

Министерство образования и науки Российской Федерации
Федеральное государственное бюджетное образовательное учреждение
высшего образования
Волгоградский государственный технический университет

Факультет Электроники и вычислительной техники

Кафедра Системы автоматизированного проектирования и ПК

Согласовано

(должность гл. специалиста предприятия)

(подпись) _____
(инициалы, фамилия)
«_____» _____ 2017

Утверждаю

Зав. кафедрой САПР и ПК, д.т.н.,

??.

(подпись) М. В. Щербаков
(инициалы, фамилия)
«_____» _____ 2017

ПОЯСНИТЕЛЬНАЯ ЗАПИСКА

к _____ выпускной работе бакалавра _____ на тему
(наименование вида работы)

Портирование сверточной нейросети на ARM архитектуру с
ограниченными вычислительными ресурсами и ресурсами памяти

Автор _____ Мельников Тимофей Алексеевич
(подпись и дата подписания) (фамилия, имя, отчество)

Обозначение ВСТАВИТЬ КОД-81
(код документа)

Группа ИВТ-461
(шифр группы)

Направление ??..??..?? Автоматизированные системы управления
(код по ОККО, наименование направления, программы)

Руководитель работы _____ А. В. Катаев
(подпись и дата подписания) (инициалы и фамилия)

Консультанты по разделам:

_____ (краткое наименование раздела)	_____ (подпись и дата подписания)	_____ (инициалы и фамилия)
_____ (краткое наименование раздела)	_____ (подпись и дата подписания)	_____ (инициалы и фамилия)
_____ (краткое наименование раздела)	_____ (подпись и дата подписания)	_____ (инициалы и фамилия)

Нормоконтролер _____ ????? ?????????????
(подпись и дата подписания) (инициалы и фамилия)

Волгоград, 2017

Министерство образования и науки Российской Федерации
Федеральное государственное бюджетное образовательное учреждение
высшего образования
Волгоградский государственный технический университет

Кафедра Системы автоматизированного проектирования и ПК

Утверждаю

Зав. кафедрой САПР и ПК, д.т.н.,

??.

(подпись) М. В. Щербаков
(инициалы, фамилия)
«_____» _____ 2017

Задание на _____ выпускную работу бакалавра

(наименование вида работы)

Студент _____ Мельников Тимофей Алексеевич

(фамилия, имя, отчество)

Код кафедры _____ ?? ?? Группа _____ ИВТ-461

Тема Портирование сверточной нейросети на ARM архитектуру с ограниченными вычислительными ресурсами и ресурсами памяти

Утверждена приказом по университету от «??» ?????? 201? № ????–ст

Срок представления готовой работы _____

(дата, подпись студента)

Исходные данные для выполнения работы

задание, выданное научным руководителем с кафедры САПР и ПК, утвержденное приказом ректора

Содержание основной части пояснительной записки

Что-то там раз

Что-то там два

Перечень графического материала

1) Графический материал раз

2) Графический материал два

ВСТАВИТЬ КОД-81

Руководитель работы _____

(подпись и дата подписания)

А. В. Катаев

(инициалы и фамилия)

Консультанты по разделам:

(краткое наименование раздела)

(подпись и дата подписания)

(инициалы и фамилия)

(краткое наименование раздела)

(подпись и дата подписания)

(инициалы и фамилия)

(краткое наименование раздела)

(подпись и дата подписания)

(инициалы и фамилия)

Министерство образования и науки Российской Федерации
Федеральное государственное бюджетное образовательное учреждение
высшего образования

Волгоградский государственный технический университет
Кафедра «Системы автоматизированного проектирования и ПК»

Утверждаю

Зав. кафедрой САПР и ПК, д.т.н.,

??.

_____	М. В. Щербаков
(подпись)	(инициалы, фамилия)
«_____»	_____ 2017

Портирование сверточной нейросети на ARM архитектуру с
ограниченными вычислительными ресурсами и ресурсами памяти
ПОЯСНИТЕЛЬНАЯ ЗАПИСКА

ВСТАВИТЬ КОД-81

Листов 34

Научный руководитель
старший преподаватель САПР и
ПК

_____ А. В. Катаев
«_____» _____ 2017

Нормоконтролер

?????, ????

_____ ????????????

«_____» _____ 2017

Исполнитель

студент группы ИВТ-461

_____ Т. А. Мельников

«_____» _____ 2017

Волгоград, 2017

Аннотация

Документ представляет собой пояснительную записку к выпускной работе бакалавра на тему «Портирование сверточной нейросети на ARM архитектуру с ограниченными вычислительными ресурсами и ресурсами памяти», выполненную студентом группы ИВТ-461, Мельниковым Тимофеем Алексеевичем.

В данной работе рассмотрена возможность реализации алгоритмов машинного обучения, в частности прямой проход сверточной нейронной сети, на устройстве с ограниченными вычислительными ресурсами и ресурсами памяти.

Объём пояснительной записки составил 34 страниц и включает 11 рисунков и 1 таблицы.

Содержание

Введение	6
1 Обзор фреймворков машинного обучения	8
1.1 Caffe	8
1.1.1 Основные характеристики Caffe	8
1.1.2 Архитектура Caffe	9
1.1.3 Приемущества Caffe	10
1.2 Torch7	11
1.2.1 Основные характеристики Torch7	11
1.2.2 Структуры используемых типов данных	13
1.2.3 Пакеты Torch7	14
1.3 Darknet	16
1.3.1 Основные характеристики Darknet	16
1.3.2 Используемые структуры данных	17
1.4 Сравнение фреймворков машинного обучения	19
2 Используемые алгоритмы и модели	21
2.1 Теоретические основы нейронных сетей	21
2.1.1 Нейронные сети: основные положения	21
2.1.2 Алгоритм обратного распространения ошибки	30
2.1.3 Сверточные нейронные сети	30
2.1.4 Обнаружение объектов с применением подхода YOLO	30
2.2 Оптимизация работы с памятью	30
2.2.1 Бинаризация весов	30
2.2.2 Оптимизация работы со слоями	30
3 Проектирование системы	31
Заключение	32
Список использованных источников	33
Приложение А — Техническое задание	34

Введение

Задачи обработки и анализа аналоговой информации являются одними из самых сложных в IT-индустрии. Долгое время такие задачи решались эвристическими алгоритмами, которые требовали огромных аппаратных ресурсов при малой точности результата. На протяжении последних десяти лет стремительно растет и развивается прикладная область математики цель которой, изучение и развитие искусственных нейронных сетей. Актуальность разработок и исследований в данной области оправдывается применением НС в различных сферах деятельности. Это автоматизация процессов анализа объектов, образов, уневерсализация управления, прогнозирование, создание экспертных систем, анализ неформализованной информации и многое другое. В частности, в данной дипломной работе используются нейронные сети для классификации и обнаружения объектов на изображении.

