



Министерство науки и высшего образования 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 «Информатика и системы управления»

КАФЕДРА «Программное обеспечение ЭВМ и информационные технологии»

**РАСЧЕТНО-ПОЯСНИТЕЛЬНАЯ ЗАПИСКА  
К ВЫПУСКНОЙ КВАЛИФИКАЦИОННОЙ РАБОТЕ  
НА ТЕМУ:**

Метод выделения звуков естественного языка в звучащей речи.

Студент группы **ИУ7-82Б**

\_\_\_\_\_  
(Подпись, дата)

**Левушкин И.К.**

\_\_\_\_\_  
(И.О. Фамилия)

Руководитель ВКР

\_\_\_\_\_  
(Подпись, дата)

**Градов В.М.**

\_\_\_\_\_  
(И.О. Фамилия)

Консультант

\_\_\_\_\_  
(Подпись, дата)

**Строганов Ю.В.**

\_\_\_\_\_  
(И.О. Фамилия)

Нормоконтроллер

\_\_\_\_\_  
(Подпись, дата)

**Строганов Ю.В.**

\_\_\_\_\_  
(И.О. Фамилия)

**2021 г.**

## РЕФЕРАТ

Расчетно-пояснительная записка 64 с., 22 рис., 6 табл., 27 ист., 1 прил.  
РАСПОЗНАВАНИЕ РЕЧИ, СЕГМЕНТАЦИЯ РЕЧЕВОГО СИГНАЛА,  
ВЕЙВЛЕТ-ПРЕОБРАЗОВАНИЕ, МЕТОД ОПОРНЫХ ВЕКТОРОВ

Объектом исследования являются системы распознавания речи. Предметом исследования - методы автоматической контекстно-независимой временной сегментации речевого сигнала. Цель работы - разработка метода выделения звуков естественного языка в звучащей речи. Задачи, решаемые в работе:

- анализ предметной области - описание основных этапов распознавания речевого сигнала, обоснование актуальности задачи сегментации речевого сигнала, проведение анализа существующих решений сегментации речевого сигнала;
- проведение анализа существующих методов выделения признаков речевого сигнала и алгоритмов классификации;
- выбор методов наиболее оптимальных для поставленной цели;
- проектирование метода выделения звуков естественного языка в звучащей речи, разработка алгоритма, реализующего данный метод;
- проектирование системы, проверяющей работоспособность метода;
- формирование выборки для проведения эксперимента;
- проведение исследования разработанного метода на данной выборке.

Область применения - системы распознавания речи. Предлагаемые направления развития:

- подготовка большего количества обучающей выборки, имеющей меньшую погрешность измерений для повышения точности метода;
- применение альтернативных алгоритмов классификации, позволяющих повысить качество распознавания с применением большего количества данных.

Поставленная цель была достигнута: метод выделения звуков естественного языка в звучащей речи был разработан, реализован и протестирован. Были рассмотрены существующие недостатки метода и предложены пути дальнейшего развития.

# СОДЕРЖАНИЕ

<b>ВВЕДЕНИЕ</b>	<b>5</b>
<b>1 Аналитическая часть</b>	<b>7</b>
1.1 Анализ предметной области . . . . .	7
1.1.1 Структура систем распознавания речи . . . . .	7
1.1.2 Сегментация речевого сигнала . . . . .	8
1.2 Методы выделения признаков речевого сигнала . . . . .	9
1.2.1 Преобразование Фурье . . . . .	9
1.2.2 Оконное преобразование Фурье . . . . .	10
1.2.3 Вейвлет-преобразование . . . . .	11
1.2.4 Преобразование Гильберта-Хуанга . . . . .	15
1.3 Алгоритмы классификации . . . . .	20
1.3.1 Формальная постановка задачи классификации . . . . .	21
1.3.2 Наивный байесовский классификатор . . . . .	22
1.3.3 Метод k-ближайших соседей . . . . .	23
1.3.4 Метод деревьев решений (Decision Trees, DT) . . . . .	24
1.3.5 Метод опорных векторов . . . . .	25
1.3.6 Методы на основе искусственных нейронных сетей . . . . .	26
1.4 Постановка задачи . . . . .	28
1.5 Выводы . . . . .	29
<b>2 Конструкторская часть</b>	<b>32</b>
2.1 Метод выделения звуков естественного языка в звучащей речи . . . . .	32
2.2 Описание алгоритма метода выделения звуков естественного языка в звучащей речи . . . . .	32
2.3 Выделение признаков речевого сигнала . . . . .	33
2.4 Формирование данных для обучения модели классификации . . . . .	34
2.5 Обучение модели классификации . . . . .	35
2.6 Структура ПО . . . . .	35
2.7 Описание формата входных и выходных данных . . . . .	38
2.8 Выводы . . . . .	39
<b>3 Технологическая часть</b>	<b>40</b>
3.1 Выбор языка программирования . . . . .	40
3.2 Выбор используемых библиотек . . . . .	40
3.3 Выбор среды разработки . . . . .	41
3.4 Пользовательский интерфейс . . . . .	42
3.5 Реализация . . . . .	44

