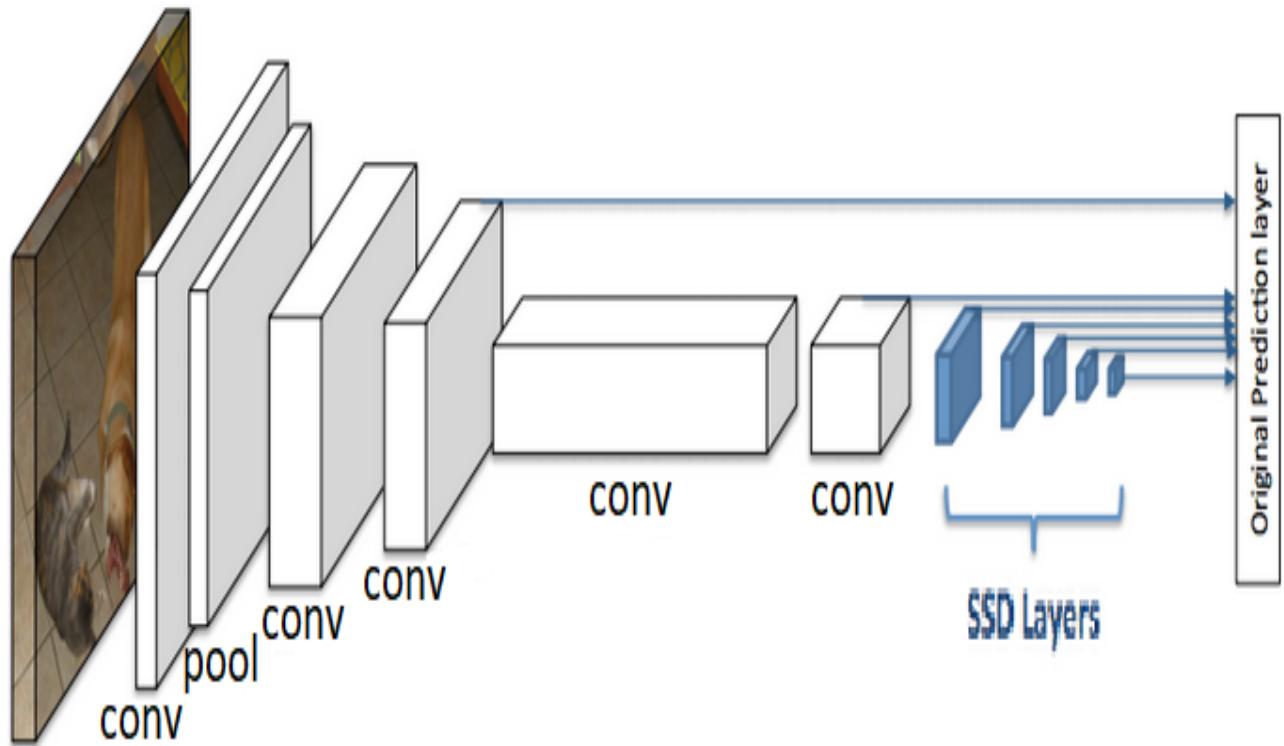


Рис. 1. Сеть с SSD-детектором

Белые прямоугольники на рис. 1 соответствуют базовой модели, а последние слои, представленные синими прямоугольниками, отвечают голове SSD.

В отличие от метода скользящего окна, SSD-алгоритм использует разделение изображения на сетку, каждая ячейка которой производит нахождение объектов в своей области. В случае отсутствия объекта он рассматривается в качестве фонового класса, а его местоположение игнорируется. Расположение и форма объекта выводится самостоятельно каждой ячейкой.

Для случая, когда в одной ячейке могут содержаться сразу несколько объектов, вводятся такие понятия, как якорная рамка (фиксированная рамка) и receptive field.

1. Якорная рамка (anchor box). Каждой клетке сетки может соответствовать несколько якорных/приоритетных рамок. Метод SSD опирается на фазу сопоставления в процессе обучения. По сути происходит предсказание класса детектируемого объекта и его местоположения. На практике каждый якорный блок определяется соотношением сторон и уровнем масштабирования.

Описываемая архитектура позволяет учесть, что форма предметов необязательно может быть квадратной, за счет предварительно заданных соотношений сторон для anchor box. Размер якорной рамки необязательно должен совпадать с размером ячейки сетки. С этой целью вводится параметр *zooms* – значение, определяющее, насколько должны увеличиваться или уменьшаться рамки по отношению к каждой ячейке сетки.

2. Рецептивное поле. Рецептивным полем называется такая область входного пространства, на которую воздействует конкретная функция CNN. Оно является главной предпосылкой SSD-архитектур, поскольку позволяет обнаруживать объекты разных масштабов, а также определять более узкую рамку.

Говоря о методах, позволяющих обнаруживать объекты в один прогон, следует также упомянуть алгоритм YOLO (You Only Look Once).

3 НЕЙРОСЕТЕВЫЕ АРХИТЕКТУРЫ В ДЕТЕКТИРОВАНИИ

3.1 Архитектура Transformer

Большое разнообразие нейросетевых архитектур для решения задач компьютерного зрения позволяет находить наилучшие комбинации, которые улучшают точность алгоритмов и скорость их работы. Одним из таких инструментов стала архитектура Transformer-моделей.

Изначально данный алгоритм использовался для решения задач обработки естественного языка. Однако модель трансформера оказалась довольно универсальной и масштабируемой, что позволило использовать ее и для других проблем глубокого обучения.

Предпосылкой создания такой архитектуры стало несовершенство использования рекуррентных нейронных сетей (RNN), которые являлись основным методом работы с последовательностями. Среди их минусов можно выделить проблемы, описанные ниже.

1. В RNN-сетях информация о последовательности находится в закрытом состоянии. При этом обновление такого состояния происходит с каждым новым шагом. При необходимости модели вернуться к той информации, которая была получена многими шагами ранее, данные должны быть сохранены внутри самого скрытого состояния без изменений. Это приводит к необходимости выбора: либо хранить большой объем информации для скрытого состояния, либо отказаться от ее хранения совсем, следовательно, потерять эти данные.

2. Обучение RNN-сетей трудно поддается процессу распараллеливания. Это связано с тем, что для получения скрытого состояния на $i + 1$ -ом шаге возникает необходимость получения состояния на i -ом шаге. В результате для батча примеров длиной в 1000 обработки будет составлять 1000 последовательных действий. Такая процедура занимает довольно много времени, а также демонстрирует низкую эффективность при работе на GPU, которые используются при параллельных вычислениях.

Таким образом, описанные проблемы осложняют использование рекур-

рентных сетей для действительно длинных последовательностей. Требуется подход, который мог бы считывать последовательность таким образом, чтобы вернуться к любому прошедшему моменту можно было за определенное время, сохранив исходную информацию. Этим свойством обладает структура self-attention («внимание к себе»), которая лежит в основе архитектуры Transformer.

Большая часть моделей, используемых для преобразования последовательностей, имеют структуру кодер-декодер [2]. Здесь кодировщик отвечает за отображений входной последовательности представлений символов (x_1, \dots, x_n) в последовательность непрерывных представлений $z = (z_1, \dots, z_n)$. С учетом z декодер затем генерирует вывод последовательности (y_1, \dots, y_m) символов по одному элементу за раз. Transformer следует этой общей архитектуре. Общая схема сети Transformer представлена на рис. 2.