Наиболее существенным недостатком НС является их требовательность к вычислительным ресурсам и ресурсам памяти. Частично данная проблема решается использованием сверточных нейронных сетей, которые, в виду особенностям логики работы, позволяют в разы сократить ресурсы потребляемые нейронной сетью.

Однако, не только искусственные нейронные сети являются трендом IT-индустрии, активно развивается концепция интернета вещей. Диапазон встраиваемых технологий простирается от концепции умных зданий до промышленной консолидации. Совмещение встраиваемых систем и искусственных нейронных сетей позволяет иначе взглянуть на решение нетривиальных задач, таких как автономное управление автомобилем.

В связи с вышесказанным целью данной дипломной работы является внедрение фреймворка машинного обучения на embedded систему C.H.I.P. и последующая оптимизация его работы. На основе проделанной работы необходимо сделать вывод о эффективности и рентабельности данного решения.

ВСТАВИТЬ КОД-81

Для достижения поставленной цели необходимо решить следующие задачи:

- Изучить фреймворки глубокого машинного обучения
- Разработать консольное приложение для реализации прямого прохода нейронной сети
- Оптимизировать использование оперативной памяти и реализовать загрузку весов по мере использования
- Разработать клиент-серверное приложение, демонстрирующее результат работы

В первом разделе пояснительной записки описаны фреймворки машинного обучения. Далее приведено обоснование выбора фреймворка darknet.

Во втором разделе описаны используемые модели нейронных сетей и алгоритм прямого прохода.

Третьей раздел посвящен разворачиванию фреймворка на устройстве С.Н.І.Р. и оптимизации работы алгоритма прямого прохода. Так же описана разработка клиент-серверной части для визуализации работы приложения.

1 Обзор фреймворков машинного обучения

Данный раздел содержит справочную информацию, технические особенности и функциональные возможности фреймворков глубокого машинного обучения и их сравнение. Раздел содержит обоснование выбора фреймворка Darknet для встраивания и оптимизации на мобильном ПК C.H.I.P.

Из всего множества фреймворков были выделены Caffe, Torch7, Darknet, как наиболее зрелые, функционально полные и широко используемые.

1.1 Caffe

Caffe представляет собой фреймворк, разработанный учеными и практиками, с прозрачной и гибкой архитектурой для глубокого обучения и построения эталонных моделей. Фреймворк распространяется под BSD-лицензией и является C++ библиотекой. Так же реализованы python и MATLAB обертки для универсализации обучения и развертывания глубоких моделей. Caffe используется на промышленных компаниях и в медиацинтрах, обрабатывая 40 миллионов изображений в день на Titan GPU (примерно 2.5 миллисекунд на изображение).

Caffe поддерживается и разрабатывается университетом Беркли, а именно центром BVLC.

1.1.1 Основные характеристики Caffe

Caffe представляет полный набор инструментов для обучения, тестирования, настройки и разработки моделей с подробной документацией и примерами. Поэтому обучиться использовать

фреймворк можно довольно быстро. Возможность использования GPU делает Caffe одним из самых быстрых фреймворков, что позволяет его использовать в промышленном секторе. Такие показатели достигнуты благодаря особенностям описанным ниже.

Caffe является модульным программным обеспечением. Что позволяет легко добавлять новые форматы данных, слои и функции потерь. В фреймворке уже реализовано множество слоев и функций потерь, что позволяет реализовать нейронную сеть для задач различных предметных областей и категорий.

В Caffe представление и реализация разделены. Для описания модели в Caffe используется конфигурационный файл в формате protobuf. Caffe поддерживает сетевые архитектуры в форме произвольно ориентированных ациклических графов. Важной деталью является то, что после создания экземпляра модели Caffe выделяется ровно столько памяти, сколько необходимо для работы сериализованной нейронной сети и для хранения адреса объекта [1].

В Caffe используется полное тестовое покрытие. Каждый модуль имеет собственный набор тестов. Модуль будет принят, только после прохождения всего набора тестов. Это позволяет эффективно оптимизировать модули и гарантирует стабильную работу фреймворка.

Caffe содержит предворительно обученные модели для академических целей и некоммерческого использования. Доступны сверточные НС с архитектурой "AlexNet" и вариации данной НС, обученные на базе данных ImageNet[2]. Так же доступны рекуррентные модели[3].

1.1.2 Архитектура Caffe

Caffe сохраняет и передает данные в четырехмерных массивах, которые названы блобами. Блобы представляют унифицированный интерфейс для работы памятью, содержащий пакеты изображений (или других данных), параметров или обновлений параметров. Блобы

скрывают вычислительные издержки смешанной работы CPU и GPU, выполняя синхронизацию по мере необходимости. Память выделяется по требованию (лениво), что позволяет эффективней ее использовать. Модели сохраняются как буфер, использующий протокол Google (Google Protocol Buffers), который имеет ряд достоинств: минимальный размер строки при сериализации, эффективная сериализация, высокая читабельность в текстовом виде и удобные интерфейсы работы на нескольких языках. Необходимые для обучения огромные массивы данных хранятся в базах данных LevelDB. Google Protocol Buffers и LevelDB обеспечивают пропускную способность в 150 Мб/с.

Слой в Caffe представляет собой структуру соответствующую формальному определению слоя: он принимает на вход один или несколько блобов и выдает один или несколько блобов результатом. Caffe предоставляет полный набор типов слоев для глубокого обучения, включая сверточный, pooling слой, inner products слой, нелиности, такие как выпрямленная линейная и логическая, слои потерь, таких как softmax и hinge. Настройка слоя требует минимальных усилий в виду композиционного построения сетей.

