

Санкт–Петербургский государственный университет
Кафедра компьютерного моделирования и многопроцессорных систем

Мирошниченко Александр Сергеевич

Выпускная квалификационная работа бакалавра

**Разработка системы распознавания речевых команд при
помощи методов машинного обучения**

Направление 01.03.02

«Прикладная математика и информатика»

Научный руководитель:
кандидат физ.-мат. наук,
доцент
Козынченко В. А.

Рецензент:
доктор физ-мат. наук,
профессор
Котина Е. Д.

Санкт-Петербург

2021 г.

Содержание

Введение	4
Постановка задачи	6
Обзор литературы	8
Глава 1. Теоретические сведения	10
1.1. Речь как объект распознавания	10
1.2. Речь в компьютерном представлении	11
1.3. Общая схема алгоритма распознавания	12
Глава 2. Описание решения	13
2.1. Предобработка	13
2.1.1 Нормализация сигнала	13
2.1.2 Удаление постоянной составляющей	13
2.1.3 Выделение начальной и конечной точек слова	14
2.2. Выделение речевых признаков	17
2.2.1 Общая схема алгоритма	17
2.2.2 Разбиение сигнала на фреймы	18
2.2.3 Получение спектра мощности сигнала	19
2.2.4 Расчёт энергий спектра мощности при помощи мел- фильтров	19
2.2.5 Логарифмирование энергий спектра мощности сигнала	22
2.2.6 Получение кепстральных коэффициентов	23
2.2.7 Приведение данных к одной размерности	24
2.3. Распознавание речевых команд	25
2.3.1 Многослойный персептрон	25
2.3.2 Свёрточная нейронная сеть	25
2.3.3 Входные и выходные данные модели	26
2.3.4 Архитектура модели нейронной сети	26

2.4. Программная реализация	28
2.5. Описание датасета	30
Глава 3. Вычислительные эксперименты	32
3.1. Описание экспериментов	32
3.2. Результаты экспериментов	34
Выводы	39
Заключение	40
Список литературы	41
Приложение	43

Введение

В современном компьютеризированном мире огромное значение имеет взаимодействие человека с компьютером - ввод и вывод информации с устройства в понятной для человека форме. Один из способов внести информацию в компьютер - записать речь через микрофон, после чего можно обрабатывать данные в памяти, которые ее представляют. Обработка может быть совершенно разной, но особенно важно распознавать слова, которые произнёс человек и давать им представление в виде текста. То есть, давать такое представление речи, как если бы она была не произнесена голосом, а напечатана при помощи клавиатуры.

Решение такой задачи может быть использовано для различных целей. К примеру, можно переписываться с другим человеком по интернету текстовыми сообщениями и при этом вообще не прикасаться к клавиатуре, управлять различными компьютерными интерфейсами при помощи голосовых команд и т.д.

Данная проблема была актуальна со времён появления компьютеров и остаётся таковой по сей день. Особенно актуальна она стала в последнее время, когда появились качественные микрофоны, возросла мощность вычислительных устройств, увеличилось количество информации, которой обмениваются люди через интернет.

Изначально, для решения данной задачи применялись такие алгоритмы, как скрытые Марковские модели, методы динамического программирования, методы дискриминантного анализа, основанные на Байесовской дискриминации и другие. Но с появлением нейронных сетей и многочисленных экспериментов с их использованием выяснилось, что задачу распознавания речи можно решать и при помощи нейросетевого подхода. И хоть сами нейронные сети появились ещё в прошлом веке, их популярность возросла только в последнее время, в связи с ростом мощности компьютеров.

Краткое содержание глав:

В главе 1 рассмотрены теоретические сведения о звуке, речевых признаках и общая схема алгоритма распознавания звука.

В главе 2 рассмотрен алгоритм распознавания речи, его программная реализация и датасет. Описаны подзадачи: предобработка звукового сигнала, выделение речевых признаков, распознавание речи.

В главе 3 приведены результаты вычислительных экспериментов.

Постановка задачи

В данной работе ставится задача распознавания речи. Под распознаванием речи понимается нахождение отображения множества речевых команд во множество соответствующих им скалярных меток.

Пусть $X = \{x_0, \dots, x_{p-1}, \dots\}$ - множество речевых команд. Оно состоит из векторов амплитуд $x_i = \{x_i^0, \dots, x_i^{k_i}\}$, $i = \overline{0, \infty}$. Здесь k_i - количество записанных амплитуд в i -м объекте. Само по себе множество команд бесконечно, однако известны значения первых p элементов, обозначаемых через $\hat{X} = \{x_0, \dots, x_{p-1}\}$.

Пусть $Y = \{y_0, \dots, y_{p-1}, \dots\}$ - множество скалярных меток, а $\hat{Y} = \{y_0, \dots, y_{p-1}\}$ - множество известных скалярных меток для первых p речевых команд. Таким образом известно отображение $\hat{Z} = \hat{X} \rightarrow \hat{Y}$, описываемое парами значений $\hat{Z} = \{(x_0, y_0), \dots, (x_{p-1}, y_{p-1})\}$.

Необходимо разработать алгоритм, который бы строил отображение $Z = X \rightarrow Y$.

Поставленная задача разбивается на следующие подзадачи:

- изучение математических методов и моделей для решения поставленной задачи распознавания речи;
- разработка и реализация алгоритма предобработки исходных данных;
- разработка и реализация алгоритма распознавания речевых команд;
- проведение вычислительных экспериментов и исследование работоспособности и эффективности алгоритма распознавания речевых команд;
- реализация математической модели программными методами;

- выводы об эффективности исследуемого подхода к решению задачи распознавания речи для прикладного применения.

Под множеством речевых команд понимается набор слов-команд для управления медиаплеером. Базовый набор команд составляет необходимый минимум для управления программным интерфейсом. В рамках данной работы сам интерфейс управления медиаплеером не рассматривается.

Обзор литературы

Далее в хронологическом порядке описаны наиболее важные моменты и работы, связанные с поставленной задачей.

Первые шаги в распознавании речи были сделаны в 1952 году. Тогда трое исследователей Bell Labs - Стивен Балашек, Рулон Биддалф и Кей Дэвис - представили публике первый в истории аппарат Audrey, способный распознавать человеческую речь [1]. Это была система, позволявшая распознавать только цифры.

В 1957 году Фрэнк Розенблатт предложил математическую модель восприятия информации мозгом - персептрон. В своей статье [2], а позднее и в своей книге [3], он описал принципы работы модели. Этот момент можно считать появлением нейронных сетей как таковых. На сегодняшний день их развитие ушло вперёд, однако принципы, заложенные Розенблаттом, являются основополагающими в этой области.

В работе 1963 года [4], которую опубликовали Богерт, Хили и Тьюки, впервые описан кепстр для анализа геологических данных. Позднее этот подход будет использован в анализе речевого сигнала. В десятой главе книги Оппенгейма [8] описано его применение в распознавании речи.

В 1967 году впервые применяется спектр для анализа звуковых сигналов [5].

В 1971 году появилась одна из первых моделей распознавания большого словаря речевых команд. Она была представлена на конкурсе DARPA и носила название Harpy [7]. Методы, которые были использованы при её реализации актуальны и по сей день.

Для выделения речевых признаков впервые в 1980 году использован алгоритм мел-частотных кепстральных коэффициентов [9]. Он основан на кепстре, который Богерт, Хили и Тьюки использовали для геоанализа.

В 1995 году Ян ЛеКун предложил модель свёрточной нейронной сети

для распознавания изображений и речи [11]. Все современные архитектуры нейронных сетей в своём большинстве - модификации свёрточной сети.

Глава 1. Теоретические сведения

В этой главе рассмотрены теоретические сведения о человеческой речи, речевых признаках и общая схема алгоритма распознавания речи.