3.6	Выводы . . . . .	52
<b>4</b>	<b>Исследовательская часть</b>	<b>53</b>
4.1	Формирование обучающей и тестовой выборки данных . . . . .	53
4.2	Метрики оценки качества моделей классификации . . . . .	53
4.2.1	precision и recall . . . . .	54
4.2.2	Коэффициент корреляции Мэтьюса . . . . .	55
4.2.3	AUC-ROC . . . . .	55
4.3	Поиск оптимальных гиперпараметров метода . . . . .	56
4.3.1	GridSearchCV с применением ККМ . . . . .	56
4.4	Оценка качества модели классификации . . . . .	58
4.5	Выводы . . . . .	60
	<b>ЗАКЛЮЧЕНИЕ</b>	<b>61</b>
	<b>СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ</b>	<b>62</b>
	<b>ПРИЛОЖЕНИЕ А</b>	<b>65</b>

## ВВЕДЕНИЕ

Распознавание речи — автоматический процесс преобразования речевого сигнала в цифровую информацию, является одной из актуальных и сложных задач настоящего времени, имеющая множество успешных применений в различных сферах бизнеса таких как:

- Телефония: автоматизация обработки входящих и исходящих звонков [2];
- Бытовая техника: голосовое управление бытовой техникой [2];
- Десктопы, ноутбуки и мобильные устройства: голосовой ввод в играх и приложениях [2];
- Автомобили: голосовое управление в салоне автомобиля [2];
- Социальные сервисы для людей с ограниченными возможностями [3].

Здесь задействованы достижения весьма различных областей: от компьютерной лингвистики до цифровой обработки сигналов.

Несмотря на значительный прогресс в точности и скорости современных систем, их результаты все еще заметно далеки от желаемых, поскольку человеческая речь представляет собой сложный для распознавания нестационарный нелинейный сигнал.

В системах распознавания речи первоочередной задачей, требующей решения, является задача автоматической сегментации в соответствии с фонетической транскрипцией языка. Под сегментацией понимается операция выделения звуков естественного языка (основных лингвистических элементов языка [1]) в звучащей речи.

**Целью** работы является разработка метода выделения звуков естественного языка в звучащей речи. Для достижения поставленной цели необходимо решить следующие **задачи**:

- провести анализ предметной области:
  - описать основные этапы распознавания речевого сигнала;
  - обосновать актуальность задачи сегментации речевого сигнала;
  - провести анализ существующих решений сегментации речевого сигнала.
- провести анализ существующих методов выделения признаков речевого сигнала и алгоритмов классификации;

- выбрать методы наиболее оптимальные для поставленной цели;
- на их основе спроектировать метод выделения звуков естественного языка в звучащей речи и разработать алгоритм, реализующий данный метод;
- спроектировать систему для проверки работоспособности метода;
- сформировать выборку для проведения эксперимента;
- провести исследование разработанного метода на данной выборке;

## 1 Аналитическая часть

В данном разделе проведен анализ предметной области - описаны основные этапы распознавания речевого сигнала, дано обоснование актуальности задачи сегментации речевого сигнала, проведен анализ существующих решений сегментации речевого сигнала.

Проанализированы существующие методы выделения признаков речевого сигнала и алгоритмов классификации.

**Объектом исследования** являются системы распознавания речи. **Предметом исследования** - методы автоматической контекстно-независимой временной сегментации речевого сигнала.

### 1.1 Анализ предметной области

**Распознавание речи** — процесс преобразования речевого сигнала в цифровую информацию.

Задачей распознавания является сопоставление набору акустических признаков речевого сигнала или наблюдений  $X(x_1, \dots, x_n)$  последовательности слов  $W(w_1, \dots, w_k)$ , имеющих наибольшую вероятность правдоподобия среди всех кандидатов. Для этого используется формула Байеса:

$$W = \operatorname{argmax} \left[ \frac{P(W) * P(X|W)}{P(X)} \right] \quad (1)$$

Причем, в процессе распознавания вероятность уже полученных признаков  $P(X)$  не подлежит оптимизации и знаменатель в формуле не используется:

$$W = \operatorname{argmax} [P(W) * P(X|W)] \quad (2)$$

#### 1.1.1 Структура систем распознавания речи

Системы распознавания речи впервые появились в 1952 году. С тех пор методы распознавания не раз менялись. Ранее использовались такие методы и алгоритмы, как:

- Динамическое программирование (Dynamic Time Warping) - временные динамические алгоритмы, выполняющие классификацию на основе сравнения с эталоном;
- Методы дискриминантного анализа, основанные на Байесовской дискриминации (Bayesian discrimination);

- Скрытые Марковские Модели (Hidden Markov Model).
- Нейронные сети (Neural Networks).