Левая часть схемы представляет собой структуру энкодера. Он является стеком из идентичных слоев, каждый из которых имеет два подслоя [3]. Первый из них – multi-head self-attention («механизм самоконтроля с несколькими головами»). Второй – это простая позиционная полносвязная нейронная сеть прямого распространения. После прохождения каждого из этих слоев выход и вход складываются (такой подход получил название residual connection), а затем следует слой layer normalization (блок «Add & Norm» на рис. 2).

Декодерная часть также представляется несколькими слоями. В дополнение к двум подуровням в каждом слое кодера, декодер вставляет третий подуровень, в котором используются выходы энкодера. Аналогично кодирующющей части, происходит использование остаточных соединений для каждого из подслоев с последующей нормализацией слоя.

Проведем более детальный разбор основных составляющих архитектуры Transformer.

Слой внимания. Как было сказано выше, в сети трансформера используется слой self-attention. Данный слой отличается от обычного внимания тем, что на выходе получаются другие представления элементов последовательности, которая была подана изначально. При этом происходит взаимодействие элементов между собой.

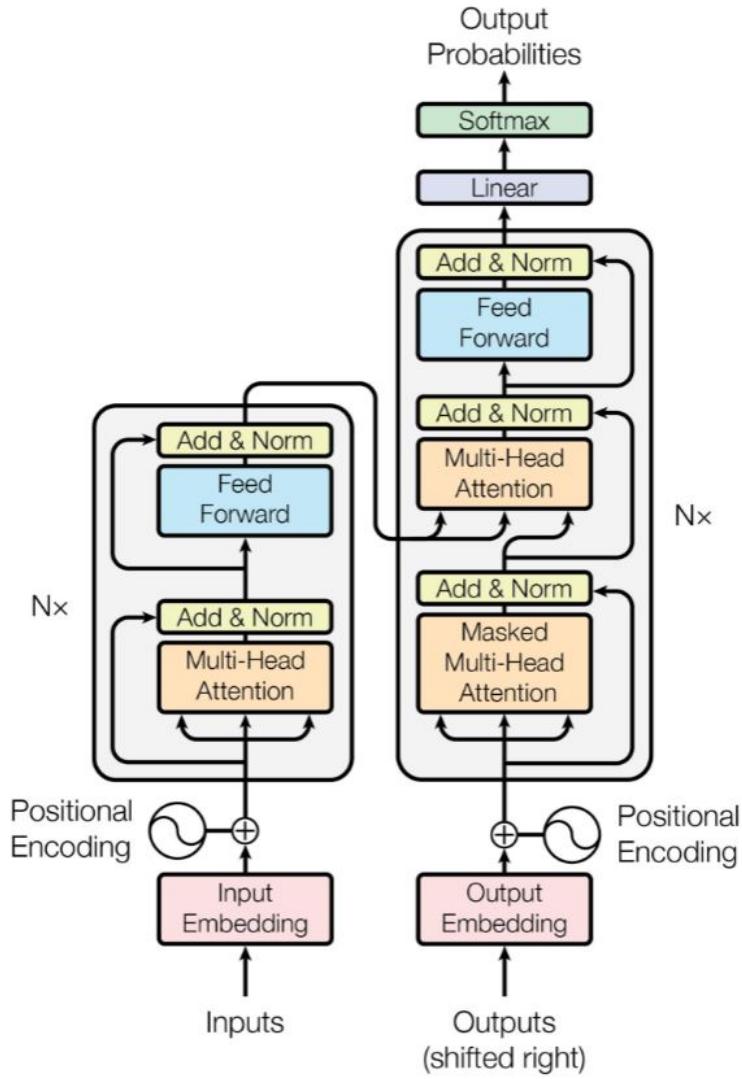


Рис. 2. Архитектура Transformer-сети

Более детально, при вычислении внимания используется обучение трех матриц W_Q, W_K, W_V .Происходит умножение каждой из этих матриц на представление x_i для каждого элемента последовательности на входе. При этом получаются вектора q_i, k_i, v_i , где i соответствует номеру элемента, причем

- 1) q_i является запросом к базе данных;
- 2) k_i – это ключи, хранящиеся в базе; по ним будет происходить поиск;
- 3) v_i отвечает самим значениям.

Понять, насколько близко значение ключа к запросу можно с использованием соотношения 3.2:

$$self - attention - weights_i = \text{softmax} \left(\frac{q_i k_1 T}{C}, \frac{q_i k_2 T}{C}, \dots \right), \quad (3.2)$$

где C – нормировочное константное значение. При этом в качестве C может

быть выбран квадратный корень $\sqrt{d_k}$ из размерности ключей и значений.

Суммирование v_i с новыми полученными коэффициентами являются выходом для *self – attention*-слоя. В матричном виде это можно представить выражением 3.3:

$$self - attention(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V. \quad (3.3)$$

Слой Multi-head attention. Конкретный набор значений Q, K, V отвечает конкретному виду соотношений для токенов, данные матрицы получают при этом ограниченных набор данных из входных. Решением данной проблемы стал multi-head-attention подход, который заключается в том, что один слой внимания заменяется сразу несколькими параллельными, в каждом из которых используются свои веса. После этого результаты объединяются. Схематично данный слой представлен на рис. 3.

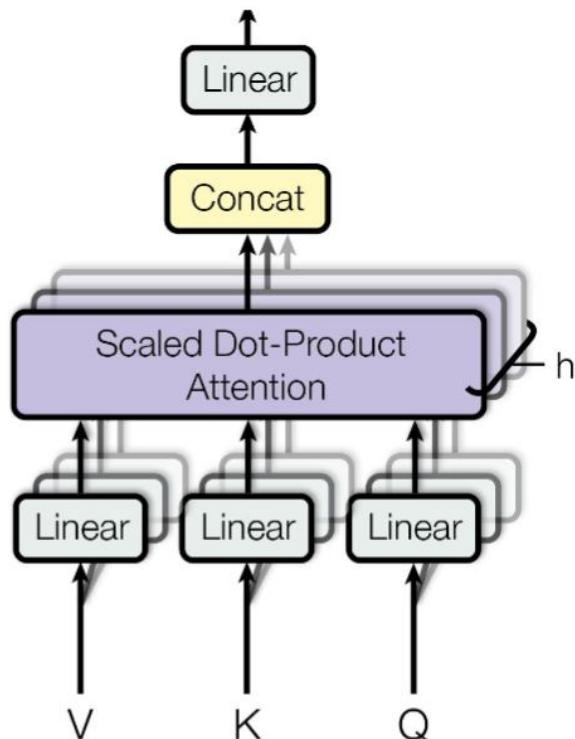


Рис. 3. Слой multi-head attention

Полносвязный слой. Следующая часть блока сети Transformer – полносвязная нейронная сеть прямого распространения (feed-forward network, FFN). Она включает в себя два полносвязных слоя, каждый из которых может независимо применяться к любому элементу для последовательности

входа. В последнее время для архитектур выходное значение после прохождения первого слоя (промежуточное значение) может превышать выходы для блока в четыре раза. В связи с этим следует учесть, что для больших моделей FFN может работать значительно дольше, чем self-attention.