1.1 Речь как объект распознавания

Речь представляет собой акустическую волну, которая излучается системой органов: лёгкими, бронхами и трахеей, а затем преобразуется в головном тракте. Если предположить, что источники возбуждения и форма головного тракта относительно независимы, то речевой аппарат человека можно представить в виде совокупности генератора тоновых сигналов и фильтров. На рисунке 1 изображена схема речеобразования.

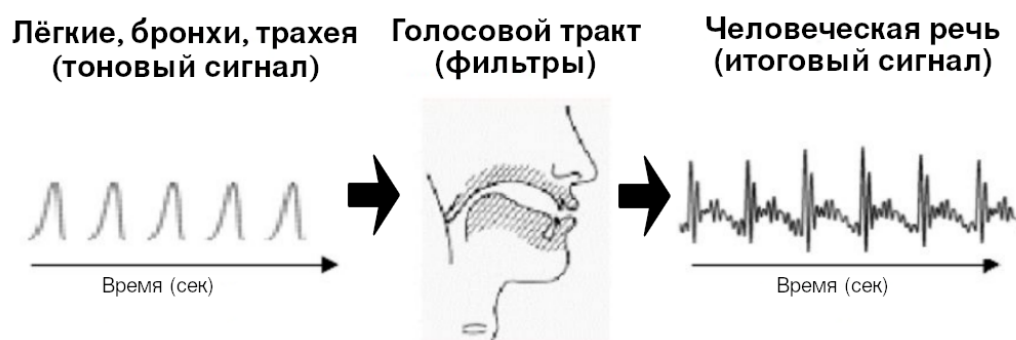


Рис. 1: Речеобразование

Распознавание речи происходит на основе анализа определённых признаков, которые причастны определенным звукам. За создание этих особенностей отвечает голосовой тракт человека. Поэтому одной из подзадач предобработки звука является выделение той части сигнала, которая была сформирована именно голосовым трактом.

Произнесённая речь поступает в приёмники звука. В случае человеческого уха происходит следующее. Звуковые волны проходят через наружное ухо в среднее и вызывают вибрацию барабанной перепонки. Колебания с

барабанной перепонки передаются на маленькие слуховые косточки в среднем ухе. А со слуховых косточек - во внутреннее ухо. Когда эти колебания достигают улитки, они воздействуют на специальные клетки - волосковые. Волосковые клетки преобразуют колебания в электрические нервные импульсы. Слуховой нерв соединяет улитку с центрами слуха в головном мозге. Когда электрические нервные импульсы достигают головного мозга, они воспринимаются как звук и обрабатываются.

В случае с компьютером и подключённым к нему микрофоном, принцип схожий. Звуковые волны преобразуются микрофоном в электрический сигнал, который далее обрабатывается АЦП в составе звуковой платы компьютера. На выходе - цифровые данные, подлежащие обработке в рамках поставленной задачи.

Речь можно представить в виде последовательности предложений, а их в свою очередь в виде последовательности слов. Слова состоят из фонем. В общем случае речь непрерывна, то есть слова не отделяются друг от друга паузами, за исключением требующих этого особенностей произносимой речи. В этой работе не рассматривается общий случай. Здесь рассмотрен частный случай, когда речь состоит из отдельных слов, отделяемых друг от друга тишиной. Под тишиной понимается не полное отсутствие звука, а временные промежутки, в течение которых говорящий молчит, делает паузы. Этот выбор был обусловлен командным типом системы распознавания, которая работает с отдельными словами.

1.2 Речь в компьютерном представлении

Каждая речевая команда с точки зрения звуковой записи в компьютере - это набор амплитудных значений, полученных с микрофона, преобразованных АЦП в числовые и записанных в звуковой файл формата wav.

1.3 Общая схема алгоритма распознавания

На рисунке 2 представлена общая схема работы алгоритма распознавания речевых команд. Алгоритм разделен на два основных блока, обозначенных на рисунке как Блок 1 и Блок 2. На вход алгоритму поступает набор амплитудных значений и соответствующая частота дискретизации в формате wav. На выходе алгоритма - текстовое представление команды.

- Блок 1 - блок, отвечающий за первоначальную обработку данных и приведение их к единому формату.
- Блок 2 - блок, отвечающий за классификацию унифицированных данных. Этот блок может работать в двух режимах: обучения и непосредственной работы. В режиме обучения происходит корректировка параметров алгоритма, которые влияют на конечный результат. В режиме непосредственной работы этого не происходит.

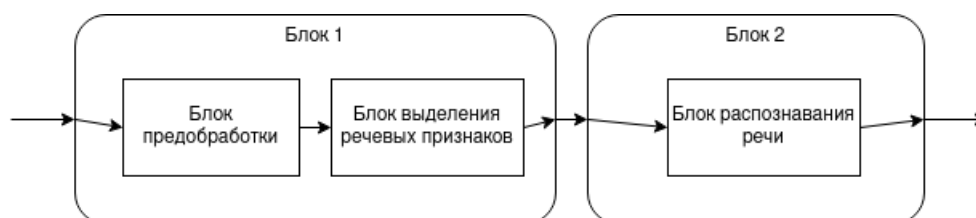


Рис. 2: Общая схема работы алгоритма распознавания речевых команд

Глава 2. Описание решения

В этой главе рассмотрен весь алгоритм распознавания речевых команд. Для удобства понимания названия параграфов расположены в том же порядке, что и этапы в самом алгоритме.

2.1 Предобработка

Каждая команда записана в звуковой wav файл. В каждом файле - набор амплитудных значений (далее - звуковой сигнал), которые были получены в результате записи команды дикторами.

2.1.1 Нормализация сигнала

Сначала проводится нормализация амплитуд. Каждое значение амплитуды приводится к такому значению, чтобы максимум среди всех амплитудных значений, содержащихся в звуковом сигнале был равен единице по формуле:

$$\bar{x}_i = \frac{x_i}{\max_j |x_j|}, \quad i = \overline{0, p-1}, \quad j \in [0, p-1], \quad (1)$$

где x - значение амплитуды, \bar{x} - новое значение амплитуды, p - количество амплитудных значений, содержащихся в звуковом сигнале.

Таким образом, все значения амплитуд принимают значения в диапазоне $[0, 1]$.

2.1.2 Удаление постоянной составляющей

Постоянная составляющая (DC-offset) - это смещение амплитуды сигнала на некоторую постоянную величину. Она возникает в цифровом сигнале, полученном с АЦП, из-за разницы напряжения между звуковой картой и устройством ввода. Данный эффект является помехой, от которой нужно избавиться. Для этого необходимо вычесть из каждого значения амплитуды

среднее арифметическое всех значений амплитуд по формуле:

$$\bar{x}_i = x_i - \sum_{j=0}^{p-1} x_j, \quad i = \overline{0, p-1}, \quad (2)$$

где x - значение амплитуды, полученное на этапе нормализации, \bar{x} - новое значение амплитуды, p - количество амплитудных значений, содержащихся в звуковом сигнале.

Все \bar{x}_i переобозначаются как x_i .

2.1.3 Выделение начальной и конечной точек слова

Каждый звуковой сигнал содержит в себе помимо фрагментов речевой команды ещё и фрагменты тишины. Очень важно отделить фрагменты речевой команды от фрагментов тишины, так как вся информация о команде содержится именно во фрагментах речевой команды.

Для того, чтобы выделить речевую команду и «обрезать» тишину в начале и в конце записи, используется алгоритм, описанный в статье [6]. Каждый звуковой сигнал разбивается на фреймы - наборы амплитуд, каждый длительностью 20 мс. Начала фреймов расположены с периодичностью 10 мс. Таким образом, фреймы пересекаются между собой. Это обеспечивает целостность обработки звукового сигнала, позволяя обрабатывать важные фонемообразующие особенности.