В настоящее время, перечисленные выше методы как правило комбинируются. Их сочетание позволяет получить более высокое качество распознавания, чем использование каждой модели отдельно.

Системы распознавания речи имеют следующие основные модули [4], [5]:

- Акустическая модель - функция, принимающая на вход признаки на небольшом участке акустического сигнала (фрейме) и выдающая распределение вероятностей различных фонем на этом фрейме. Таким образом, акустическая модель дает возможность по звуку восстановить, что было произнесено — с той или иной степенью уверенности.;
- Языковая модель - позволяет узнать, какие последовательности слов в языке более вероятны, а какие менее. Здесь в самом простом случае требуется предсказать следующее слово по известным предыдущим словам;
- Декодер - компонент, который сопоставляет входной речевой поток с информацией, хранящейся в акустических и языковых моделях, и определяет наиболее вероятную последовательность слов, которая и является конечным результатом распознавания.

### **1.1.2 Сегментация речевого сигнала**

Для эффективной работы системы акустическая модель должна получать на вход признаки, которые бы наилучшим образом описывали особенности данного сигнала. Но даже при наличии возможности использовать наиболее информативные методы выделения признаков речевого сигнала, результат распознавания зачастую далек от совершенства из-за сложной структуры речевого сигнала, недостаточной изученности механизмов речеобразования и речевосприятия.

Чтобы повысить точность акустической модели, помимо основных признаков необходимо также подавать на вход «метапризнаки» сигнала.

Одними из таких признаков могут служить границы между основными лингвистическими элементами языка. Другими словами, возникает потребность в решении задачи временной сегментации речевого сигнала.

Рассмотрев существующие решения сегментации речевого сигнала ([7],



[8], [9]) были выделены два основных этапа в задаче сегментации:

- **выделение признаков сигнала;**
- **определение межфонемных переходов на основе выделенных признаков сигнала.**

Существует два основных типа алгоритмов сегментации речи.

- **Контекстно-зависимая сегментация** - к этому типу относятся алгоритмы, которые производят сегментацию речи при условии, что известна последовательность фонем данной фразы;
- **Контекстно-независимая сегментация** - другой тип алгоритмов, не использующий априорной информации о фразе, при этом границы сегментов определяются по степени изменения акустических характеристик сигнала.

При автоматической сегментации желательно использовать только общие характеристики речевого сигнала, поскольку обычно на этом этапе нет конкретной информации о содержании речевого высказывания, поэтому в данной работе **предметом исследования** являются методы автоматической контекстно-независимой временной сегментации речевого сигнала.

## **1.2 Методы выделения признаков речевого сигнала**

Сегментация речевого сигнала подразумевает выделение участков сигнала, соответствующим отдельным структурным единицам. Для решения данной задачи необходимо выбрать метод выделения акустических характеристик сигнала на заданных интервалах времени. На основе этих признаков можно будет оценить степень изменения сигнала во времени и таким образом определить границы между основными лингвистическими элементами языка.

Ниже приведен анализ наиболее распространенных методов выделения признаков речевого сигнала.

### **1.2.1 Преобразование Фурье**

Преобразование Фурье — это функция, описывающая амплитуду и фазу каждой синусоиды, соответствующей определённой частоте, где амплитуда представляет собой высоту кривой, а фаза - начальную точку синусоиды [10]. Эта новая функция описывает коэффициенты при разложении исходной функции на элементарные составляющие — гармонические колебания с разными частотами (подобно тому, как музыкальный аккорд может быть выражен в виде амплитуд нот, которые его составляют). Преобразование Фурье функции  $f$

вещественной переменной является интегральным и задаётся с помощью следующей формулы:

$$F(\omega) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f(t) e^{-it\omega} dt, \quad (3)$$

где

- $f(t)$  - исходный сигнал;
- $F(t)$  - Фурье-образ функции  $f(t)$  или частотный спектр функции  $f(t)$ ;
- $\omega$  - круговая частота.

### 1.2.2 Оконное преобразование Фурье

Классическое преобразование Фурье имеет дело со спектром сигнала, взятым во всем диапазоне существования переменной. Нередко интерес представляет только локальное распределение частот, в то время как требуется сохранить изначальную переменную (обычно время).