Сеть Vision transformer

Одним из наиболее значимых примеров того, где нашли применение Transformer-модели, стало компьютерное зрение. Архитектура, называемая Vision transformer, в свое время продемонстрировала наилучшие результаты в плане качества для задач классификации изображений. Этому поспособствовала реализация self-attention-идеи для изображений, которые были разделены на множество квадратных сегментов. В данном подходе изображения представляются в виде визуальных токенов, после чего подаются на вход сети.

Общий алгоритм работы сети для решения задач компьютерного зрения представляется следующим образом:

- 1) получение карты признаков (feature map) для исходного изображения;
- 2) конвертация карт признаков в визуальные токены;
- 3) подача токенов на вход сети трансформера.

Выход сети может быть использован для решения задач классификации. В том случае, если объединить выход с картой признаков, то может быть реализована задача сегментации.

Архитектура Vision transformer (ViT) представляет собой комбинацию трех частей – токенизатора, трансформера и проектора.

Первый компонент, токенизатор, как следует из названия, отвечает за извлечение токенов. Второй осуществляет прохождение через сеть, которая была описана ранее, а третий производит сложение feature map и выхода трансформера. Схема ViT представлена на рис. 4.

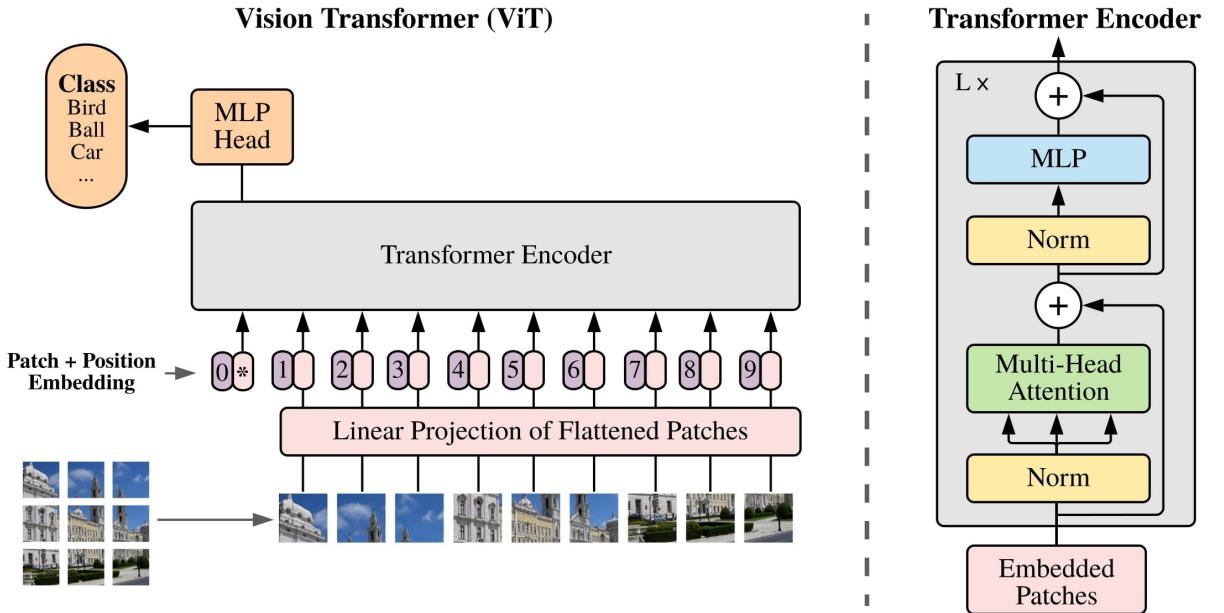


Рис. 4. Схема сети ViT

Сеть DETR

Еще одной модификацией, которая используется для детектирования, является применение сети DETR (Object Detection with Transformers). В этом случае так же, как и для сети ViT, трансформер применяется не к самому изображению, а к признакам, которые были выделены сверточной сетью. Схема работы для сети DETR представлена на рис. 5.

Однако данный подход имеет ряд существенных недостатков:

- 1) самая важная проблема – долгое и сложное обучение сети;
- 2) сеть DETR не является универсальной; существует ряд задач, в которых другие подходы проявляют себя значительно лучше

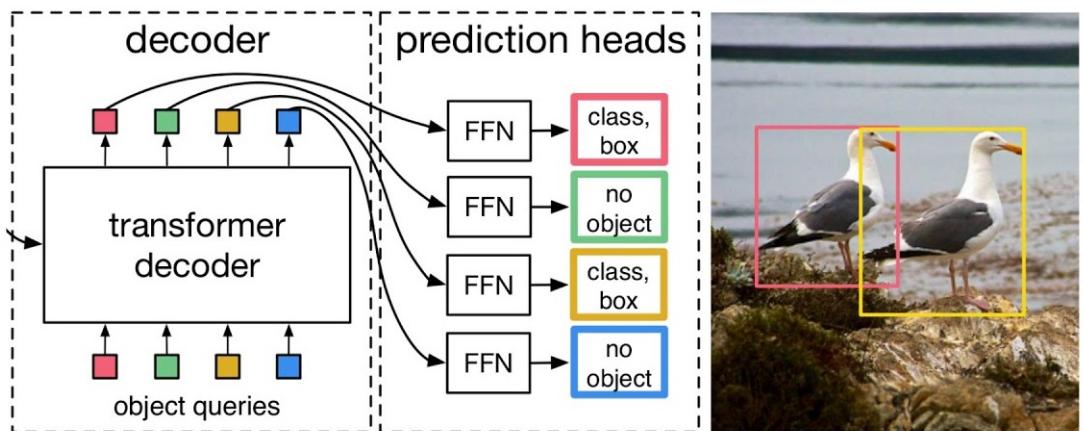


Рис. 5. Схема сети DETR

3.2 Архитектура YOLO

Другой архитектурой, применяемой при решении задачи детекции, является YOLO-сеть («You Only Look Once»). Данный алгоритм показал отличные возможности по скорости работы, что позволяет распознавать образы на изображении в режиме реального времени. Дословно аббревиатуру YOLO можно перевести как «стоит лишь раз взглянуть», такое название дано сети, поскольку изображение проходит через сверточную сеть только однажды. Другие же алгоритмы предусматривают многократное прохождение через CNN, что замедляет их работу.

Архитектура сети YOLO представлена на рис. 6.