Затем для каждого фрейма вычисляется мгновенная энергия:

$$E_k = \sum_{m=0}^{N-1} x_{k_m}^2, \quad k = \overline{0, z-1}, \quad (3)$$

где z - количество фреймов для конкретного звукового сигнала, N - длина одного фрейма (количество амплитуд в одном фрейме).

Вычисление функции мгновенной энергии имеет значительный недостаток. У неё велика чувствительность к относительно большим значениям

амплитуды из-за возведения их во вторую степень. Это ведёт к искажению соотношений амплитудных значений звукового сигнала друг к другу. Поэтому функция мгновенной энергии переопределяется как:

$$E_k = \sum_{m=0}^{N-1} |x_{k_m}|, \quad k = \overline{0, z-1}, \quad (4)$$

где z - количество фреймов для конкретного звукового сигнала, N - длина одного фрейма.

После того, как посчитаны мгновенные энергии для каждого фрейма, вычисляются нижнее и верхнее пороговые значения:

$$\begin{aligned} I_1 &= 0.03 \cdot (MX - MN) + MN \\ I_2 &= 4 \cdot MN \\ ITL &= \min(I_1, I_2) \\ ITU &= 10 \cdot ITL, \end{aligned} \quad (5)$$

где MN , MX - минимум и максимум мгновенной энергии среди всех фреймов соответственно, ITL , ITU - нижнее и верхнее пороговое значение.

Происходит поиск фрейма, с которого начинается речевая команда, начиная с самого первого фрейма. Фрейм, в котором значение мгновенной энергии превышает ITL , предварительно помечается как начало речевой команды. Затем, начиная с этого помеченного фрейма, происходит поиск фрейма, в котором значение мгновенной энергии превышает ITU . Если значение мгновенной энергии для какого-то фрейма во время последнего поиска меньше ITL , то этот фрейм становится предварительным началом речевой команды.

Аналогично происходит поиск конца речевой команды в звуковом сигнале, но поиск по фреймам происходит не с начала сигнала, а с конца.

После этого этапа имеются два предварительно помеченных фрейма m_{begin} , m_{end} - начало и конец речевой команды в звуковом сигнале соответ-

ственно.

Функция мгновенной энергии, определённая формулой (4), хорошо справляется с отделением звонких звуков от тишины. Но глухие она отделяет плохо. Поэтому используется вторая характеристика для доопределения начала и конца слова - число переходов через ноль. Это количество таких случаев, когда соседние значения амплитуд имеют противоположные знаки. Определяется формулой:

$$Z_k = \frac{1}{2} \sum_{m=1}^{N-1} |\text{sgn}(x_{k_{m-1}}) - \text{sgn}(x_{k_m})|, \quad k = \overline{0, z-1}, \quad (6)$$

где z - количество фреймов для конкретного звукового сигнала, N - длина одного фрейма.

Подразумевается, что первые 100 мс звукового сигнала - это тишина, и речевая команда начинается позднее.

Вычисляется среднее значение переходов через ноль в течение первых 100 мс (7) и среднее квадратическое отклонение количества переходов через ноль в течение первых 100 мс (8):

$$IZC = \frac{1}{z} \sum_{k=0}^{z-1} Z_k \quad (7)$$

$$\sigma_{IZC} = \sqrt{\frac{1}{z} \sum_{k=0}^{z-1} (Z_k - IZC)^2}, \quad (8)$$

где z - количество фреймов для конкретного звукового сигнала. Затем вычисляется значение пороговой функции числа переходов через ноль по формуле:

$$IZCT = \min(IF, IZC + 2\sigma_{IZC}), \quad (9)$$

где IF - фиксированное количество переходов через ноль. В данном случае оно составляет 25 пересечений за 10 мс, то есть $IF = 2.5$.

Далее происходит уточнение точек начала и конца речевой команды в звуковом сигнале. Начиная от фрейма m_{begin} влево происходит поиск фреймов, у которых число переходов через ноль выше порогового значения. Поиск происходит на расстоянии 25 фреймов, так как производится уточнение границ слова. Если пороговое значение было превышено 3 или более раз, то фрейм, где это произошло впервые, помечается как начало речевой команды и обозначается как r_{begin} . Если пороговое значение было превышено менее 3-х раз, то метка m_{begin} переобозначается как r_{begin} .

Аналогично от фрейма m_{end} происходит поиск вправо для уточнения точки конца речевой команды, которая обозначается как r_{end} .

В результате уточнения точек начала и конца речевой команды имеются 2 помеченных фрейма - r_{begin}, r_{end} . Сигнал обрезается, и в нем остаётся только речевая команда в виде набора фреймов $[r_{begin}, \dots, r_{end}]$. Количество получившихся фреймов обозначается как u , а сами фреймы перенумеруются следующим образом: $[r_0, \dots, r_{u-1}]$.

2.2 Выделение речевых признаков

2.2.1 Общая схема алгоритма

Для того, чтобы выделить речевые признаки, используется алгоритм мел-частотных кепстральных коэффициентов [9] (далее - MFCC). Он является одним из стандартных подходов к решению поставленной задачи. Состоит MFCC из нескольких шагов:

1. Для каждого звукового сигнала проделать шаги:
 - a. Разбить сигнал на фреймы.
 - b. Для каждого фрейма проделать шаги:
 - I. Получить спектр мощности сигнала.
 - II. Составить набор мел-фильтров.

- III. Рассчитать энергии спектра мощности при помощи мел-фильтров.
- IV. Прологарифмировать энергии спектра мощности, полученные при помощи мел-фильтров.
- V. Получить вектор кепстральных коэффициентов, применяя дискретное косинусное преобразование к результату предыдущего шага.

с. Объединить MFCC векторы коэффициентов в матрицу.

2. Привести матрицы коэффициентов MFCC каждого звукового сигнала к одной размерности. Для этого дополнить их нулями слева до необходимой длины. Таким образом, размерности всех матриц приводятся к максимальной из всех размерностей.

2.2.2 Разбиение сигнала на фреймы

Амплитуда сигнала постоянно изменяется. Для упрощения предполагается, что в течение коротких промежутков времени амплитуда сигнала является статистически постоянной. Поэтому сигнал разделяется на фреймы длительностью 20-40 мс (в рамках данной работы - 20 мс). Если взять длину фрейма короче этой длительности, то количество амплитудных значений в каждом фрейме будет недостаточным для построения достоверной спектральной оценки для всего сигнала. Если взять длину фрейма длиннее чем 20-40 мс, то сигнал уже не будет статистически постоянным на этом промежутке времени.

Так как на выходе алгоритма выделения начальной и конечной точек слова получается набор фреймов необходимой длительности, то шаг разбиения сигнала на фреймы опускается.

2.2.3 Получение спектра мощности сигнала

Следующий шаг - вычисление спектра мощности сигнала для каждого фрейма. Этот шаг обоснован особенностью восприятия звука человеком. В человеческом ухе улитка вибрирует в разных зонах в зависимости от частоты входного звукового сигнала. И в зависимости от того, какая зона (волосковых клеток) в улитке вибрирует, соответствующие нервы информируют мозг об определённых частотах, которые содержатся в звуковом сигнале. Анализ спектра мощности сигнала в числовом представлении позволяет провести аналогичное разложение сигнала на частоты, которые в нём присутствуют.

Спектр мощности сигнала S_{k_m} в k -ом фрейме - результат дискретного преобразования Фурье:

$$S_{k_m} = \sum_{n=0}^{N-1} x_{k_n} \cdot e^{\frac{-2\pi i}{N} mn}, \quad m = \overline{0, N-1}, \quad k = \overline{0, u-1}, \quad (10)$$

где u - количество фреймов, получившееся после выделения начальной и конечной точек речевой команды, N - длина одного фрейма, i - мнимая единица.