Caffe обеспечивает функциональность для любого направленного ациклического графа слоев, позволяя корректно выполнять прямой и обратный проход. Модели Caffe — это сквозные системы машинного обучения.[1]

1.1.3 Преимущества Caffe

От других современных фреймворков глубокого обучения Caffe отличается следующими качествами:

- Реализация полностью основана на C++, что облегчает интеграцию с встраиваемыми системами. CPU режим позволяет использовать фреймворк без специализированного GPU.

– Готовые модели позволяют не тратить время и ресурсы на обучение. Важным пунктом является подробная документация для сериализации и использования моделей.

1.2 Torch7

Torch7 — это универсальный математический фреймворк и библиотека машинного обучения, которая имеет оболочку для языка программирования Lua. Его цель — предоставить гибкую среду для проектирования и обучения моделей глубокого обучения. Гибкость достигается с помощью Lua, так как он является очень легким скриптовым языком. Эффективная реализация низкоуровневых числовых процедур, используя OpenMP и CUDA, позволяет фреймворку достиг высокой производительности. Фреймворк имеет простой Lua-интерфейс, что позволяет легко подключать его к стороннему программному обеспечению.

1.2.1 Основные характеристики Torch7

Структура фреймворка имеет три основных преимущества:

- она облегчает разработку численных методов;
- фреймворк легко расширяем (включая использование сторонних библиотек);
- высокая скорость работы фреймворка.

Второе преимущество достигается за счет выбранных разработчиками технологий. Скриптовый (интерпретируемый) язык с хорошим API-интерфейсом для C обеспечивает фреймворку гибкость в разработке и не накладывает ограничения на его расширяемость. Так как, язык высокого уровня делает процесс разработки программы более простым и понятным, чем язык низкого уровня. К тому же,

интерпретируемость позволяет быстро и легко реализовывать различные идеи в интерактивном режиме. Хороший API-интерфейс сохраняет функциональные возможности из разных библиотек, так как становится прослойкой между универсальной структурой на языке Lua и различными структурами используемых библиотек на языке C.

Высокая скорость работы достигается благодаря компилятору JIT (Just In Time). На данный момент Lua является самым быстрым интерпретируемым языком. Lua разрабатывался для легкого внедрения в приложения, написанные на C. Поэтому представляет большое C-API на основе виртуального стека, для передачи значений между Lua и C. Это унифицирует интерфейс для C/C++ и делает обертывание библиотек тривиальным [4].

Lua предназначен для использования в качестве мощного, легкого скриптового языка обладающими всеми необходимыми выразительными средствами. Он реализован как библиотека, которая написана на чистом C (точнее на подмножестве ANSI C и C++). Lua сочетает простой процедурный синтаксис с мощными конструкциями описания данных на основе ассоциативных массивов и расширяемой семантики. Lua динамически типизируется, выполняется путем интерпретации байт-кода для виртуальной машины на основе регистров и имеет автоматическое управление памятью с инкрементной сборкой мусора, что делает его идеальным для настройки, написания сценариев и быстрого прототипирования [5].

Lua предлагает хорошую поддержку объектно-ориентированного программирования, функционального программирования и программирования, управляемого данными. Основным типом Lua является таблица, которая реализует ассоциативные массивы очень эффективным способом. Ассоциативный массив — это массив, который может индексироваться не только числами, но и любыми другими типами данных языка. Таблицы не имеют фиксированного размера, они динамически изменяемы и могут использоваться как "виртуальные таблицы" над другой таблицей, что позволяет имитировать парадигмы объектно-ориентированного программирования. Таблицы являются единственным, но очень мощным механизмом структурирования данных

в Lua. Torch7 использует таблицы для простого, равномерного и эффективного представления обычных массивов, таблиц символов, кортежей, очередей и других структур данных. Lua также использует таблицы для представления пакетов.

Lua и Python очень схожи как по структурированию данных, так и по стилю программирования. Если говорить о популярности в сообществе, то Python опережает Lua из-за огромного количества предоставляемых библиотек. Однако разработчики выбрали Lua по ряду других причин, которые, в виду специфики фреймворка, являются ключевыми. Во-первых, интеграция Lua с C очень проста. За несколько часов любая библиотека на C или C++ может стать библиотекой Lua. Во-вторых, Lua предоставляет эффективные возможности встраивания. Что бы преобразовать прототип в финальный продукт требуется не много дополнительной работы. В-третьих, Lua обладает высокой производительностью благодаря интерпритатору LuaJIT, который выдает производительность на уровне C. Еще одним преимуществом Lua является переносимость. Lua написан на чистом ANSI C, его можно скомпилировать для любых устройств (сотовые телефоны, встроенные процессоры в FPGA, процессоры DSP и др.).

1.2.2 Структуры используемых типов данных

Ключевой сущностью в Torch7 является класс Tensor, поставляемый автономной C-библиотекой Tensor. Данный класс расширяет базовый набор типов Lua, чтобы реализовать эффективную работу с многомерными массивами. Большинство пакетов Torch7 или сторонних пакетов, зависящих от Torch7, реализуют собственный класс Tensor для представления сигналов, изображений, видео и других объектов, что упрощает интегрирование различных библиотек. Библиотека Torch Tensor предоставляет множество классических операций (включая операции линейной алгебры), которые реализованы и оптимизированы на C, используются SSE инструкции для Intel

платформ. Опционально можно использовать высокопроизводительные реализации операций линейной алгебры в библиотеке BLAS. Так же данная библиотека поддерживает инструкции OpenMP и вычисления на CUDA GPU.