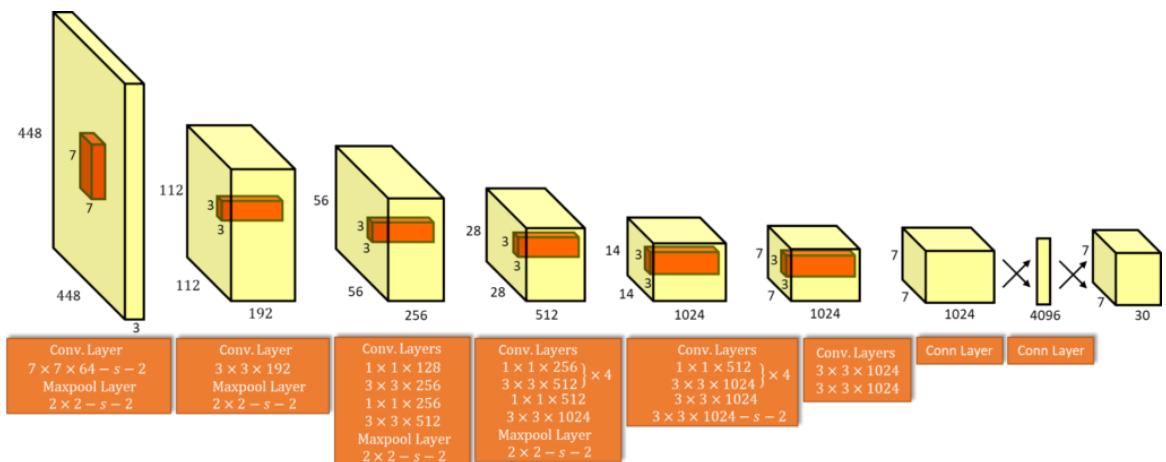


Рис. 6. Архитектура сети YOLO

Как было сказано выше, YOLO использует сверточные нейронные сети для извлечения признаков из изображений. Обычно используются предобученные сверточные сети, такие как Darknet, для извлечения признаков и последующего применения полносвязных слоев для предсказания ограничивающих рамок и классов. Данная модель также использует несколько слоев обнаружения на разных масштабах для обнаружения объектов разного размера. Более крупные объекты обнаруживаются на более высоких слоях, а меньшие объекты – на более низких слоях. Это позволяет YOLO эффективно обнаруживать объекты разного масштаба на изображении.

В методе YOLO результаты обнаружения представляются в виде тензора размером $7 \times 7 \times 1024$. Оценка вероятности нахождения объекта конкретного класса в текущем обрамляющем прямоугольнике – это произведение оценки

вероятности нахождения объекта в ячейке и оценки вероятности для конкретного класса.

Основная идея YOLO-архитектуры заключается в детектировании как задаче регрессии. Изначально изображение с использованием сети разбивается на несколько ячеек. При его прохождении через CNN может быть получено как местоположение объекта, так и его категория и соответствующая доверительная вероятность. Таким образом, имеется сетка, каждый элемент (ячейка) которой должна определить следующее:

- 1) координаты t_x, t_y, t_w, t_h ограничивающей рамки, а также значение, определяющее, присутствует ли искомый объект в ячейке (так называемый «индекс объектности»);
- 2) несколько вероятностных классов.

Определение данных значений происходит методом k -средних. Пусть координаты (c_x, c_y) соответствуют левому верхнему углу. Тогда координаты центра (b_x, b_y) ограничивающей рамки, а также ее ширина b_w и высота b_h могут быть предсказаны, исходя из выражений (3.4):

$$\begin{aligned} b_x &= \sigma(t_x) + c_x, \\ b_y &= \sigma(t_y) + c_y, \\ b_w &= p_w e^{t_w}, \\ b_h &= p_h e^{t_h}, \end{aligned} \tag{3.4}$$

где p_w и p_h – ширина и высота привязки, $\sigma(\cdot)$ – сигмоидная функция активации. Графически смысл используемых обозначений представлен ниже на рис. 7.

При использовании данной архитектуры может возникнуть следующая проблема: получение нескольких предсказанных ограничивающих рамок для одного класса. Можно было бы оставлять только одну рамку, соответствующую наибольшей вероятности нахождения в ней объекта, однако такой вариант неприменим в том случае, если на изображении имеется несколько экземпляров одного класса. С целью решения данной проблемы используется метод «немаксимального подавления». Изначально выбирается рамка с наибольшей вероятностью и сравнивается со всеми другими для конкретно-

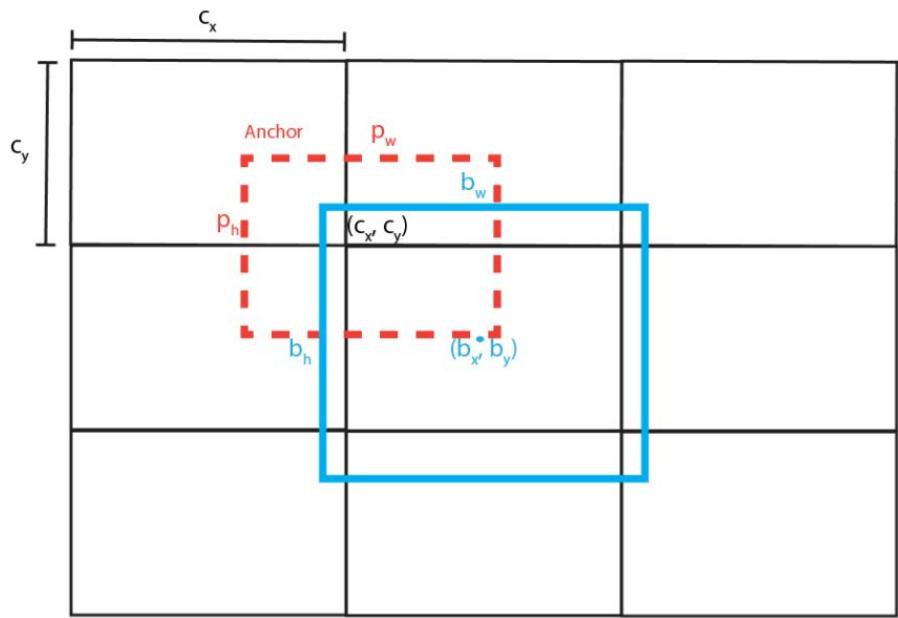


Рис. 7. Архитектура сети YOLO

го класса. Данная процедура использует IoU-метрику. В том случае, когда значение IoU меньше заданного порога (к примеру, 0.5), ограничивающий прямоугольник, соответствующий меньшей вероятности, исключается из рассмотрения (ситуация, когда два прямоугольника указывают на один и тот же объект).

4 МОДЕЛИРОВАНИЕ НЕЙРОННОЙ СЕТИ

4.1 Описание используемого датасета

Реализация задачи детекции в данной работе производилась на готовом датасете Car Object Detection, выгруженного с сайта, предлагающего готовые наборы данных для машинного обучения <https://www.kaggle.com> [4]. Указанный набор содержит медиафайлы, снятые в условиях дорожной обстановки. На кадрах фиксировались автомобили. Исходные данные представляют собой два видеоряда, каждый из которых был разрезан на стоп-кадры (которые и являются исходными изображениями). Общее количество изображений в датасете равно 1001.

Суть задачи сводится к тому, чтобы протестировать различные нейросетевые архитектуры, для обнаружения на изображении автомобилей. В датасете также имеются тренировочные изображения – каждый автомобиль на них отмечен своей ограничивающей рамкой.

Набор представляет собой изображения трех типов:

- 1) кадры, на которых имеется только один автомобиль;
- 2) кадры, на которых имеется несколько автомобилей;
- 3) кадры, на которых нет автомобилей.

Максимальное количество машин на изображении равно семи. Таким образом, задача усложняется тем, что исходное количество детектируемых на изображении объектов заранее неизвестно.

4.2 Предобработка датасета для сети ViT

У трансформеров есть несколько ограничений при использовании. Одно из них связано с ограниченным фиксированным выходом из нейронной сети. Для того, чтобы избежать этой проблемы количество выходных объектов K устанавливается значительно больше, чем максимальное ожидаемое количество объектов на изображениях.

Для возможности сопоставления истинных меток с предсказанными нейронной сетью, датасет был преобразован способом, описанным ниже.