2.2.4 Расчёт энергий спектра мощности при помощи мел-фильтров

Спектр мощности звукового сигнала содержит в себе излишнюю информацию, которая не требуется для алгоритма распознавания речи. В частности, улитка человеческого уха не может отличить две близкие по значениям частоты во входном звуковом сигнале. Данный эффект более выражен для более высоких частотных значений. В связи с этим, спектр мощности делится на определённые частотные интервалы, для каждого из которых рассчитывается энергия. Это производится при помощи мел-фильтров, выраженных в виде оконной функции. Частотный интервал первого фильтра - очень узкий. Он даёт представление о существующем количестве энергии около частоты 0

Гц. Более высоким частотным значениям соответствуют более широкие интервалы фильтров. Это объясняется тем, что для более высоких частотных значений разница между амплитудами соседних частот становится все менее значимой. Мел-шкала позволяет точно рассчитать ширины частотных интервалов для фильтров.

Мел - психофизическая единица высоты звука. Используемые формулы для перевода частоты из шкалы Герц в шкалу мел и обратно описаны в книге [10] на стр. 150:

$$mel(hz) = 1127 \cdot \ln\left(1 + \frac{hz}{700}\right) \quad (11)$$

$$hz(mel) = 700 \cdot \left(e^{\frac{mel}{1127}} - 1\right) \quad (12)$$

На рисунке 3 изображено сравнение шкал мел и Гц.

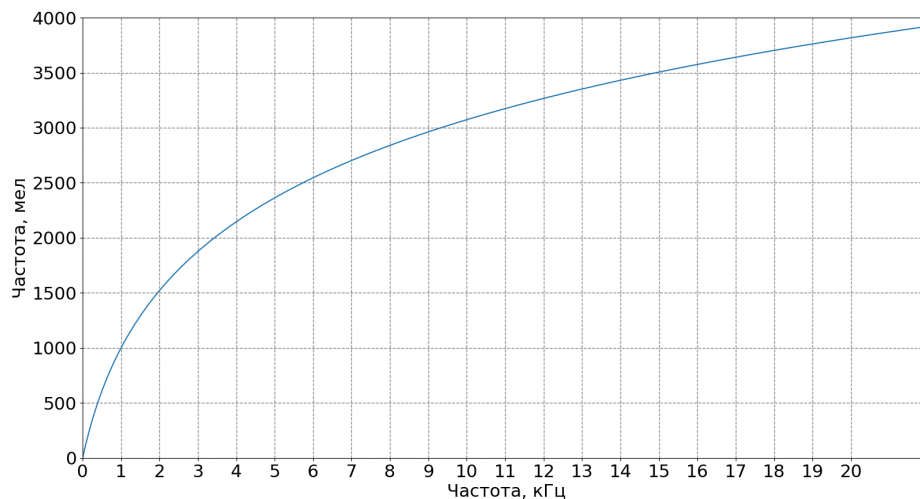


Рис. 3: Сравнение шкал мел и Гц

Составляются треугольные мел-фильтры в виде оконной функции:

$$H(v, b) = \begin{cases} 0, & b < f(v) \\ \frac{b - f(v)}{f(v+1) - f(v)}, & f(v) \leq b < f(v+1) \\ \frac{f(v+2) - b}{f(v+2) - f(v+1)}, & f(v+1) \leq b < f(v+2) \\ 0, & f(v+2) < b \end{cases}, \quad (13)$$

для которой f определяется как:

$$f(a) = \frac{N}{w}hz(mel(f_{min}) + a \frac{mel(f_{max}) - mel(f_{min})}{Q + 1}), \quad (14)$$

где mel и hz - функции, определённые формулами (11) и (12) соответственно, w - частота дискретизации звукового сигнала, f_{min} , f_{max} - нижний и верхний пороги частотного диапазона соответственно, N - длина одного фрейма, Q - количество мел-фильтров. Параметр Q рекомендовано выбирать в диапазоне 20-40 (в рамках данной работы - 26) [9].

На рисунке 4 в качестве примера изображена оконная функция для звукового сигнала с параметрами $w = 44100$ Гц, $f_{min} = 0$ Гц, $f_{max} = 22050$ Гц, $N = 1024$, $Q = 13$.

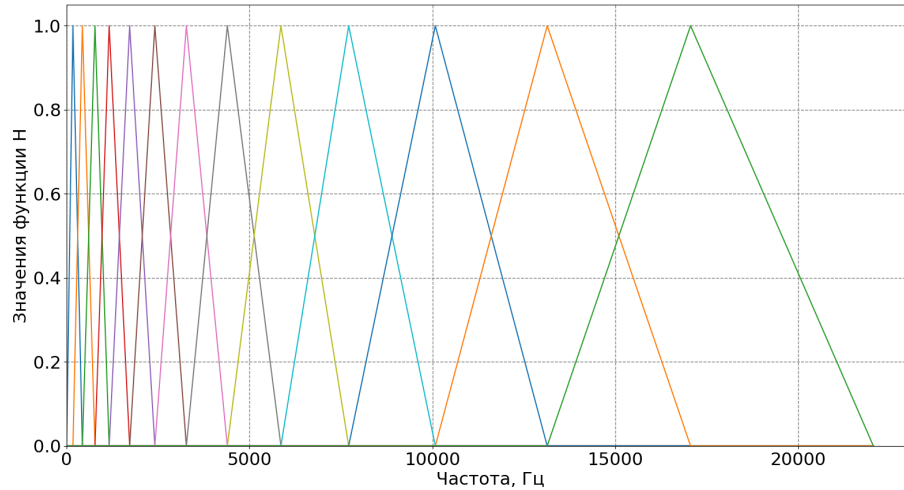


Рис. 4: Оконная функция

Считаются энергии спектра мощности при помощи мел-фильтров:

$$W_{k_q} = \sum_{m=0}^{N-1} |S_{k_m}|^2 \cdot H(q, m), \quad q = \overline{0, Q-1}, \quad k = \overline{0, u-1}, \quad (15)$$

где u - количество фреймов, получившееся после выделения начальной и конечной точек речевой команды, Q - выбранное количество мел-фильтров.

2.2.5 Логарифмирование энергий спектра мощности сигнала

После получения энергий спектра мощности сигнала при помощи мел-фильтров, происходит логарифмирование последних. Обосновано это тем, что громкость речи нелинейно зависит от количества содержащейся в ней энергии. Чтобы удвоить громкость произносимой речи, человеку необходимо приложить в 8 раз больше энергии. Это означает, что фрагменты речи, имеющие большое различие в соответствующих им значениях энергий могут иметь примерно одинаковую громкость, если речь начинается громко. Операция логарифмирования производит корреляцию между громкостью и энергией звукового сигнала более похожей на человеческую.

Логарифмирование энергий спектра мощности, полученных при помо-

щи мел-фильтров:

$$V_{k_q} = \ln(W_{k_q}), \quad q = \overline{0, Q-1}, \quad k = \overline{0, u-1}, \quad (16)$$

где u - количество фреймов, получившееся после выделения начальной и конечной точек слова, Q - выбранное количество мел-фильтров.

2.2.6 Получение кепстральных коэффициентов

Последним шагом является вычисление дискретного косинусного преобразования от прологарифмированных энергий спектра мощности сигнала. Так как мел-фильтры пересекаются между собой, значения энергий, полученных при помощи соответствующих мел-фильтров, сильно коррелируют друг с другом. Дискретное косинусное преобразование декоррелирует их. Коэффициенты этого преобразования называются кепстральными, так как все вместе они образуют кепстр сигнала. В качестве результата берутся первые 12-13 (в рамках данной работы - 12) этих коэффициентов. Исследователями в этой области эмпирическим путём показано, что использование последующих коэффициентов ухудшает точность распознавания в сравнении с использованием первых 12-13.

Стоит отметить, что кепстр определяется в виде прямого преобразования Фурье от логарифма спектра мощности сигнала:

$$C(x(t)) = F(\ln[F(x[t])]), \quad (17)$$

где $x(t)$ - входной сигнал, F - функция прямого преобразования Фурье. В книге [8] можно найти описание кепстра на стр. 367-380.