С позиций точного представления произвольных сигналов и функций, преобразование Фурье имеет ряд недостатков, которые привели к появлению оконного преобразования Фурье и стимулировали развитие вейвлет-преобразования [11]:

- ограниченная информативность анализа нестационарных сигналов и практически полное отсутствие возможностей анализа их особенностей (сингулярностей), так как в частотной области происходит «размазывание» особенностей сигналов (разрывов, ступеней, пиков и т. п.) по всему частотному диапазону спектра;
- появление эффекта Гиббса на скачках функций, при усечении сигналов и при вырезке отрезков сигналов для локального детального анализа;
- гармонический характер базисных функций, определенных в интервале от  $-\infty$  до  $+\infty$ .

Неспособность преобразования Фурье осуществлять временную локализацию сингулярностей сигналов может быть частично устранена введением в преобразование так называемой движущейся оконной функции, имеющей компактный носитель. Использование оконной функции позволяет представлять результат преобразования в виде функции двух переменных — частоты и временного положения окна.

Оконное преобразование Фурье имеет следующий вид [12]:

$$F(t, \omega) = \int_{-\infty}^{\infty} f(\tau)W(\tau - t)e^{-i\omega\tau} d\tau, \quad (4)$$

где

- $F(t, \omega)$  даёт распределение частот части оригинального сигнала  $f(t)$  в окрестности момента времени  $t$ ;
- $W(\tau - t)$  - некоторая оконная функция.

### 1.2.3 Вейвлет-преобразование

Широко используемое преобразование Фурье для анализа сигналов, как непрерывных, так и дискретных, оказывается недостаточно эффективным при обработке сложных сигналов.

Например, Фурье спектры для сигналов из двух синусоид, которые с разными частотами, первый из которых, представляющий собой сумму синусоид, а второй представляет собой, последовательно следующие друг за другом синусоиды, одинаковы и будут выглядеть как два пика на двух фиксированных частотах (рисунок 1). Из этого следует, преобразование Фурье в своём обычном виде не приспособлено для анализа нестационарных сигналов, так как теряется информация о временных характеристиках сигнала.

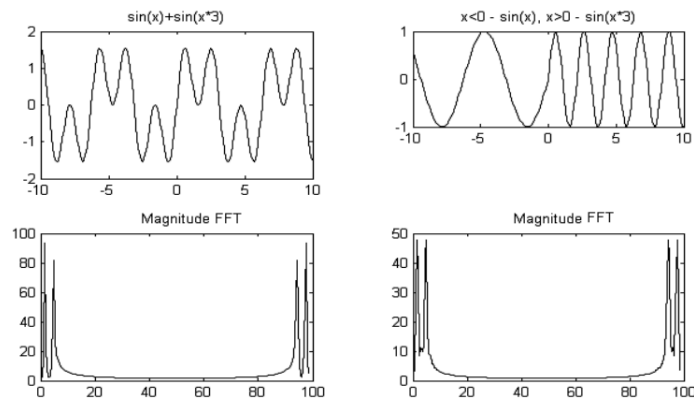


Рисунок 1 – Пример неинформативности преобразования Фурье

Речевой сигнал является примером нестационарного процесса, в котором информативным является сам факт изменения его частотно-временных характеристик.

Для анализа таких процессов требуются базисные функции, способные выявлять в исследуемом сигнале как частотные, так и его временные характеристики, то есть функции со свойствами частотно-временной локализации. Такие возможности предоставляют вейвлеты, являющиеся обобщением спектрального анализа [10].

Ниже приведена общая формула вейвлет-преобразования функции  $f(t)$ .

$$W_{\Psi}(x, a) = \frac{1}{\sqrt{a}} \int_{-\infty}^{\infty} f(t) \Psi\left(\frac{t-x}{a}\right) dt, \quad (5)$$

где

- $\Psi\left(\frac{t-x}{a}\right)$  - вейвлет;
- $a$  - масштабный коэффициент;
- $x$  - параметры сдвига.

Вейвлеты – функции двух аргументов – масштаба и сдвига. В отличие от стандартного преобразования Фурье, они позволяют обрабатывать сигнал одновременно в физическом – время, координата; и частотном пространствах.

Выбор анализирующего вейвлета определяется тем, какую информацию необходимо извлечь из сигнала. Каждый вейвлет имеет характерные особенности во временном и в частотном пространстве, поэтому иногда с помощью разных вейвлетов можно полнее выявить и подчеркнуть те или иные свойства анализируемого сигнала.

Таким образом, вейвлет-преобразование обеспечивает двумерное представление исследуемого сигнала в частотной области в плоскости частота-положение. Аналогом частоты при этом является масштаб аргумента базисной функции (чаще всего – времени), а положение характеризуется её сдвигом. Это позволяет найти особенности сигналов, одновременно локализуя их на временной шкале. Другими словами, вейвлет-анализ можно охарактеризовать как спектральный анализ локальных возмущений.