1. Координаты каждой ограничивающей рамки были преобразованы из абсолютных размеров, измеряемых в пикселях, в относительные по отношению к размерам изображения;
2. Полученные координаты переведены из описания левого верхнего угла (x_{min}, y_{min}) и правого нижнего (x_{max}, y_{max}) ограничивающей рамки в описание центра (x_c, y_c) и размеров изображения (h, w) ;
3. Для каждого изображения все соответствующие ему ограничивающие рамки были объединены в массив, который дополнили нулевыми ограничивающими рамками до размера массива равного K ;
4. Каждой ограничивающей рамке была добавлена еще одна координата, соответствующая нахождению в ней объекта поиска, т.е. для ненулевых рамок 1, для остальных – 0.

Таким образом размерность метки для одного изображения равняется $(K, 5)$, при этом некоторые из рамок будут заполнены нулями.

4.3 Венгерский алгоритм

Постановка задачи о назначениях

Задача о назначениях – это задача оптимизации, заключающаяся в поиске оптимального соответствия между двумя наборами элементов. В общем виде задача о назначениях формулируется следующим образом.

Имеется n работ и n кандидатов для их выполнения. Затраты i кандидата на выполнение j работы равны c_{ij} . Каждый кандидат может быть назначен только на одну работу, и каждая работа может быть выполнена только одним кандидатом. Требуется найти назначение кандидатов на работы, при котором суммарные затраты на выполнение работ минимальны.

Для построения математической модели вводятся переменные x_{ij} , принимающие значение 1 или 0 в зависимости от того, назначен или нет i механизму на j работу соответственно.

Тогда условие о том, что каждый кандидат выполняет только одну работу, запишется в виде

$$\sum_{j=1}^n x_{ij} = 1, \quad \forall i = \overline{1, n}. \quad (4.5)$$

Условие о том, что каждая работа может выполняться одним кандидатом, запишется в виде

$$\sum_{i=1}^n x_{ij} = 1, \quad \forall j = \overline{1, n}. \quad (4.6)$$

Целевая функция примет вид

$$F = \sum_{i=1}^n \sum_{j=1}^n c_{ij} x_{ij}. \quad (4.7)$$

В функцию войдут только те значения c_{ij} , для которых x_{ij} отличны от 0, т.е. затраты, соответствующие назначенным работам.

В результате математическую модель задачи о назначениях можно описать следующей системой уравнений

$$\begin{cases} F = \sum_{i=1}^n \sum_{j=1}^n c_{ij}x_{ij} \rightarrow \min, \\ \sum_{j=1}^n x_{ij} = 1, \quad \forall i = \overline{1, n}, \\ \sum_{i=1}^n x_{ij} = 1, \quad \forall j = \overline{1, n}, \\ x_{ij} \in \{0, 1\}, \quad \forall i = \overline{1, n}, \forall j = \overline{1, n}. \end{cases} \quad (4.8)$$

Алгоритм решения

При решении задачи используется тот факт, что оптимальность не нарушается при уменьшении или увеличении строк и столбцов матрицы $C = \{c_{ij}\}$ на одно и тоже число. Основная идея венгресского алгоритма решения задачи о назначениях состоит в переходе от матрицы C к матрице C' такой, которая состоит из неотрицательных элементов и n независимых нулей (никакие два не принадлежат одной строке или одному столбцу матрицы).

Алгоритм состоит из следующих этапов.

1. В каждой строке таблицы C найти наименьший элемент и вычесть его из всех элементов данной строки.

$$c_{ij} = c_{ij} - \min_j(c_{ij}), \quad \forall i = \overline{1, n} \quad (4.9)$$

2. В каждом столбце полученной таблицы найти наименьший элемент и вычесть его из всех элементов данного столбца.

$$c_{ij} = c_{ij} - \min_i(c_{ij}), \quad \forall j = \overline{1, n} \quad (4.10)$$

3. Найти строку, содержащую наименьшее количество нулей. Отметить один из нулей в этой строке. Остальные нули в этой строке и столбце зачеркиваем. Затем снова найти строку с наименьшим количеством незачеркнутых нулей и повторить процедуру среди незачеркнутых нулей.
4. Повторить действия из 3 пункта n раз. Если в результате количество отмеченных нулей совпадает с n , то получено оптимальное решение. В противном случае перейти к шагу 5
5. Отметить все строчки, не содержащие отмеченные нули. Отметить все столбцы с нулями в этих строчках. Отметить все строчки, содержащие

отмеченные нули. Повторять последние 2 шага, пока новые строки и столбцы не перестанут отмечаться. Провести линии через все отмеченные столбцы и неотмеченные строчки. Таким образом, на данном шаге покроются все нули минимальным количеством линий.

6. Из непокрытых линиями элементов матрицы найти наименьший, вычесть его из непокрытых элементов и прибавить к элементам на пересечении линий. После проделанного количества нулей, отмечаемых на шаге 3 должно вырасти минимум на 1. Вернуться к шагу 3.

В результате данного алгоритма матрица C будет преобразована в C' , а решением $X = \{x_{ij}\}$ будет матрица, элементы которой равны единице, если $c'_{ij} = \text{ноль}$, отмеченный в пункте 3.

4.4 Функция потерь сети ViT и сопоставление истинных и предсказанных меток

В качестве функции потерь была выбрана следующая функция:

$$L(y_{true}, y_{pred}) = \lambda_1 BinCrossEntropy(p_{true}, p_{pred}) + \\ + \{p_{true} = 1\} L_{bbox}(bbox_{true}, bbox_{pred}), \quad (4.11)$$

где $y_{true} = (bbox_{true}, p_{true})$ – истинная метка, $y_{pred} = (bbox_{pred}, p_{pred})$ – метка, предсказанная нейросетью, $bbox_{true}$, $bbox_{pred}$ – ограничивающие рамки объекта, p_{true} , p_{pred} – вероятность нахождения объекта в ограничивающей рамке, λ_1 – гиперпараметр.

Функции $BinCrossEntropy(p_{true}, p_{pred})$ и $L_{bbox}(bbox_{true}, bbox_{pred})$ определяются следующим образом:

$$BinCrossEntropy(p_{true}, p_{pred}) = -p_{true} \log(p_{pred}) - \\ - (1 - p_{true}) \log(1 - p_{pred}), \quad (4.12)$$

$$L_{bbox}(bbox_{true}, bbox_{pred}) = -\lambda_2 \log(IoU(bbox_{true}, bbox_{pred})) + \\ + \lambda_3 \|bbox_{true} - bbox_{pred}\|_1, \quad (4.13)$$

где λ_2 , λ_3 – гиперпараметры, $\|\cdot\|_1$ – L1 норма.

IoU (Intersection over Union) – площадь пересечения ограничивающих рамок, деленная на площадь их объединения.

Для того чтобы определить соответствие между объектами на входном и выходном изображениях, используется венгерский алгоритм. Он находит наилучшее соответствие между объектами таким образом, чтобы минимизировать общую ошибку. Его использование реализовано следующим образом:

1. Для каждого y_{true} и y_{pred} вычисляется значение функции потерь L ;
2. На полученной матрице C размера $K \times K$, используется венгерский алгоритм, описанный в главе 4.3.