Однако в алгоритме MFCC используется дискретное косинусное преобразование как частный случай прямого преобразования Фурье по вышеописанным причинам.

Получение итоговых коэффициентов MFCC происходит по формуле:

$$c_{k_n} = \sum_{m=0}^{Q-1} V_{k_m} \cos\left(\frac{\pi}{Q}\left(m + \frac{1}{2}\right)n\right), \quad n = \overline{0, d-1}, d \leq Q, \quad k = \overline{0, u-1}, \quad (18)$$

где u - количество фреймов, получившееся после выделения начальной и конечной точек слова, Q - выбранное количество мел-фильтров, d - выбранное количество коэффициентов MFCC (в рамках данной работы - 12) [9], π - математическая константа.

Таким образом, строится матрица для каждой звуковой дорожки следующего вида:

$$C_{u \times d} = \begin{pmatrix} c_{00} & c_{01} & \dots & c_{0(d-1)} \\ c_{10} & c_{11} & \dots & c_{1(d-1)} \\ \vdots & \vdots & \ddots & \vdots \\ c_{(u-1)0} & c_{(u-1)1} & \dots & c_{(u-1)(d-1)} \end{pmatrix}$$

где u - количество фреймов, получившееся после выделения начальной и конечной точек слова, d - выбранное количество коэффициентов MFCC.

2.2.7 Приведение данных к одной размерности

Среди всех значений u существует максимальное u_{max} . Для каждой матрицы, соответствующей определённой звуковой дорожке, сделаем следующее:

- если $u < u_{max}$, то соответствующая матрица C дополняется нулями слева:

$$C_{u_{max} \times d} = \begin{pmatrix} 0 & \dots & 0 & c_{00} & c_{01} & \dots & c_{0(d-1)} \\ 0 & \dots & 0 & c_{10} & c_{11} & \dots & c_{1(d-1)} \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & c_{(u-1)0} & c_{(u-1)1} & \dots & c_{(u-1)(d-1)} \end{pmatrix}$$

- если $u = u_{max}$, то матрица C остаётся без изменений.

Таким образом, все матрицы приводятся к одной размерности $u_{max} \times d$.

2.3 Распознавание речевых команд

Для распознавания речевых команд используются технологии нейронных сетей. В данной работе рассматриваются два типа нейронных сетей: многослойный персептрон и свёрточная сеть. Производится сравнение производительности этих двух типов при различных параметрах обучения.

2.3.1 Многослойный персептрон

Этот вид нейронной сети описан в третьей части книги Фрэнка Розенблатта [3], который первый предложил модель персептрона. Она состоит из входного слоя, полносвязных слоев и выходного слоя.

2.3.2 Свёрточная нейронная сеть

Данный тип нейронной сети был предложен Яном ЛеКуном [11]. Это многослойная сеть, позволяющая обеспечить устойчивость распознавания к инвариантным изменениям данных за счёт общих весов и локального рецептивного поля.

Входной слой сети состоит из одной плоскости. Его размерность совпадает с размерностью входных данных.

Последующие слои - свёрточные. Каждый свёрточный слой состоит из нескольких плоскостей нейронов, которые известны как карты признаков. Каждый нейрон в свёрточном слое соединён с небольшой областью предыдущего слоя. В этом заключается принцип локального рецептивного поля.

После каждого свёрточного слоя, который получил локальные признаки, следует пулинг слой. Его задача состоит в понижении размерности данных.

После всех свёрточных слоёв следует выпрямляющий слой, который преобразует данные к вектору.

Далее следуют полносвязные слои. Иногда добавляется дропаут слой, который отключает некоторые случайные нейроны в процессе обучения. Это позволяет бороться с переобучением сети.

Выходной слой имеет размерность требуемых выходных данных.

2.3.3 Входные и выходные данные модели

Входной тензор модели - матрица MFCC коэффициентов для соответствующей команды. Его размерность - $u_{max} \times d$. Для составленного датасета $u_{max} = 400, d = 13$.

Выходной тензор имеет размерность $g \times 1$, где g - количество возможных команд для распознавания. Для составленного датасета $g = 11$. Скалярная метка распознанной команды соответствует индексу максимального элемента выходного тензора. Индексация в выходном тензоре начинается с 0.

2.3.4 Архитектура модели нейронной сети

Архитектура свёрточной нейронной сети приведена на рисунке 6. Как видно из рисунка, сеть включает в себя:

- Входной слой InputLayer (размерность входных данных - 400×13)
- Слой свёртки Conv2D (32 нейрона, размерность ядра свёрки - 5×5 , функция активации - ReLu)
- Слой пулинга AveragePooling2D (размерность пула - 2×2)
- Слой свёртки Conv2D (64 нейрона, размерность ядра свёрки - 5×5 , функция активации - ReLu)
- Слой пулинга AveragePooling2D (размерность пула - 2×2)

- Выпрямляющий слой Flatten
- Полносвязный слой Dense (128 нейронов, функция активации - ReLu)
- Слой дропаута Dropout (процент исключения случайных нейронов от общего числа в слое - 30%)
- Выходной полносвязный слой Dense (11 нейронов, функция активации - Softmax)

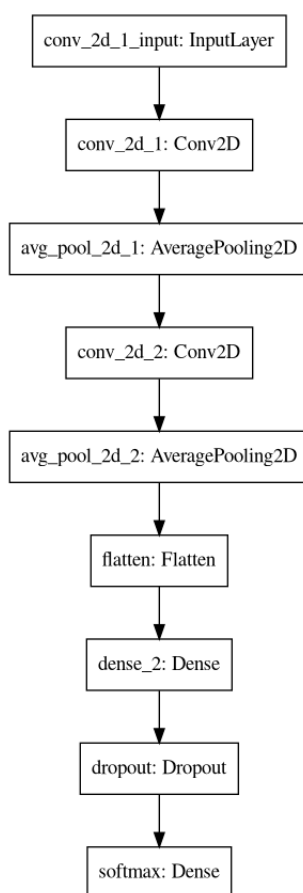


Рис. 5: Структура свёрточной нейронной сети

Архитектура свёрточной нейронной сети приведена на рисунке 6. Как видно из рисунка, сеть включает в себя:

- Входной слой InputLayer (размерность входных данных - 400×13)
- Выпрямляющий слой Flatten

- Полносвязный слой Dense (256 нейронов, функция активации - ReLu, регуляризатор ядра - L2 с параметром $\lambda = 0.00001$)
- Полносвязный слой Dense (128 нейронов, функция активации - ReLu, регуляризатор ядра - L2 с параметром $\lambda = 0.00001$)
- Полносвязный слой Dense (128 нейронов, функция активации - ReLu, регуляризатор ядра - L2 с параметром $\lambda = 0.00001$)
- Полносвязный слой Dense (64 нейрона, функция активации - ReLu, регуляризатор ядра - L2 с параметром $\lambda = 0.00001$)
- Выходной полносвязный слой Dense(11 нейронов, функция активации - Softmax)

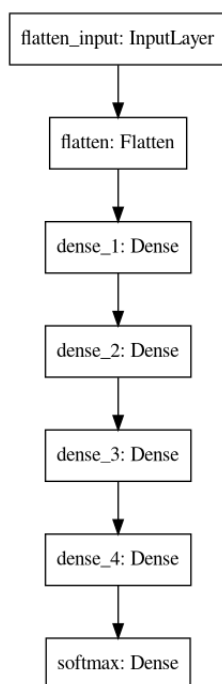


Рис. 6: Структура многослойного персептрона

2.4 Программная реализация

Построение программного интерфейса для блоков предобработки и распознавания команд было произведено при помощи языка программирования Python 3.8.

Блок предобработки был реализован при помощи библиотек `numpy`, `scipy`.