1.2.3 Пакеты Torch7

На данный момент Torch7 имеет 7 основных пакетов:

- torch: основной пакет Torch7. Обеспечивает фреймворк классом Tensor, облегчает сереализацию и другие базовые функции;
- lap и plot: представляют стандартные функции для создания, преобрзования и визуализации объектов Tensor. Пример работы показан на рисунке 1

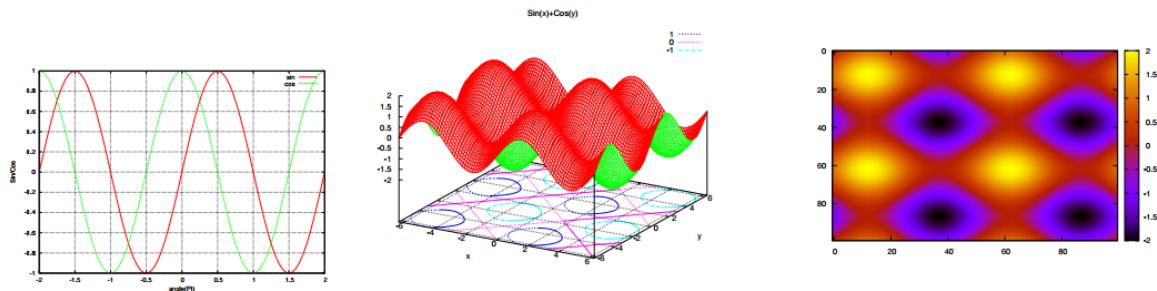


Рисунок 1 — Графики, полученные с помощью пакета plot фреймворка Torch7. Слева: простые синусоидальные функций. В центре: Поверхность, хранящаяся в 2D Tensor. Справа: Матричный график, построенный с использованием карты тепла

- qt: предоставляет интерфейс работы Torch7 с Qt. Реализует конвертацию Tensor в QImage и наоборот. Отлично подходит для быстрого создания интерактивных демонстраци с кроссплатформенным графическим интерфейсом.

- nn: предоставляет набор стандартных модулей для создания нейронной сети. В пакет так же входит набор контейнерных модулей, которые можно использовать для определения произвольно направленных графов. Явное описание графа позволяет избежать

сложности с анализатором графов или любого другого компилятора промежуточного уровня.

На практике нейронная сеть представляет собой последовательные графы, либо графы с шаблонными витвлениями и рекурсиями. На рисунке 2 показано создание многослойного перцептрона, используя пакет nn.

```
1 mlp = nn.Sequential()
2 mlp.add(nn.Linear(100,1000))
3 mlp.add(nn.Tanh())
4 mlp.add(nn.Linear(1000,10))
5 mlp.add(nn.SoftMax())
```

Рисунок 2 — Создание многослойного перцептрона, используя пакет nn

Каждый модуль или контейнер имеет стандартные функции для вычисления выходного состояния, обратного распространения производных входов и внутренних параметров. Для нейронной сети, приведенной на рисунке 2, вызов этих функций показан на рисунке 3.

```
1 Y = mlp.forward(X)           -- вычисление активации Y = f(X)
2 E = loss.forward(Y,T)        -- вычислить функцию потерь E = l(Y,T)
3 dE_dY = loss.updateGradInput(Y,T) -- вычислить градиент dE/dY = dl(Y,T)/dY
4 dE_dX = mlp.updateGradInput(X,dE_dY) -- вычислить ошибку, вплоть до dE/dx
5 mlp.accGradParameters(X,dE_dY) -- вычислить градиенты по весам: dE/dW
```

Рисунок 3 — Вычисление выходного состояния, обратного распространения производных входов и внутренних параметров

- `image`: пакет обработки изображений. Данный пакет предоставляет стандартные функции работы с изображениями (сохранение, загрузка, масштабирование, вращение, конвертация цветовых пространств, свертка и др.).

- `optim`: компактный пакет, который обеспечивает фреймворк методами оптимизации. В него входят реализация наклонного спуска, сопряженного градиента и алгоритма Бroyдена — Флетчера — Гольдфарба — Шанно (BFGS).

- `unsup`: содержит алгоритмы обучения без учителя, такие как K-means, разреженное кодирование и автокодеры.

В дополнение к основным доступен постоянно растущий список сторонних пакетов. К примеру, `mattorch`, который обеспечивает

двухсторонний интерфейс между матричным форматом Matlab и форматом Tensor или parallel, который предоставляет функции разветвления и исполнения Lua-кода на локальных или удаленных машинах, используя механизм сериализации Torch7. Этот список постоянно растет, поскольку Lua упрощает интерфейс любой библиотеки C.

1.3 Darknet

Darknet является фреймворком машинного обучения с открытым исходным кодом, написанным на C и CUDA. Он прост в установке и поддерживает вычисления как на центральном процессоре, так и на графическом.

1.3.1 Основные характеристики Darknet

Darknet один из немногих фреймворков машинного обучения, который не имеет обязательных зависимостей. Что позволяет быстро разворачивать его на встраиваемых системах. На ряду с встроенным функционалом, Darknet поставляется с двумя опциональными зависимостями:

- OpenCV — для предоставления более широкого спектра поддерживаемых форматов изображений;
- CUDA — для вычислений на GPU.

Обе не являются обязательными для установки фреймворка.

Еще одним важным преимуществом фреймворка является независимость от архитектуры системы. Darknet полностью написан на C, что делает его универсальным, а его интеграцию в встраиваемые системы или в специализированное оборудование простой и понятной.

В оригинальном виде фреймворк, поставляемый разработчиками, представляет консольное приложения для работы с нейронными сетями. С помощью него можно проектировать, обучать, тестировать нейронные сети типовых топологий. В список функций так же входит визуализация модели классификации и обучение рекуррентных моделей. Однако, конфигурация файлов исходных кодов спроектирована специально для предоставления возможности компиляции необходимых модулей в библиотеку. Поэтому фреймворк можно встраивать как нативную библиотеку в любой пользовательский проект.

Важной особенностью фреймворка является оптимизация работы с памятью и с вычислительными ресурсами. Это позволяет работать с визуальными задачами даже на устройствах с ограниченными ресурсами памяти. Darknet имеет две эффективные реализации сверточных нейронных сетей: сети с бинарными весами и XNOR-сети. В сетях с бинарными весами фильтры аппроксимируются двоичными значениями, что приводит к экономии памяти в 32 раза. В XNOR-сетях как фильтры, так и входные данные для сверточных слоев являются двоичными. XNOR-сети реализуют свертки, используя в основном бинарные операции. Это приводит к ускорению сверточных операций в 58 раз и экономии памяти в 32 раза. Данная оптимизация позволяет запускать современные нейронные сети на центральных процессорах в режиме реального времени. Если говорить о точности работы, то классификация модели AlexNet на 2.9 % меньше у сети с бинарными весами по сравнению с оригинальной реализацией. Метод используемый в сетях с бинарными весами и XNOR-сетях превосходит новейшие сетевые методы бинаризации (BinaryConnect и BinaryNets) на 16 % (тест проводился на классификацию, используя модель ImageNet)[6].