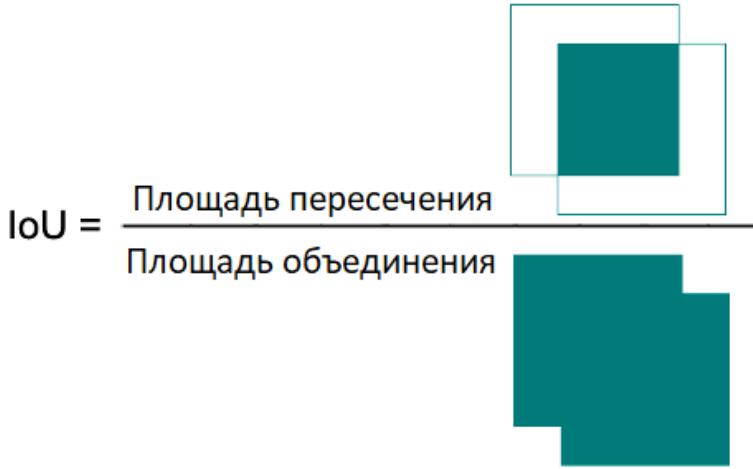


Рис. 8. IoU

4.5 Метрики сети ViT

Поскольку сети ViT используются нестандартные выходы из сети, то для оценки качества сети необходимо было разработать особые метрики. Отдельно стоит сказать, что все метрики расчитываются для перестановки выходов, оптимальной для функции потерь данной сети.

Использовались следующие метрики:

- Точность нахождения объектов – метрика, описывающая долю правильно предсказанных p_{pred} . Формула, по которой вычисляется метрика, представлена ниже.

$$P = \frac{\sum_{i=1}^K \delta_i}{K}, \quad (4.14)$$

где δ_i – символ Кронекера

$$\delta_i = \begin{cases} 1, & \text{если } (p_{pred})_i = (p_{true})_i, \\ 0, & \text{иначе.} \end{cases} \quad (4.15)$$

Таким образом, данная метрика является поэлементной точностью, но только для вероятности нахождения объекта

- Среднее IoU – среднее арифметическое значение метрики IoU среди выходов, где $p_{true} = 1$. Если данное условие ни разу не выполняется, то считаем, что значение метрики равно единице. Алгоритм расчета метрики IoU для одной ограничивающей рамки:
 1. Перевести координаты ограничивающей рамки из описания центра объекта, ширины и высоты к координатам левого верхнего и правого

нижнего углов:

$$bbox_{old} = (x_c, y_c, h, w), \quad (4.16)$$

$$bbox_{new} = (x_{min}, y_{min}, x_{max}, y_{max}). \quad (4.17)$$

2. Найти координаты левого верхнего и правого нижнего угла пересечения предсказанной ограничивающей рамки с истинной. Для этого для каждой из координат находим максимум из минимальных координат и минимум из максимальных:

$$x_{min}^{inter} = \max(x_{min}^{true}, x_{min}^{pred}), \quad (4.18)$$

$$y_{min}^{inter} = \max(y_{min}^{true}, y_{min}^{pred}), \quad (4.19)$$

$$x_{max}^{inter} = \min(x_{max}^{true}, x_{max}^{pred}), \quad (4.20)$$

$$y_{max}^{inter} = \min(y_{max}^{true}, y_{max}^{pred}). \quad (4.21)$$

3. Вычислить ширину и высоту пересечения:

$$h^{inter} = x_{max}^{inter} - x_{min}^{inter}, \quad (4.22)$$

$$w^{inter} = y_{max}^{inter} - y_{min}^{inter}. \quad (4.23)$$

4. Если ширина или высота получились меньше либо равны нуля, то площадь пересечения равняется нулю, поскольку это значит, что ограничивающие рамки не пересекаются. Иначе просто вычисляем площадь пересечения.

$$S^{inter} = \begin{cases} h^{inter}w^{inter}, & \text{если } h^{inter} > 0, w^{inter} > 0; \\ 0, & \text{иначе.} \end{cases} \quad (4.24)$$

5. Найти площадь объединения. Для этого необходимо посчитать площади каждой ограничивающей рамки, сложить и вычесть площадь пересечения.

$$S^{union} = h^{true}w^{true} + h^{pred}w^{pred} - S^{inter}. \quad (4.25)$$

6. Вычислить значение метрики IoU, поделив площадь пересечения на площадь объединения.

$$IoU = \frac{S^{inter}}{S^{union}}. \quad (4.26)$$

После вычисления метрики IoU для каждой ограничивающей рамки рассчитывается среднее IoU по формуле

$$IoU_{avg} = \frac{\sum_{i=1}^K \delta_i IoU_i}{\sum_{i=1}^K \delta_i}, \quad (4.27)$$

где δ_i – символ Кронекера

$$\delta_i = \begin{cases} 1, & \text{если } (p_{true})_i = 1, \\ 0, & \text{иначе.} \end{cases} \quad (4.28)$$

Совмешая данные две метрики, в дальнейшем строилась оценка обучения сети ViT.

4.6 Слои и структуры сети ViT

Полносвязный слой

Полносвязные слои (Fully Connected или Dense layers) - это тип слоев в нейронных сетях, который применяет линейную функцию к каждому входному элементу (нейрону) и создает новые признаки с заданной размерностью.

Пусть вход в слой имеет размерность (N, M) , где N - это количество примеров (например, изображений), а M - это размерность каждого примера (например, количество пикселей в изображении). Полносвязный слой преобразовывает каждый пример размерности $(1, M)$ в новый пример размерности $(1, K)$, где K - это заданное число новых признаков. Этот процесс происходит путем умножения входа на матрицу весов размерности (M, K) , после чего к результату прибавляется вектор смещения размерности $(1, K)$. В общем случае, формула для вычисления выхода из полносвязного слоя будет выглядеть следующим образом:

$$y = Wx + b,$$

где x – входной тензор, W – матрица весов, b – вектор смещения, а y – выходной тензор.

После вычисления выхода из полносвязного слоя, к результату может быть применена некоторая нелинейная функция активации, такая как ReLU, Sigmoid, Tanh и другие. Нелинейная функция активации обычно используется для того, чтобы преобразовать выход полносвязного слоя в нелинейный диапазон значений и добавить нелинейность в модель, что может улучшить ее способность обобщения и улучшить качество предсказаний.

Слой Dropout

Слой Dropout - это техника регуляризации, которая используется в нейронных сетях для предотвращения переобучения. Идея Dropout заключается в том, чтобы случайно удалять (выключать) некоторые нейроны в каждом слое, во время обучения сети.

Каждый нейрон в слое Dropout имеет вероятность p (обычно 0.5) быть

выключенным в каждом проходе обучения. При этом, выходы всех нейронов будут умножены на $1/(1-p)$, чтобы уравнять суммарную величину выходных данных слоя.

Таким образом, когда происходит проход обучения сети, Dropout выключает случайно выбранные нейроны в каждом проходе. Это заставляет нейроны в слое работать более самостоятельно, а не перенаправляться на зависимости, которые могут быть связаны с конкретными нейронами. Dropout также заставляет сеть учитывать более разнообразное представление входных данных, что обычно приводит к улучшению обобщающей способности модели.

После обучения, Dropout отключается, и все нейроны начинают работать с соответствующими весами. В этом случае, устанавливается режим "inference" (применения), и каждый нейрон является активным.

Слой Dropout часто используется вместе с полносвязными слоями (Dense layers) или сверточными слоями (Convolutional layers) в глубоких нейронных сетях, особенно в случае, когда обучение сети сталкивается с проблемой переобучения. Dropout позволяет контролировать сложность модели, предотвращает переобучение и улучшает ее обобщающую способность.