Блок распознавания был реализован при помощи библиотек `jupyterlab`, `Keras`, `numpy`. Для анализа результатов обучения нейронных сетей были использованы библиотеки `pandas`, `matplotlib`, `seaborn`.

Программный комплекс разделен на две основные части:

- Корневой файл «`preprocessing.py`», в котором реализована предобработка. Весь вспомогательный функционал вынесен в отдельный модуль «`helpers`».
- Корневой файл «`boxu.ipynb`», в котором реализовано обучение и тестирование нейронных сетей. Весь вспомогательный функционал также вынесен в отдельный модуль «`helpers`».

В корневой директории «`recorded_audio`» хранятся записанные шестью разными дикторами звуковые файлы, содержащие команды. Каждый файл содержит в своём названии скалярную метку речевой команды, которая в нём записана.

В корневую директорию «`data`», в процессе предобработки файлов из директории «`recorded_audio`», сохраняются файлы:

- «`speaker $\{i\}$ _data.npy`»¹
- «`speaker $\{i\}$ _labels.npy`»¹,

где $i = \overline{1, 6}$ - номер диктора. Их описание приводится в разделе 2.5.

В корневой директории «`logs`» в процессе предобработки данных создаётся файл «`log.log`», в который записывается краткий отчёт об обработке каждого звукового файла.

В корневой директории «`model_checkpoints`» в процессе обучения нейронной сети создаются дочерние директории «`experiment $\{i\}$ »1, $i = \overline{1, 3}$ -`

номер эксперимента. В дочерних директориях создаются свои дочерние директории «nn_type». Параметр nn_type принимает значения «cnn» при обучении свёрточной сети и «mlp» при обучении многослойного персептрона. В директорию «nn_type» сохраняются веса сети с наилучшей точностью распознавания на валидационных данных.

Для создания датасета было разработано веб-приложение при помощи языков Javascript, HTML, CSS и фреймворка NodeJS. Оно позволяет записывать речевые команды и сохранять их в файлы в нужном для программы предобработки wav-формате с названиями, содержащими соответствующие скалярные метки. В распространённых программах записи звука много времени тратится на сохранение с правильными параметрами звуковой дорожки. Сравнительно с ними созданное веб-приложение позволило в короткие сроки записать необходимое количество дикторов и сократило временные затраты на запись образцов речи в несколько раз. Код веб-приложения доступен в приложении.

2.5 Описание датасета

Составленный датасет состоит из 11 команд, записанных шестью дикторами. Каждый диктор работал с командами : «back», «down», «menu», «off», «on», «open», «play», «power», «stop», «up», «volume». В таблице 1 указаны типы голосов дикторов и данные о количестве записанных команд.

¹Фигурные скобки не являются частью имени файла или директории.

Диктор	Тип голоса	Кол-во звук. дорожек на каждую команду	Сумм. кол-во звук. дорожек
speaker1	Мужской	50	550
speaker2	Мужской	40	440
speaker3	Мужской	40	440
speaker4	Мужской	40	440
speaker5	Мужской	50	550
speaker6	Женский	50	550

Таблица 1: Типы дикторов и данные о количестве записанных команд

Датасет предварительно разделяется на тренировочную и тестовую части. На тренировочную часть отводится 70% данных каждого диктора, на тестовую часть - 30%.

Матрицы коэффициентов MFCC, полученных на этапе 2.2.6 объединяются в массив и записываются в файл «speaker $\{i\}$ _data.npy»¹ в виде numpy массива. В файл «speaker $\{i\}$ _labels.npy»¹ записываются скалярные метки команд в виде массива, в том же порядке, что и соответствующие им матрицы коэффициентов MFCC. Здесь $i = \overline{1,6}$ - номер диктора. Соответствие произнесённых диктором команд и их скалярных меток представлено в таблице 2.

Команда	back	down	menu	off	on	open	play	power	stop	up	volume
Скалярная метка	0	1	2	3	4	5	6	7	8	9	10

Таблица 2: Соответствие произнесённых диктором команд и их скалярных меток

¹Фигурные скобки не являются частью имени файла или директории.

Глава 3. Вычислительные эксперименты

Используемые обозначения:

all_speakers = дикторы speaker1, speaker2, speaker3, speaker4, speaker5, speaker6;

all_male_speakers = дикторы speaker1, speaker2, speaker3, speaker4, speaker5.

3.1 Описание экспериментов

Для каждого из двух типов нейронных сетей, многослойного персептрона и свёрточной сети, проводится три эксперимента:

1. нейронная сеть обучается на первом дикторе с мужским голосом, тестирование производится сначала на каждом дикторе, а потом на всех вместе;
2. нейронная сеть обучается на всех дикторах с мужским голосом, тестирование производится сначала на каждом дикторе, а потом на всех вместе;
3. нейронная сеть обучается на всех дикторах, тестирование производится сначала на каждом дикторе, а потом на всех вместе.

Обучение каждого из двух типов нейронных сетей производится со следующими параметрами:

- после каждой эпохи тренировочные данные перемешиваются;
- 15% тренировочных данных в каждой эпохе - валидационные;
- метрика для оценки эффективности - точность (в библиотеке Keras называется accuracy);
- функция потерь - категориальная кросс-энтропия (в библиотеке Keras называется categorical_crossentropy);

- алгоритм оптимизации - Adam;
- количество эпох обучения сети - 50;
- если в течение 20 эпох значение функции потерь на валидационных данных не уменьшается, то происходит ранняя остановка обучения сети.

В случае обучения на `all_speakers`, помимо тестирования строится матрица путаницы (`confusion matrix`) для каждого из четырёх пороговых значений: 0.5, 0.6, 0.7, 0.8. Анализ матриц путаницы лежит за рамками данной работы. Используя его, можно подобрать оптимальное пороговое значение для более удобного использования интерфейса распознавания речевых команд уже при его внедрении в программное обеспечение медиаплеера.

Смысл порогового значения следующий: если максимальный элемент в выходном тензоре нейронной сети ниже порогового значения, то команда распознаётся как тишина. В этом случае интерфейс медиаплеера просит повторить команду ещё раз. Если максимальный элемент в выходном тензоре нейронной сети выше или равен пороговому значению, то команда распознаётся как соответствующая индексу этого максимального элемента.

3.2 Результаты экспериментов

Графики обучения многослойного персептрона приведены на рисунках 7, 8, 9.

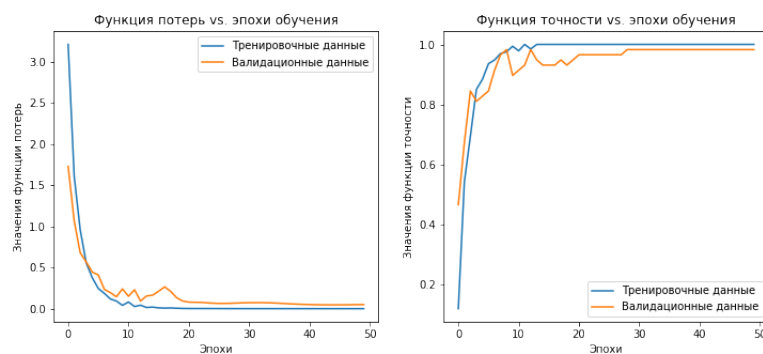


Рис. 7: Графики функций потерь и точности многослойного персептрона в течение обучения на speaker1