1.3.2 Используемые структуры данных

Ключевыми типами данных в Darknet являются структуры `network` и `layer`. Структура `layer` представляет собой объект для

параметров слоя сети. Данная структура имеет общий интерфейс для всех типов слоев, поэтому обладает большим набором параметров. Для расчета выходов и градиента слоев, структура предоставляет два указателя на функции `forward` и `backward` соответственно. Реализации данных функций находятся в файлах исходных кодов у каждого типа слоя. Такая модульная структура позволяет быстро добавлять новые типы слоев и компактно реализовывать операции работы с нейронной сетью. В целом, слои представляют двунаправленный связанный список, что соответствует логике работы с нейронными сетями.

Структура `network` определяет абстрактную модель для хранения внутренних параметров нейронной сети. Как и Caffe, Darknet разделяет представление и реализацию. Это реализуется разделением данных модели на конфигурационный файл, в котором определены внутренние параметры, и на бинарный файл с весами модели. Конфигурационный файл имеет строковый формат и представляет собой описание параметров нейронной сети, параметров обучения, параметров слоев и их последовательность. Формат конфигурационного файла представлен на рисунке 4



Рисунок 4 — Формат конфигурационного файла нейронной сети

Основной структурой данных в фреймворке является динамический массив. Веса, изображения, строковые таблицы хранятся в одномерном массиве, который обернут в структуру соответствующего типа данных. Данный подход позволяет сократить издержки работы с памятью.

1.4 Сравнение фреймворков машинного обучения

Для использования сверточной нейронной сети на системе с ограниченными вычислительными ресурсами и ресурсами памяти необходимо, что бы фреймворк, поставляющий данные функции удовлетворял следующим условиям:

- высокопроизводительные вычисления;
- оптимизированная работа с памятью;
- минимальное число зависимостей.

Рассмотренные выше фреймворки, используя различные технологии и алгоритмы, обеспечивают высокую производительность своих реализаций. Caffe использует библиотеку BLAS (ATLAS, Intel MKL, OpenBLAS) для векторных и матричных вычислений, Lua в совокупности с технологиями SSE, OpenMP позволяют Torch показывать высокую скорость работы, бинаризация ядер в Darknet, позволяет использовать быстрые бинарные операции для расчетов.

Если говорить о оптимизации работы с памятью, то аппроксимация фильтров и входов в Darknet позволяет значительно уменьшить объем выделяемой памяти. На рисунке 5 сравнение бинарной свертки и свертки с двойной точностью.

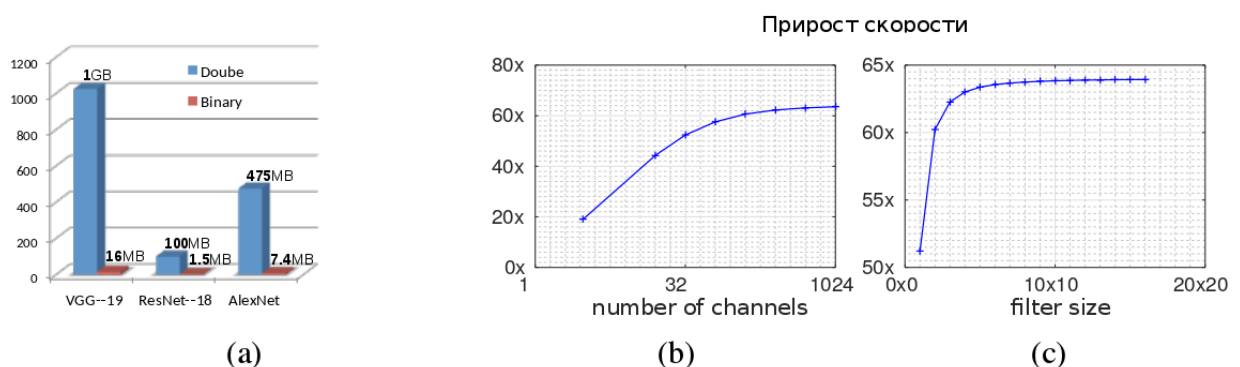


Рисунок 5 — Эффективность использования памяти и вычислений. а — выделяемая память для весов различных архитектур, б — ускорение в зависимости от числа каналов, с — ускорение в зависимости от размера фильтра

Caffe и Torch имеют достаточно большое количество зависимостей. Это объясняется желанием максимально ускорить

ВСТАВИТЬ КОД-81

процессы обучения и прохода нейронных сетей, однако накладывает ограничения на специализированное оборудование и оборудование с ограниченными запасами физической памяти.

Суммировав все показатели, можно сделать вывод, что Darknet является лучшим вариантом для разворачивания на мобильном ПК C.H.I.P.

2 Используемые алгоритмы и модели

2.1 Теоретические основы нейронных сетей

2.1.1 Нейронные сети: основные положения

Основой любой нейронной сети являются однотипные, простые элементы, которые представляют собой упрощенную модель нейронов мозга. Далее по тексту термин “нейрон” используется для определения ячейки нейронной сети — искусственного нейрона. В соответствии с клетками головного мозга, которые могут быть возбужденными или заторможенными, нейрон характеризуется состоянием в момент прохода нейронной сети. Каждый нейрон обладает набором синапсов и одним аксоном. Синапсы являются однонаправленными связями, которые связывают конкретный нейрон с выходами группы других нейронов. В свою очередь, аксон передает сигнал нейрона на синапсы нейронов, расположенных на следующем слое. На рисунке 6 представлен общий вид нейрона. Каждый синапс описывается величиной синаптической связи, иными словами, синапсы характеризуются весом w_i , который является аналогом электрической проводимости в клетках мозга.