Слой Patches

В реализации сети Transformer с архитектурой ViT каждое изображение делится на патчи. Для этого создается специальный слой Patches. На вход данный класс будет принимать только размеры самих патчей. При его вызове в него передается всё входное изображение. Далее это изображение будет разделено на патчи указанного размера и каждый из них будет вытянут в одномерный вектор признаков.

Слой PatchEncoder

Для каждого из патчей важно помнить правильное их расположение. Поэтому в сети используется слой PatchEncoder. На вход данный класс будет принимать количество патчей изображения и размерность обучаемых эмб-

дингов. Эти эмбединги необходимы для кодирования номера патча в раздelenном изображении. Каждый патч проходит через полносвязный слой с количеством выходов равным размерности эмбедингов. К полученному результату прибавляется значение эмбединга соответствующее номеру патча в изображении. Таким образом на выходе слоя PatchEncoder будет вектор с размерностью эмбедингов для каждого патча.

4.7 Реализация сети ViT

Реализовано на языке Python. Основные используемые пакеты Tensorflow и в частности его модуль Keras.

Потом описываем функцию построения модели: сначала входной слой такой-то, который подается на такие-то слои и далее Какие-то слои повторяются несколько раз, меняя только кол-во фильтров

4.8 Предобработка датасета для сети YOLO

5 РЕЗУЛЬТАТЫ МОДЕЛИРОВАНИЯ

5.1 Сеть ViT

В качестве архитектуры сети была выбрана архитектура ViT. Для проверки качества обучения сети использовались следующие метрики

1. Средний IoU (Intersection over Union) – Среднее арифметическое отношений пересечения площадей ограничивающих рамок и их объединения;
2. Точность – отношение количества правильно меток присутствия объекта к количеству всех предсказаний.

Обучение сети проводилось 800 эпох. Наилучшее достигнутое значение среднего IoU равняется 54.5%. Для точности – 98.4%. Результаты обучения сети представлены на рис. 9 и рис. 10.

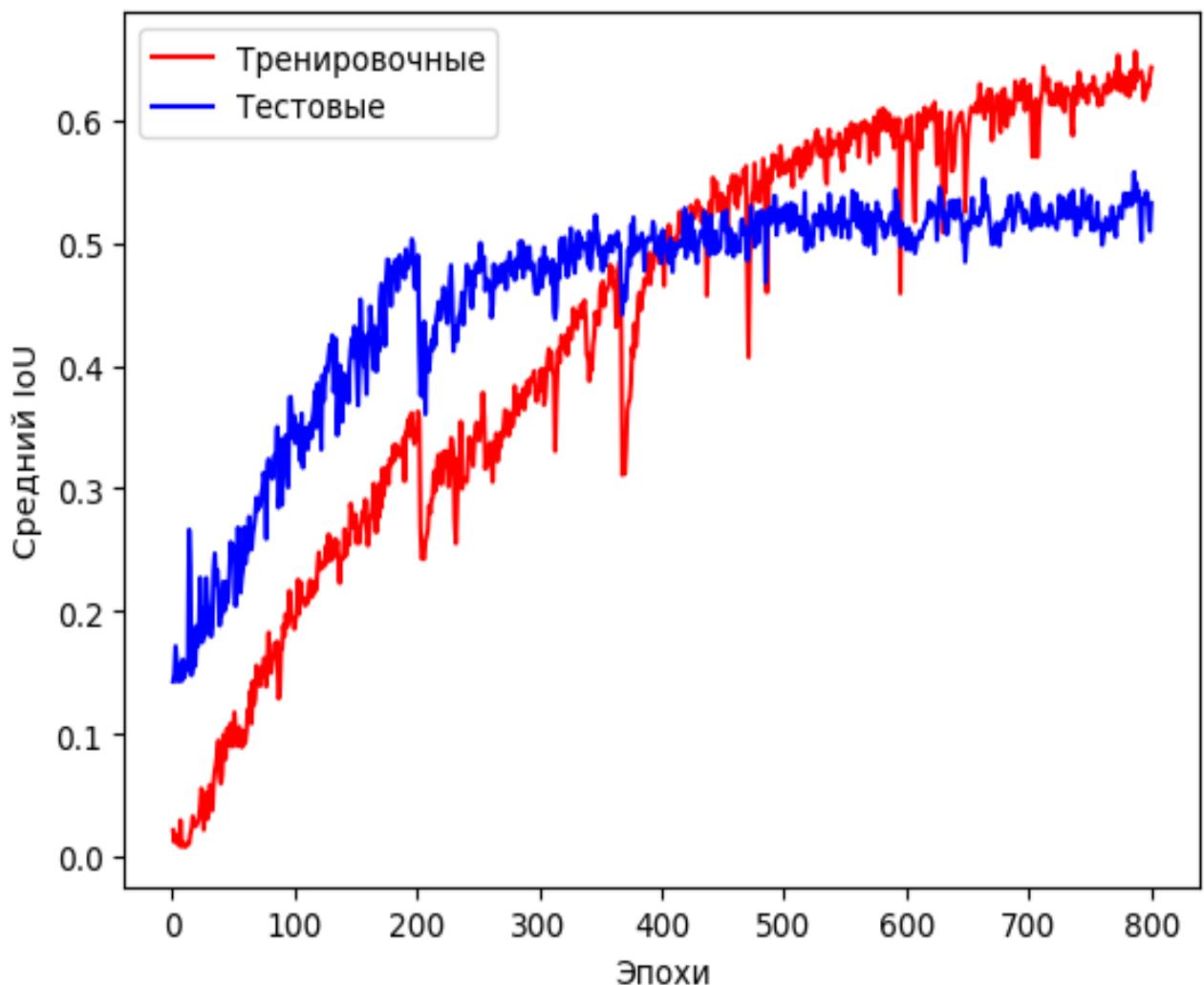


Рис. 9. График среднего IoU алгоритма ViT

Скачки, представленные на графиках, могут быть вызваны за счет использования венгерского алгоритма: поскольку у сети 10 разных выходов, каждый из которых обучается отдельно, в следствие чего в процессе обучения венгерский алгоритм возвращает разные перестановки выходов.