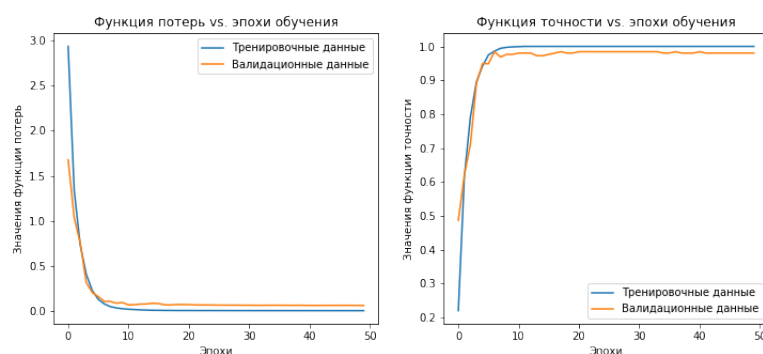


Рис. 8: Графики функций потерь и точности многослойного персептрона в течение обучения на all_male_speakers

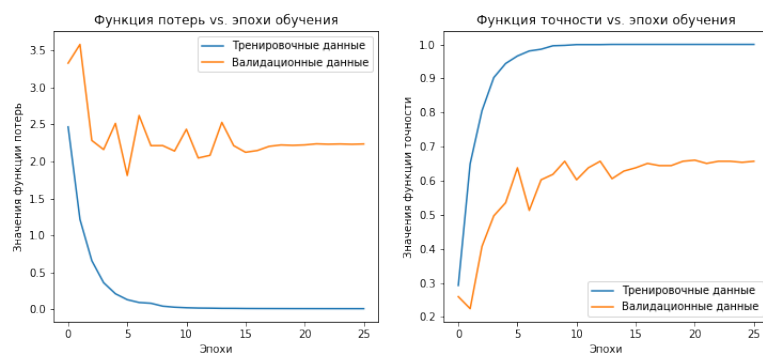


Рис. 9: Графики функций потерь и точности многослойного персептрона в течение обучения на all_speakers

Графики обучения свёрточной сети приведены на рисунках 10, 11, 12.

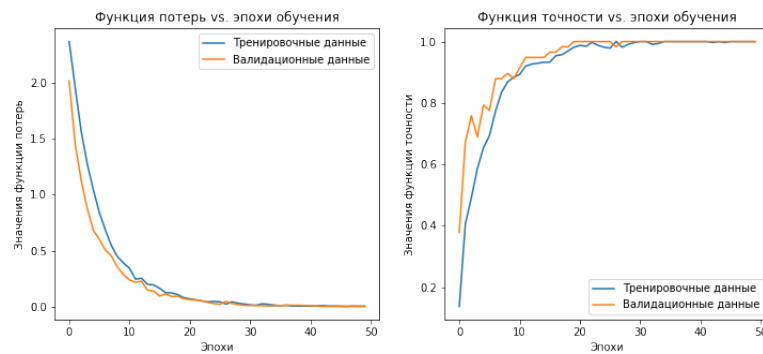


Рис. 10: Графики функций потерь и точности свёрточной сети в течение обучения на speaker1

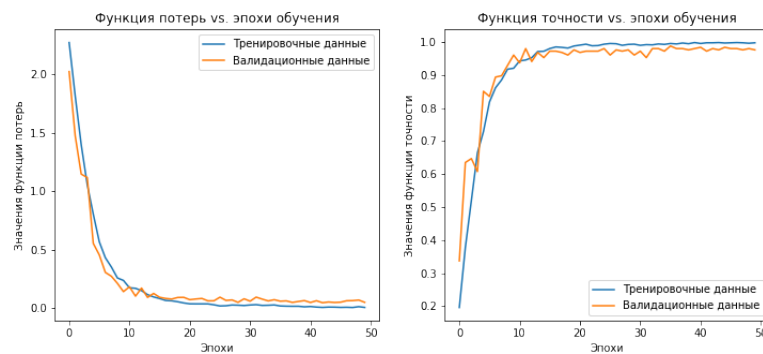


Рис. 11: Графики функций потерь и точности свёрточной сети в течение обучения на all_male_speakers

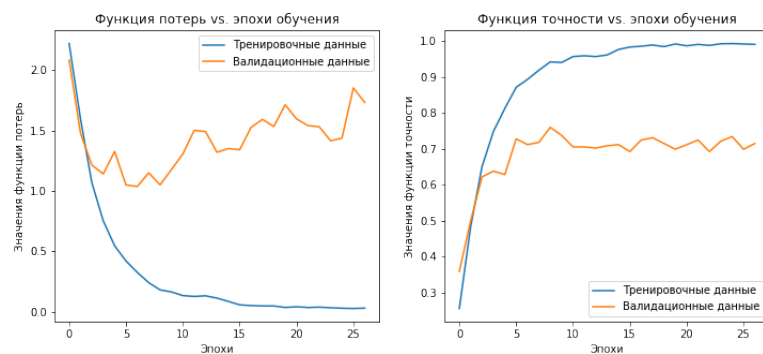


Рис. 12: Графики функций потерь и точности свёрточной сети в течение обучения на all_speakers

Матрицы путаницы для многослойного персептрона при обучении на all_speakers представлены на рисунке 13.

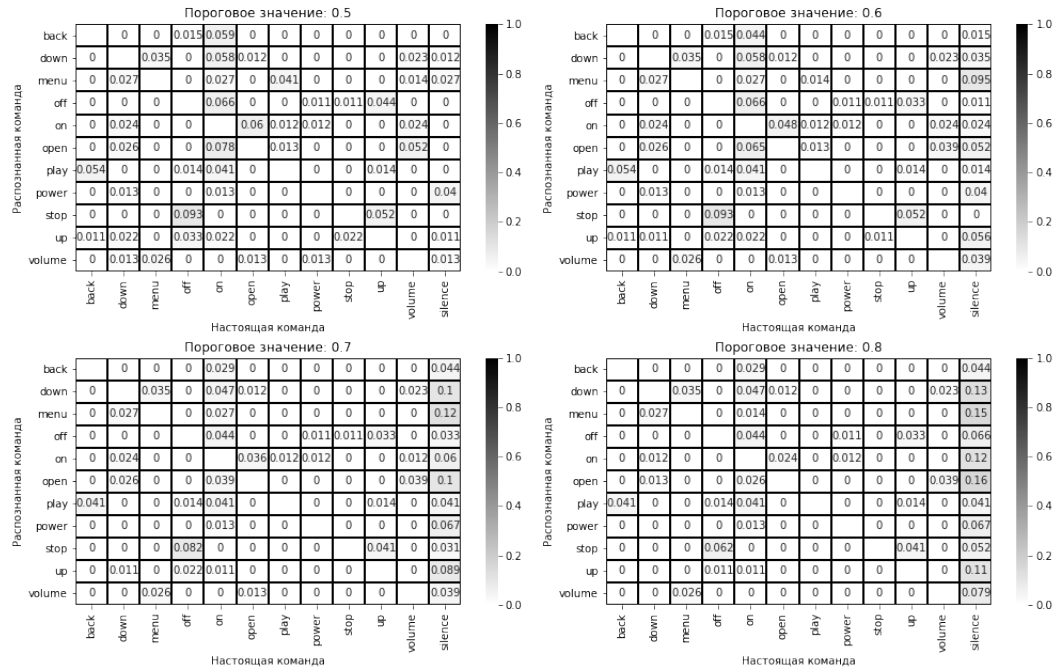


Рис. 13: Матрицы путаниц для многослойного персептрона при обучении на all_speakers

Матрицы путаницы для свёрточной сети при обучении на all_speakers представлены на рисунке 14.