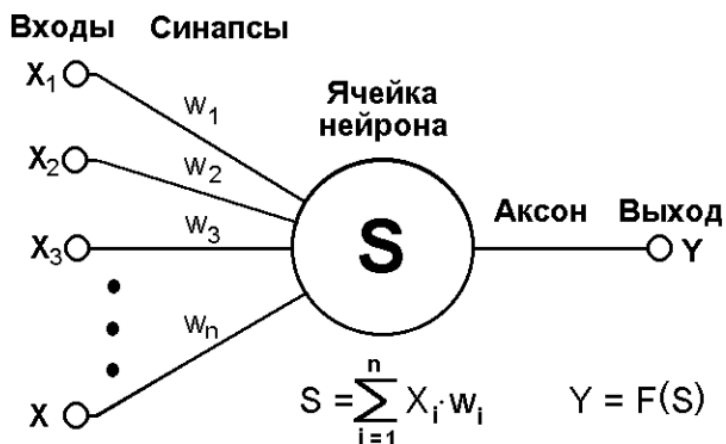


Рисунок 6 — Искусственный нейрон

Состояние нейрона в момент прохода нейронной сети определяется как взвешенная сумма его входов:

$$s = \sum_{i=1}^n x_i w_i \quad (1)$$

Выходом нейрона является функция от его состояния:

$$y = f(x) \quad (2)$$

Функция f должна обладать свойством нелинейности. Это необходимо для создания многослойных нейронных сетей. Если в НС используется пороговая функция, то смысла в ее многослойности нет, так как такая сеть эквивалентна сети с одним скрытым слоем и с весовой матрицей единственного слоя.

Нелинейная функция f именуется активационной функцией нейрона. На данный момент существует огромное количество видов активационных функций. На рисунке 7 показаны некоторые из них.

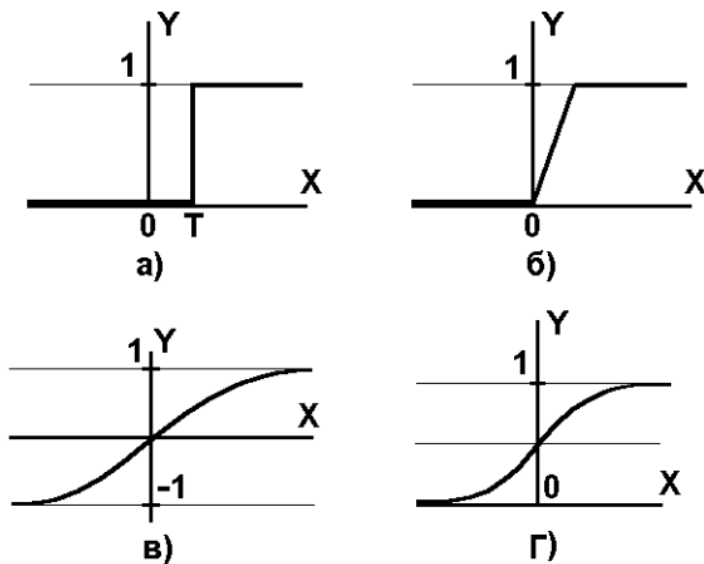


Рисунок 7 — а) функция единичного скачка; б)линейный порог (гистерезис); в) сигмоид – гиперболический тангенс; г) сигмоид – формула (3)

Одной из самых первых используемых активационных функций является логистическая функция или сигмоид (функция имеет

S-образный вид):

$$f(x) = \frac{1}{1 + e^{-\alpha x}} \quad (2)$$

Чем меньше параметр α , тем функция становится более пологой. В пределе при $\alpha = 0$ сигмоид вырождается в горизонтальную прямую в значении 0.5. Если увеличивать α , то сигмоид преобразуется в функцию единичного скачка в точке $x = 0$. Значение данной активационной функции лежит в интервале $[0, 1]$. Популярность функции обеспечивает простота ее производной, которая используется при обучении НС.

$$f'(x) = \alpha f(x)(1 - f(x)) \quad (2)$$

Логистическая функция дифференцируема на всей оси абсцисс. Это свойство используется в некоторых алгоритмах обучения. Также, сигмоид усиливает слабые сигналы лучше, чем большие, это позволяет избежать перенасыщения от больших сигналов, так как области определения больших сигналов соответствуют пологому наклону функции.

Если говорить про обработку сигналов НС, то, зачастую, они обрабатываются параллельно. Это достигается с помощью объединения большой группы нейронов в слои и соединения определенным образом нейроны разных слоев. Существуют конфигурации, где нейроны одного слоя соединены между собой. Данная конфигурация обрабатывается послойно.

На рис 8 изображена простейшая конфигурация нейронной сети — трехнейронный перцептрон. Пусть нейронной этой НС используют активационную функцию в виде скачка.

На n входов поступают некоторые сигналы, которые распространяются на три нейрона, образующие скрытый слой НС. Каждый нейрон выдает сигнал:

$$y_j = f \left[\sum_{i=1}^n x_i w_{ij} \right], j = 1 \dots 3$$

Из весовых коэффициентов можно составить матрицу W , в которой w_{ij} – вес i -того входного сигнала в j -том нейроне. Тогда, процесс прохода НС описывается в матричной форме следующим образом:

$$Y = F(XW) \quad (2)$$

где X – вектор входных сигналов, Y – вектор выходных сигналов, $F(XW)$ – активационная функция, выполняющаяся над каждым элементом вектора XW .

Теоретически количество слоев (глубина) и количество нейронов в них (высота), используемых в НС, не ограничено, но фактически ограничения накладывают вычислительные мощности устройства, на котором выполняется обработка НС. Но чем сложнее НС, тем масштабнее задачи она может решить.

Структура НС зависит от сложности задачи. Оптимальные конфигурации для некоторых типов задачи уже реализованы и описаны, например в [8],[9],[10]. Если же задача не является типовой, то разработчик самостоятельно генерирует модель, в зависимости от сложности задачи, размера обучающей выборки и вычислительных ресурсов. При этом необходимо учитывать основополагающие принципы: качество модели напрямую зависит от количества нейронов сети, плотности связей между ними и количеством слоев; сложность алгоритмов функционирования сети (например, введение нескольких типов синапсов, использование непороговых активационных функций) влияет на производительность НС. Задача поиска оптимальной конфигурации для той или иной задачи является

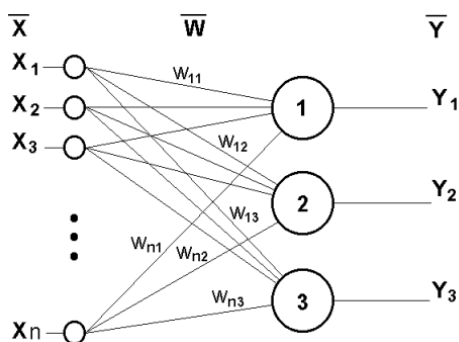


Рисунок 8 — Однослойный перцептрон

отдельным направлением нейрокомпьютерной науки. Синтез нейронной сети напрямую зависит от типа решаемой задачи, поэтому список подробных рекомендаций составить затруднительно. В большинстве случаев оптимальный вариант выбирается эмпирическим методом.