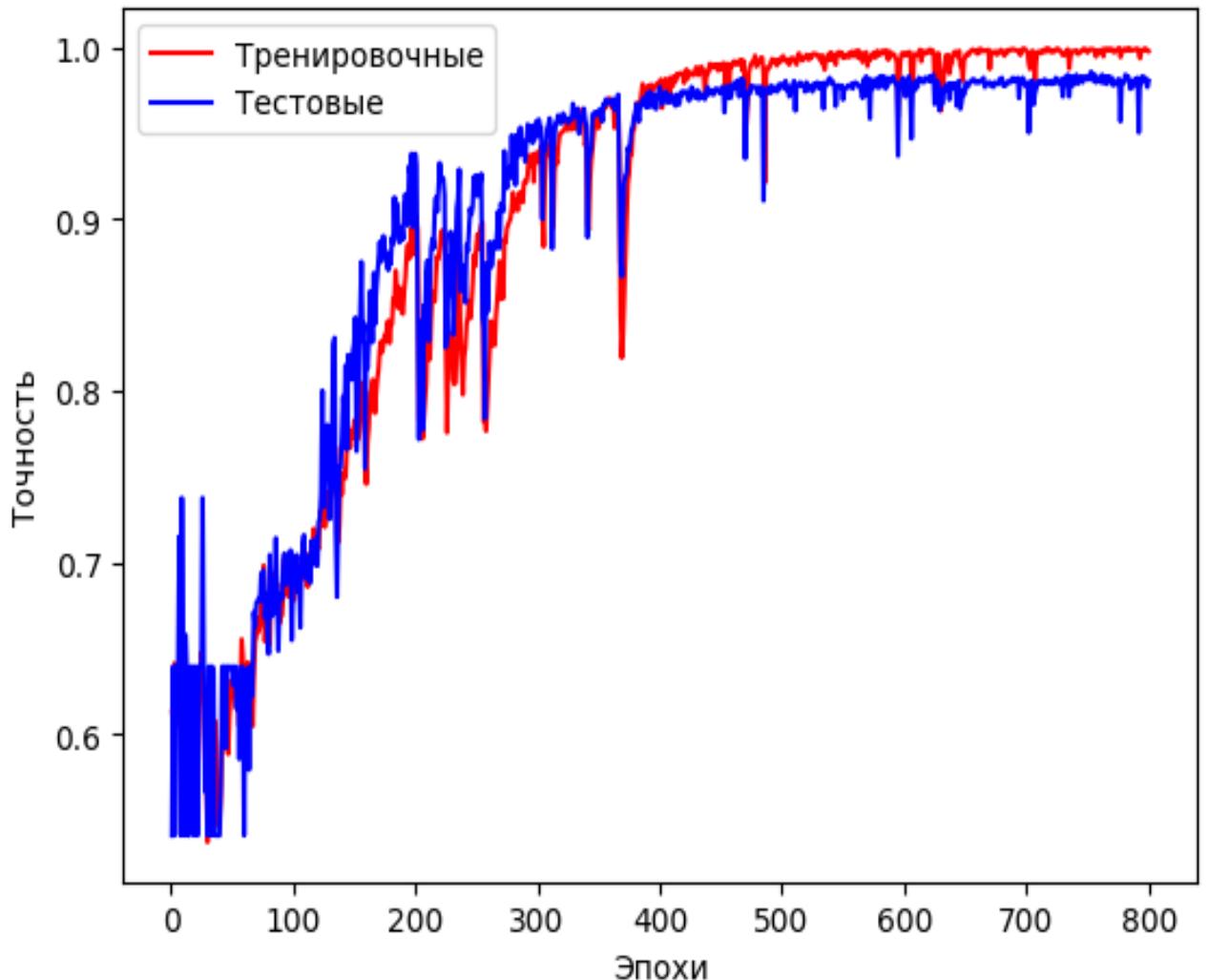
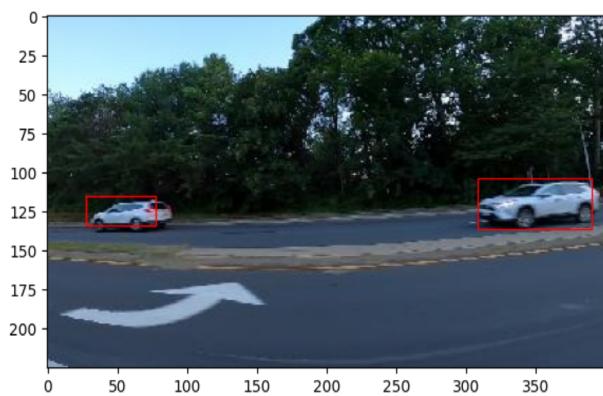


Рис. 10. График точности алгоритма ViT

Ниже на рис. 11 и рис. 12 представлены примеры обнаружения автомобилей с помощью нейронной сети и сравнение предсказанных ограничивающих рамок с истинными.

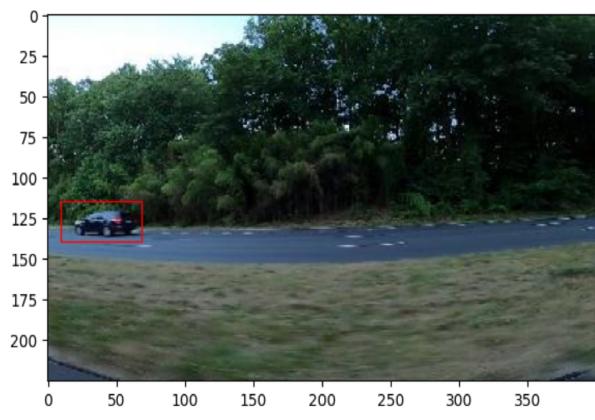


Предсказанные ограничивающие рамки

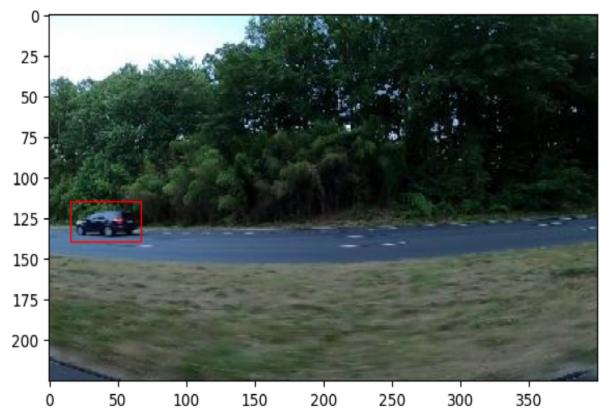


Истинные ограничивающие рамки

Рис. 11. Первый пример обнаружения объектов ViT нейросетью



Предсказанные ограничивающие рамки



Истинные ограничивающие рамки

Рис. 12. Второй пример обнаружения объектов ViT нейросетью

5.2 Сеть YOLO

ыыыыыы

ЗАКЛЮЧЕНИЕ

В настоящей работе было дано определение задаче обнаружения объектов, позволяющая выделять и классифицировать объекты на изображении. Рассмотрена нейросетевая архитектура Transformer, разобраны основные составляющие. Описан механизм внимания, способ его вычисления. Исследован слой Multi-head attention, а так же его связь с механизмом внимания. Рассмотрены 2 нейронных сети с использованием Transformer: ViT и DETR, их основные особенности и различия.

Реализована и протестирована нейросеть ViT. Были представлены графики, полученные в результате обучения.

Список использованных источников

1. Муаль М.Н.Б., Козырев Д.В., Уанкпо Г.Ж.К., Нибасумба Э. Разработка нейросетевого метода в задаче классификации и распознавания изображения // Современные информационные технологии и ИТ-образование. – 2021. – Т. 17, № 3. – С. 507-518.
2. Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, Yoshua Bengio. Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation [Электронный ресурс] // ArXiv : [сайт]. — URL: <https://arxiv.org/abs/1406.1078> (дата обращения: 28.05.2023).
3. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, Illia Polosukhin. Attention Is All You Need [Электронный ресурс] // arXiv : [сайт]. — URL: <https://arxiv.org/pdf/1706.03762.pdf> (дата обращения: 26.05.2023).
4. EDWARD ZHANG. Car Object Detection [Электронный ресурс] // Kaggle : [сайт]. — URL: <https://www.kaggle.com/datasets/sshikamaru/car-object-detection> (дата обращения: 26.05.2023).