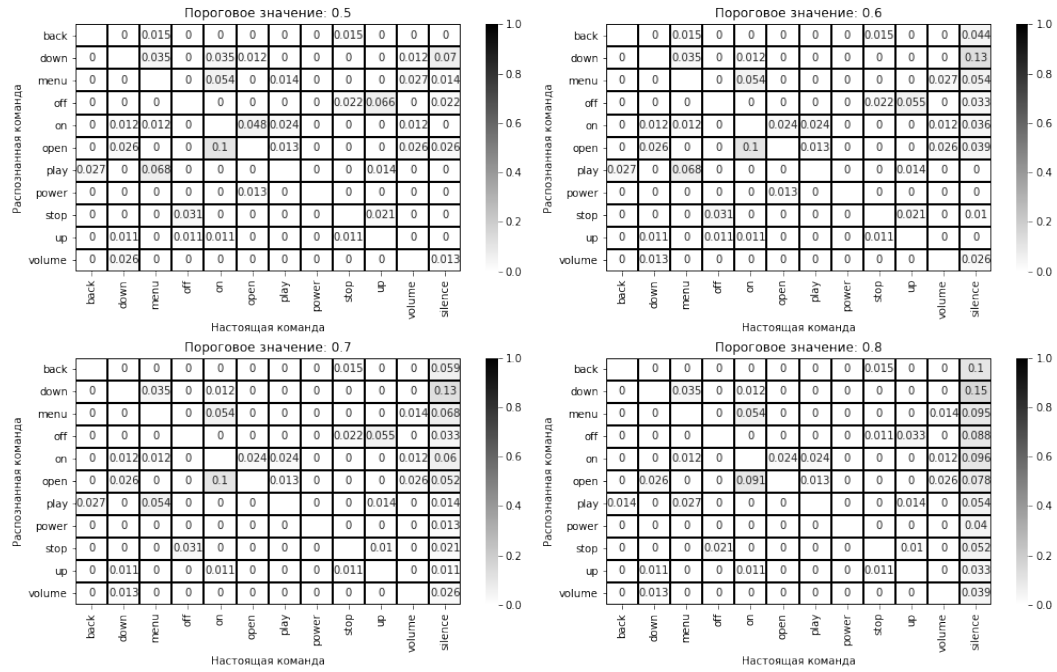


Рис. 14: Матрицы путаницы для свёрточной сети при обучении на all_speakers

Результаты тестирования представлены в таблице 3. Используемые обозначения: train_data - тренировочные данные, test_data - тестовые данные, cnn_loss - значения функции потерь для свёрточной сети, mlp_loss - значения функции потерь для многослойного персептрона, cnn_accuracy - значения функции точности для свёрточной сети, mlp_accuracy - значения функции точности для многослойного персептрона.

train_data	test_data	cnn_loss	mlp_loss	cnn_accuracy	mlp_accuracy
speaker1	speaker1	0.131	0.142	0.97	0.958
speaker1	speaker2	11.885	9.029	0.106	0.106
speaker1	speaker3	4.642	2.647	0.439	0.47
speaker1	speaker4	10.255	6.913	0.182	0.258
speaker1	speaker5	3.214	3.641	0.6	0.558
speaker1	speaker6	13.157	8.106	0.164	0.139
speaker1	all_speakers	7.023	4.956	0.429	0.43
all_male_speakers	speaker1	0.075	0.161	0.976	0.958
all_male_speakers	speaker2	0.335	0.347	0.924	0.894
all_male_speakers	speaker3	0.1	0.102	0.977	0.985
all_male_speakers	speaker4	0.231	0.433	0.932	0.909
all_male_speakers	speaker5	0.049	0.12	0.97	0.976
all_male_speakers	speaker6	11.315	11.646	0.273	0.261
all_male_speakers	all_speakers	2.217	2.339	0.831	0.819
all_speakers	speaker1	0.125	0.157	0.958	0.964
all_speakers	speaker2	0.204	0.338	0.932	0.894
all_speakers	speaker3	0.097	0.113	0.977	0.962
all_speakers	speaker4	0.257	0.418	0.932	0.864
all_speakers	speaker5	0.067	0.039	0.994	0.994
all_speakers	speaker6	1.569	2.247	0.715	0.624
all_speakers	all_speakers	0.409	0.581	0.915	0.881

Таблица 3: Результаты тестирования

Выводы

Видно, что при обучении только на одном дикторе с мужским голосом, точность распознавания на всех остальных дикторах по отдельности низкая.

При обучении на всех дикторах с мужским голосом, распознавание на дикторе с женским голосом работает лучше, чем при обучении на одном дикторе с мужским голосом, но точность все равно низкая.

При обучении на всех дикторах, распознавание на каждом даёт приемлемую точность. Однако стоит отметить, что если большая часть дикторов в тренировочной части имеет мужские голоса, то на дикторе с женским голосом распознавание работает хуже, чем на дикторах с мужским голосом.

Можно сделать вывод о том, что свёрточная нейронная сеть немного лучше справляется с поставленной задачей, чем многослойный персептрон. При этом количество весов, влияющее на время распознавания и на объем занимаемой памяти для их хранения, у свёрточной нейронной сети намного меньше из-за особенности её архитектуры.

Предложения по улучшению предложенной модели распознавания речи:

- увеличить размер датасета для обучения;
- повысить глубину предложенных архитектур нейронных сетей;
- увеличить количество коэффициентов MFCC в алгоритме выделения речевых признаков до 13;
- изменить количество мел-фильтров в алгоритме выделения речевых признаков.

В результате проведённой работы можно сделать вывод о возможности применения исследованного подхода к решению задачи распознавания речи для его использования при голосовом управлении медиаплеером.

Заключение

В данной работе:

- изучены математические методы и модели для решения поставленной задачи распознавания речи;
- произведён сбор исходных данных для обучения, тестирования и анализа эффективности алгоритма распознавания речи;
- разработан и реализован алгоритм предобработки исходных данных;
- разработан и реализован алгоритм распознавания речевых команд;
- проведены вычислительные эксперименты, в результате которых исследована работоспособность и эффективность алгоритма распознавания речевых команд;
- получен опыт реализации математических моделей программными методами;
- сделаны выводы об эффективности исследуемого подхода к решению задачи распознавания речи для прикладного применения.

Список литературы

- [1] Davis K. N., Biddulph R., Balashek S. Automatic recognition of spoken digits // The Journal of the Acoustical Society of America, 1952. Vol. 24, No 6. P. 637-642
- [2] Rosenblatt, F. The perceptron: A probabilistic model for information storage and organization in the brain // Psychological Review, 1958, Vol. 65(6), P. 386–408
- [3] Rosenblatt, F. Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms. Spartan Books, Washington DC, 1961
- [4] Bogert B. P., Healy M. J. R., Tukey J. W. The Quefrency Alanysis of Time Series for Echoes: Cepstrum, Pseudo Autocovariance, Cross-Cepstrum and Saphe Cracking // Proceedings of the Symposium on Time Series Analysis, 1963, Ch. 15, P. 209-243
- [5] Plomp R., Pols L. C. W., van der Geer J.P. Dimensional analysis of vowel spectra // The Journal of the Acoustical Society of America, 1967, Vol. 41, P. 707-712
- [6] Rabiner L.R., Sambur, M.R. An Algorithm for Determining the Endpoints of Isolated Utterances // The Bell System Technical Journal, 1975, Vol. 54, No. 2, P. 297-315
- [7] Newell A. Harpy, production systems and human cognition // Research Showcase @ Carnegie Mellon University, 1978
- [8] Оппенгейм А. В., Шафер Р. В. Цифровая обработка сигналов / Под ред. С. Я. Шаца, М.: Связь, 1979. С. 416
- [9] Davis S., Mermelstein P. Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences // IEEE

Transactions on Acoustics, Speech, and Signal Processing, 1980, Vol. ASSP-28, No. 4, P. 357-366

- [10] O'Shaughnessy D. Speech communication: human and machine, Addison-Wesley, 1987, P. 568
- [11] LeCun Y., Bengio Y.. Convolutional networks for images, speech, and time-series // The Handbook of Brain Theory and Neural Networks, 1995, MIT

Приложение

Ссылка на репозиторий, содержащий программный код веб-приложения для записи датасета: <https://gitlab.com/polotent/commandrecorder>

Ссылка на репозиторий, содержащий программный код предобработки данных, обучения и тестирования нейронной сети: <https://gitlab.com/polotent/boxy>