Используя вейвлет-преобразование, сигнал можно представить как последовательность образов с разной степенью детализации, что позволяет найти локальные особенности сигнала и классифицировать их по интенсивности.

**Многомасштабный Вейвлет-анализ** основывается на разложении сигнала по функциям, образующим ортонормированный базис. Каждую функцию можно разложить на некотором заданном уровне разрешения (масштабе)  $j_n$  в

ряд вида [10]:

$$f(x) = \sum_{k=0}^{2M-1} s_{j_n,k} \varphi_{j_n,k} + \sum_{j \geq j_n} \sum_{k=0}^{2M-1} d_{j,k} \Psi_{j,k}, \quad (6)$$

где

- $\varphi_{j_n,k}$  и  $\Psi_{j,k}$  - масштабированные и смещенные версии скейлинг-функции (масштабной функции)  $\varphi$  и материнского вейвлета  $\Psi$ ;
- $s_{j,k}$  - коэффициенты аппроксимации;
- $d_{j,k}$  - детализирующие коэффициенты.

На рисунках 2 - 5 представлены аналитические графики функций  $\varphi$  и  $\Psi$  указанных вейвлетов.

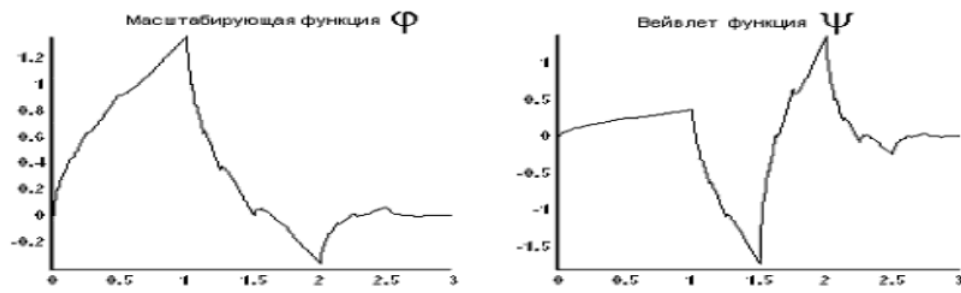


Рисунок 2 – Аналитические графики функций  $\varphi$  и  $\Psi$ , вейвлет Симлета

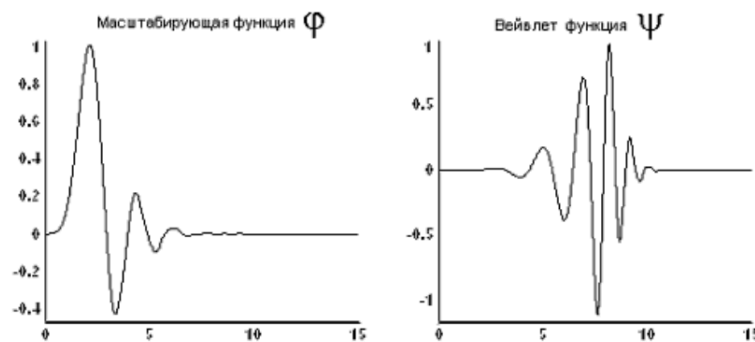


Рисунок 3 – Аналитические графики функций  $\varphi$  и  $\Psi$ , вейвлет Добеши 8

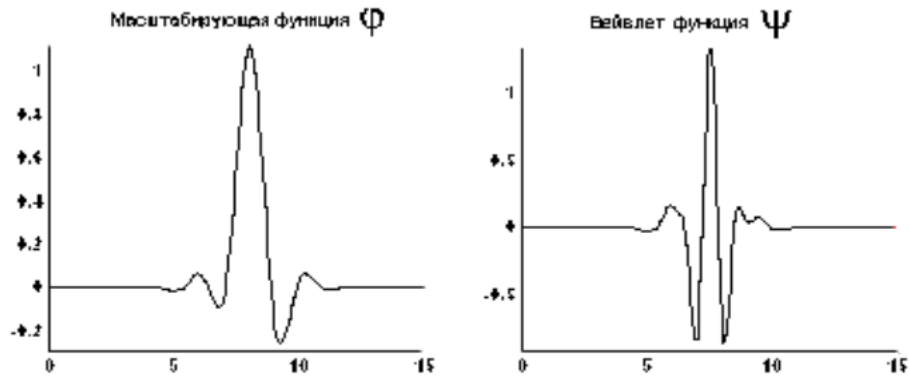


Рисунок 4 – Аналитические графики функций  $\varphi$  и  $\Psi$ , вейвлет Симлета 8