Очевидно, что функционирование нейронной сети напрямую зависит от величины синаптических связей между нейронами. Поэтому, после нахождения конфигурации нейронной сети, разработчик должен найти оптимальные значения всех переменных весовых коэффициентов (некоторые веса могут быть постоянными).

Описанный процесс называется обучением нейронной сети, он является ключевым при создании НС. От того, насколько хорошо он будет выполнен, зависит качество решений поставленных задач перед нейронными сетями. На этапе обучения кроме качества поиска весов важное место занимает такой параметр как время обучения. Эти два параметра обратно пропорциональны: чем лучше подобраны веса, тем больше затрачено времени на обучение.

Существует два варианта обучения: с учителем и без него. В первом случае, при обучении предоставляются как входные сигналы, так и желаемые выходные. Далее обучение представляет собой алгоритм подгонки весов, таким образом, чтобы желаемые выходные сигналы совпадали с выходными сигналами НС. Во втором случае, выходы генерируются нейронной сетью, а при обучении учитываются только входные и производные от них сигналы.

Существующие алгоритмы машинного обучения делятся на два типа: детерминистские и стохастические. В первом случае подбор оптимальных весов представляет собой четкую последовательность, во втором — подчинен некоторому случайному процессу.

Необходимо сказать, что среди классификаций НС важное место занимают бинарные и аналоговые сети. Первые используют двоичные сигналы, в результате чего выход каждого из нейронов принимает одно из двух значений: логический ноль ("заторможенное" состояние) или логическая единица ("возбужденное" состояние). К такой классификации относится перцептрон, описанные выше. Его активационная функция является пороговой, значение которой либо 0 либо 1. В аналоговых

сетях выходное значение нейронов может быть непрерывным, это реализуется использованием в качестве активационных непрерывные функции, например сигмоид.

Еще одна классификация разделяет НС на синхронные и асинхронные[9]. Первый случай предполагает изменение состояния только одного нейрона в каждый момент времени. Во втором случае изменение происходит одновременно у группы нейронов, как правило, у всего слоя. Ход времени в НС представлен последовательным выполнением однотипных действий над нейронами. В данной главе будут рассмотрены только синхронные НС.

Обычно, сети классифицируют по числу слоев. На 9 показана НС полученная добавлением еще одного слоя, состоящего из двух нейронов, в НС, изображенную на 8. Слои, которые не являются входными и выходными, называются скрытыми.

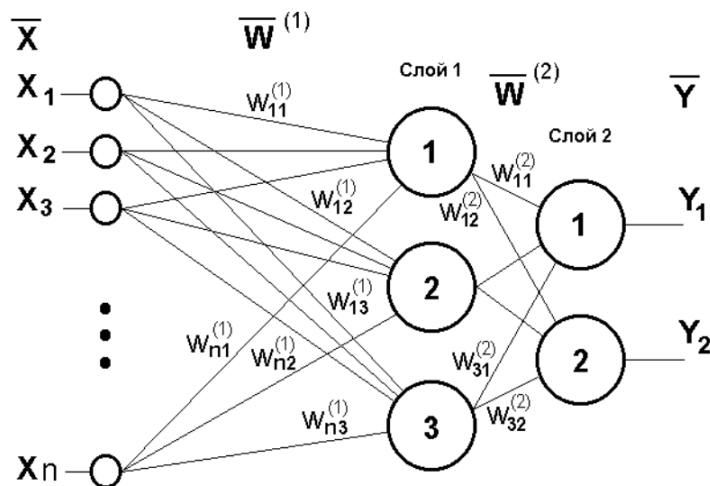


Рисунок 9 — Двухслойный перцептрон

Бывают случаи, когда нелинейность используется еще и в синаптических связях. Большинство современных сетей используют формулу (1) для вычисления значения нейрона, однако, для эффективного решения некоторых задач используется другая запись, например:

$$s = \sum_{i=1}^n x_i^2 w_i \quad (2)$$

или даже

$$s = \sum_{i=1}^n x_i^2 x_{((i+1) \bmod n)} w_i \quad (2)$$

Главное, что бы разработчик понимал, какие цели он преследует при наделении нейрона подобной связью и какие ограничения на нейрон накладываются. Введение такой нелинейности увеличивает вычислительную мощность НС, другими словами, позволяет уменьшить число нейронов и связей без потери качества работы[8].

При обучении НС учитывается не только время процесса и качество обучения. Помимо этих параметров необходимо подобрать пороговое значение T . Из рисунка 7 видно, что, в общем случае, T может принимать произвольное значение. То же самое относится и к центральной точке сигмоиды, положение которой изменяется по оси X влево или вправо. В общем случае каждая активационная функция имеет параметр, который необходимо подобрать при обучении. В связи с этим формула (1) должна выглядеть следующим образом:

$$s = \sum_{i=1}^n x_i w_i - T \quad (11)$$

Что бы добавить данное смещение, необходимо добавить еще один вход, который имеет синаптическую связь со всеми нейронами слоя. На этот вход всегда "возбужденный" сигнал. Присвоим такому входу номер 0. Тогда

$$s = \sum_{i=1}^n x_i w_i - T \quad (11)$$

где $w_0 = -T$, $x_0 = 1$.

Из чего следует, что отличие формулы (1) от формулы (12) в способе нумерации входов.

Все задачи, которая решает НС можно свести к классификации. Грубо говоря задача НС определить к какому классу принадлежит группа входных сигналов, находящихся в n -мерном пространстве. С

математической точки зрения процесс представляет собой разбиение гиперпространство гиперплоскостями на области.

К каждой области принадлежит отдельный класс. Максимальное число классов для НС перцептронного типа не превышает 2^m , где m — число выходов сети. Однако существует ограничение на формы гиперплоскостей, иначе говоря, не все нейронные сети могут разделить n -мерное пространство на необходимое количество классов.