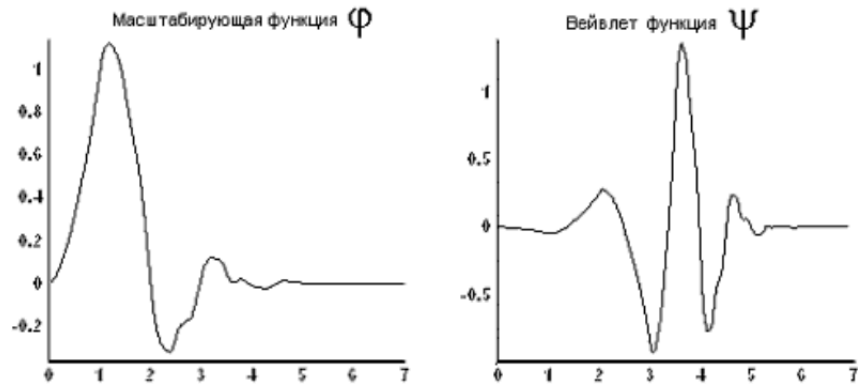


Рисунок 5 – Аналитические графики функций  $\varphi$  и  $\Psi$ , вейвлет Добеши 4

Масштабирование и смещение функций  $\varphi_{j,k}$  и  $\Psi_{j,k}$  находится следующим образом:

$$\varphi_{j,k} = 2^{j/2} \varphi(2^j x - k), \quad (7)$$

$$\Psi_{j,k} = 2^{j/2} \Psi(2^j x - k), \quad (8)$$

В свою очередь сами функции  $\varphi$  и  $\Psi$  определяются следующим образом:

$$\varphi = \sqrt{2} \sum_{k=0}^{2M-1} h_k \varphi(2x - k), \quad (9)$$

$$\Psi = \sqrt{2} \sum_{k=0}^{2M-1} g_k \Psi(2x - k), \quad (10)$$

где  $g_k = (-1)^k h_{2M-k-1}$

Таким образом, ортогональный вейвлет-анализ сводится к нахождению коэффициентов аппроксимации и детализирующих коэффициентов в разложении сигнала  $f(t)$  по формуле (5).

#### **Достоинства вейвлет-преобразования:**

- вейвлетные преобразования обладают всеми достоинствами преобразований Фурье;
- вейвлетные базисы могут быть хорошо локализованными как по частоте, так и по времени. При выделении в сигналах хорошо локализованных разномасштабных процессов можно рассматривать только те масштабные уровни разложения, которые представляют интерес;
- базисные вейвлеты могут реализоваться функциями различной гладкости.

#### **Недостатки:**

- относительная сложность преобразования.

#### **1.2.4 Преобразование Гильберта-Хуанга**

Под преобразованием Гильберта-Хуанга (Hilbert-Huang transform – ННТ) понимается метод эмпирической модовой декомпозиции (EMD) нелинейных и нестационарных процессов и Гильбертов спектральный анализ (HSA) [10]. Этот метод потенциально жизнеспособен для нелинейного и нестационарного анализа данных [9], специально для частотно-энергетических временных представлений.

EMD-HSA предложил Норден Хуанг в 1995 в США (NASA) для изучения поверхностных волн тайфунов, включая возможность на анализ произвольных временных рядов коллективом соавторов в 1998 г.. В последующие годы, активно расширяя применения алгоритма для других новых отраслей науки и техники, взамен термина EMD-HSA был принят более короткий термин преобразования ННТ.

**EMD (Empirical Mode Decomposition)** – метод разложения сигналов на функции, получившие названия внутренних или «эмпирических мод». Метод

представляет собой адаптивную итерационную вычислительную процедуру разложения исходных данных (непрерывных или дискретных сигналов) на эмпирические моды или внутренние колебания.

**Огибающие сигналов.** У каждого сигнала присутствуют локальные экстремумы: чередующиеся локальные максимумы и локальные минимумы с произвольным расположением по координатам (независимым переменным) сигналов. По этим экстремумам с использованием методов аппроксимации можно построить две огибающие сигналов: нижнюю – построенную по точкам локальных минимумов, и верхнюю – построенную по точкам локальных максимумов, а также функцию «среднего значения огибающих», которой отвечает срединная линия, расположенная в точности между нижней и верхней огибающими [10].

Модовая декомпозиция сигналов основана на предположении, что любые данные состоят из различных внутренних колебаний IMF. В любой момент времени данные могут иметь множество сосуществующих внутренних колебаний – IMFs. Каждое колебание, линейное или нелинейное, представляет собой модовую функцию, которая имеет экстремумы и нулевые пересечения. Кроме того, колебания в определенной степени «симметричны» относительно локального среднего значения. Конечные сложные данные образуются суммой модовых функций, наложенных на региональный тренд сигнала.