Например, однослойный перцептрон, с одним нейроном, изображенный на рисунке 10 не способен разделить двумерное пространство на две полуплоскости так, что бы классифицировать входные сигналы на классы А и В (см. 1).

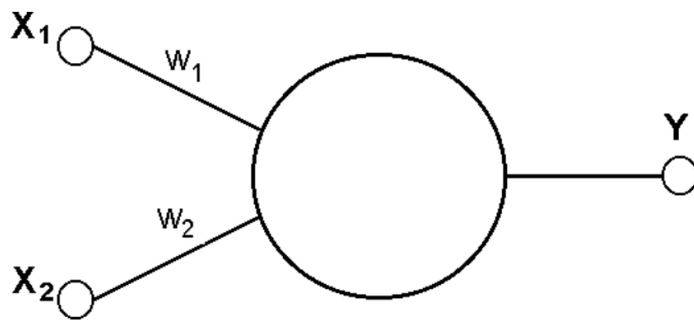


Рисунок 10 — Однонейронный перцептрон

Таблица 1 — Классификация XOR

x1	x2	A	B
0	0	•	
0	1		•
1	0		•
1	1	•	

Сеть, представленная на рисунке 10 описывает следующие уравнение:

$$x_1 w_1 + x_2 w_2 = T \quad (11)$$

Данное уравнение является прямой, которая не способна разделить плоскость таким образом, что бы группа входных сигналов x_1, x_2 принадлежали необходимым классам. На рисунке 11 показана работа НС.

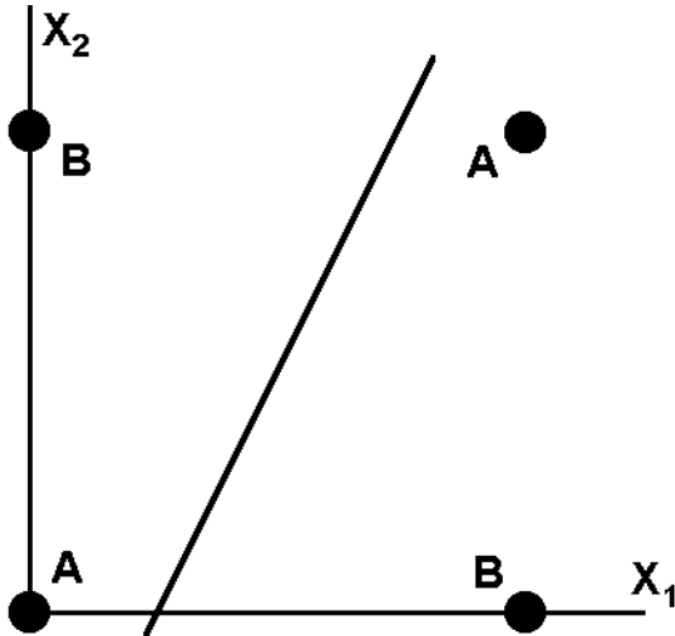


Рисунок 11 — Визуальное представление работы НС с рисунка 10

Таблица 1 является таблицей истинности для логической функции исключающего ИЛИ. Невозможность реализовать данную функцию, используя односторонний перцептрон, получила название проблемы исключающего ИЛИ.

Задачи, которые не решаются однослойной сетью, называются линейно неразделимыми[8]. Для решения таких задач используются нейронные сети с большим количеством скрытых слоев. Однако, и в таких случаях корректное разделение на классы не гарантируется. Как было сказано раньше, конфигурация НС это итеративный эмпирический процесс.

После обзора теоретических основ нейронной сети, можно более подробно рассмотреть алгоритм обучения с учителем, на основе взят перцептрон, изображенный на рисунке 8.

Алгоритм выглядит следующим образом[8]:

1) Инициализировать весовые коэффициенты случайными значениями.

ВСТАВИТЬ КОД-81

2) Подать на вход вектор входных сигналов, вычислить выходные сигналы.

3) Если выход совпадает с желаемыми значениями, перейти на шаг 4. Иначе вычислить разницу желаемым и полученным значением НС:

$$\delta = Y_l - Y \quad (11)$$

Изменить веса в соответствии с формулой:

$$w_{ij}(t+1) = w_{ij}(t) + \nu \delta x_i \quad (11)$$

где t и $t + 1$ — номера текущей и следующей итерации; ν — коэффициент скорости обучения; i — номер входа; j — номер нейрона в слое.

Веса будут увеличены, если $Y_l > I$, тем самым ошибка уменьшится. В обратном случае они будут уменьшены, и Y соответственно тоже уменьшится, приближаясь к Y_l .

4) Повтор шага 2, пока не будет достигнута желаемая точность.

На втором шаге на вход НС подаются все входные вектора из обучающей выборки в случайном порядке. Число итераций зависит от сложности задачи и конфигурации нейронной сети. Определить точное количество итераций, необходимых для корректного обучения определить невозможно.

2.1.2 Сверточные нейронные сети

2.1.3 Обнаружение объектов с применением подхода YOLO

ВСТАВИТЬ КОД-81

3 Проектирование системы

3.1 Оптимизация работы с памятью

3.1.1 Бинаризация весов

3.1.2 Оптимизация работы со слоями

ВСТАВИТЬ КОД-81

Заключение

Список использованных источников

- 1 <https://arxiv.org/pdf/1408.5093.pdf>
- 2 J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In ICML, 2014
- 3 A. Krizhevsky, I. Sutskever, and G. Hinton. ImageNet classification with deep convolutional neural networks. In NIPS, 2012
- 4 http://ronan.collobert.com/pub/matos/2011_torch7_nipsw.pdf
- 5 <http://www.lua.ru/doc/1.html>
- 6 <https://pjreddie.com/media/files/papers/xnor.pdf>
- 7 http://www.shestopaloff.ca/kyriako/Russian/Artificial_Intelligence/Some_publications.html
- 8 1. Е. Монахова, "Нейрохирурги"с Ордынки, PC Week/RE, №9, 1995.
- 9 2. Ф.Уоссермен, Нейрокомпьютерная техника, М.,Мир, 1992.
- 10 3. Итоги науки и техники: физические и математические модели нейронных сетей, том 1, М., изд. ВИНТИ, 1990.

ВСТАВИТЬ КОД-81

Приложение А
Техническое задание