Эмпирическая мода – это такая функция, которая обладает следующими свойствами:

- 1) количество экстремумов функции (максимумов и минимумов) и количество пересечений нуля не должны отличаться более чем на единицу;
- 2) в любой точке функции среднее значение огибающих, определенных локальными максимумами и локальными минимумами, должно быть нулевым.

IMF представляет собой колебательный режим, но вместо постоянной амплитуды и частоты, как в простой гармонике, у IMF могут быть переменная амплитуда и частота, как функции независимой переменной. Первое свойство гарантирует, что локальные максимумы функции всегда положительны, локальные минимумы соответственно отрицательны, а между ними всегда имеют место пересечения нулевой линии. Второе свойство гарантирует, что мгновенные



частоты функции не будут иметь нежелательных флуктуаций, являющихся результатом асимметричной формы волны.

Любую функцию и любой произвольный сигнал, изначально содержащие произвольную последовательность локальных экстремумов (минимум 2), можно разделить на семейство функций IMFs и остаточный тренд. Если данные лишены экстремумов, но содержат точки перегиба («скрытые» экстремумы наложения модовых функций и крутых трендов), то для открытия экстремумов может использоваться дифференцирование сигнала.

Схема преобразования Гильберта-Хуанга делится на две части. На первом этапе, экспериментальные данные разлагаются в ряд внутренних модовых функций (IMFs). Эта декомпозиция рассматривается как расширение данных в терминах внутренних модовых функций. Иначе, эти внутренние модовые функции представлены как базис преобразования, которое может быть линейным или нелинейным, как диктуется по условиям. Так как IMFs имеют хорошие Гильбертовы преобразования, то могут быть вычислены соответствующие мгновенные частоты. На следующем этапе локализуются любые явления, как во времени, так и на частотной оси. Локальная энергия и мгновенная частота, выведенная из IMFs, дают дистрибутивные «энергетические время-частотные» данные, и такое представление, определяемое как Гильбертов спектр:

$$r_j(t) = r_{j-1} - c_j(t) \quad (11)$$

Пусть имеется произвольный сигнал  $y(t)$ . Сущность метода EMD заключается в последовательном вычислении функций эмпирических мод  $c(t)$  и остатков  $r_j(t) = r_{j-1} - c_j(t)$ , где  $j = 1, 2, 3, \dots, n$  при  $r_0 = y(t)$ . Результатом разложения будет представление сигнала в виде суммы модовых функций и конечного остатка:

$$x(t) = \sum_{j=1}^n c_j(t) + r_n(t), \quad (12)$$

где  $n$  - количество эмпирических мод, которое устанавливается в ходе вычислений.

Алгоритм эмпирической декомпозиции сигнала складывается из следующих операций преобразования:

- 1) находится в сигнале  $y(k)$  положение всех локальных экстремумов, максимумов и минимумов процесса (номера точек  $k_{i.ext}$  экстремумов), и значения в этих точках (рисунок 6). Между этими экстремумами сосредоточена вся информация сигнала. Группируются массивы координат  $k_{i.ext}$  для максимумов и минимумов, и соответствующих им амплитудных значений  $y(k_{i.ext})$ . Число строк в массивах максимумов и минимумов не должно отличаться более чем на 1.

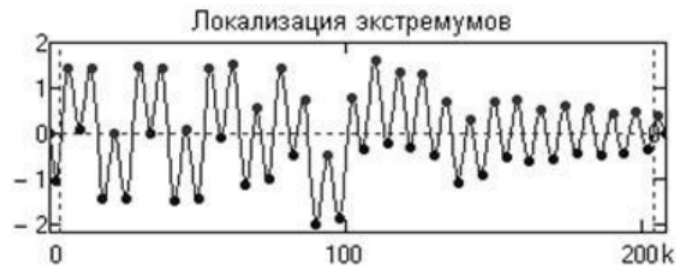


Рисунок 6 – Локализация экстремумов в сигнале

- 2) Применяя сплайны вычисляются верхняя  $u_t(k)$  и нижняя  $u_b(k)$  огибающие процесса соответственно, по максимумам и минимумам, как это показано на рисунке 7. Определяется функция средних значений  $m_1(k)$  между огибающими (рисунок 7).

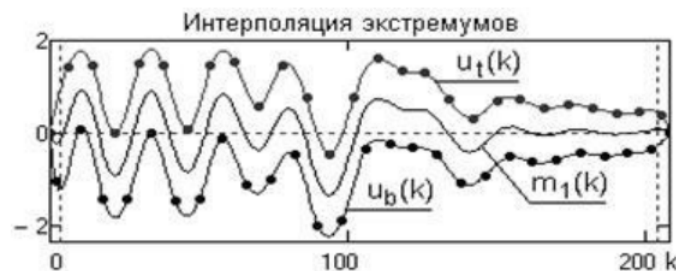


Рисунок 7 – Интерполяция экстремумов и построение огибающих

Далее, определяется функция средних значений  $m_1(k)$  между огибающими.

$$m_1(k) = \frac{u_t(k) + u_b(k)}{2}, \quad (13)$$

Разность между сигналом  $y(k)$  и функцией  $m_1(k)$  дает первую компоненту отсеивания - функцию  $h_1(k)$ , которая является первым приближением к первой функции IMF:

$$h_1(k) = y(k) - m_1(k), \quad (14)$$

- 3) Повторяются операции 1) и 2), принимая вместо  $y(k)$  функцию  $h_1(k)$ , и находят второе приближение к первой функции IMF – функцию  $h_2(k)$ .

$$h_2(k) = h_1(k) - m_2(k), \quad (15)$$

Последующие итерации выполняются аналогичным образом:

$$h_i(k) = h_{i-1}(k) - m_i(k), \quad (16)$$

По мере увеличения количества итераций функция  $m_i(k)$  стремится к нулевому значению, а функция  $h_i(k)$  – к неизменяемой форме.

Последнее значение  $h_i(k)$  итераций принимается за наиболее высокочастотную функцию  $r_1(k) = y(k) - c_1(k)$  семейства IMF, которая непосредственно входит в состав исходного сигнала  $y(k)$ . Это позволяет вычесть  $c_1(k)$  из состава сигнала и оставить в нем более низкочастотные составляющие  $r_1(k) = y(k) - c_1(k)$ ,

На рисунке 8 показано графическое представление вычитания из сигнала высокочастотной составляющей, сформированной по алгоритму, заданному (12)-(16).

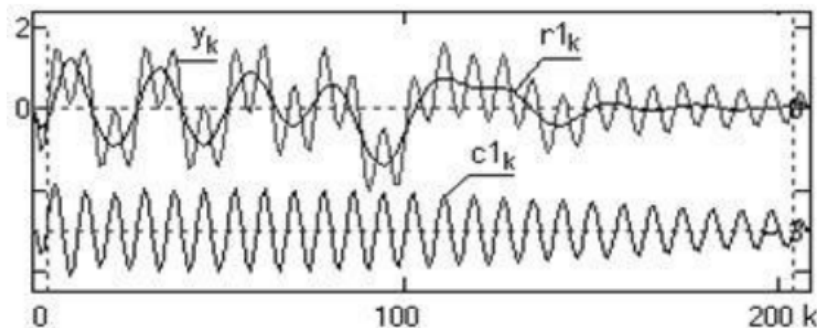


Рисунок 8 – Выявление низкочастотных составляющих в сигнале

Следующая внутренняя модовая функция находится, повторяя операции 1)-3) декомпозиции, описанные выше, с той разницей, что входным сигналом является остаток  $r_1(k)$ .

Шаги 1)-3) повторяются для всех последующих  $r_j(k)$ , и результат будет представлять последовательность вычислений:

$$r_1(k) - c_2(k) = r_2(k), \dots, r_{n-1}(k) - c_n(k) = r_n(k), \quad (17)$$

Метод EMD закончен, когда остаток, в идеале, не содержит экстремумов. Это означает, что остаток – или константа или монотонная функция. Извлечённые IMFs симметричны, имеют уникальные локальные частоты, различные IMFs не показывают ту же самую частоту в то же самое время. Другими словами, остановка декомпозиции сигнала должна происходить при максимальном «выпрямлении» остатка, то есть превращения его в тренд сигнала по интервалу задания с числом экстремумов не более 2-3 [10].

### **Достоинства преобразования Гильберта-Хуанга:**

- базис, используемый при разложении (набор эмпирических мод) конструируется непосредственно из того сигнала, с которым ведется работа, что позволяет учесть все его локальные особенности, внутреннюю структуру, наличие нежелательных особенностей (шумы, тренды, аномальные выбросы, пропущенные значения).

### **Недостатки:**

- гораздо более высокая вычислительная сложность по сравнению с вейвлет-преобразованием.

### **1.3 Алгоритмы классификации**

После получения акустических характеристик сигнала на заданных интервалах времени необходимо определить границы между основными лингвистическими элементами языка.

Большинство работ, связанных с автоматической сегментацией речевого сигнала, для решения данной задачи используют определенные закономерности, эвристики в изменении значений выделенных признаков речевого сигнала

в момент межфонемного перехода [13]. Эти закономерности чаще всего не имеют большой теоретической базы и лишь основываются на экспериментальном опыте других работ.

Так, в работах [7], [8] сегментация речевого сигнала на фрагменты, осуществляющая на базе Вейвлет-преобразования, используют тот факт, что на межфонемных переходах сигнал претерпевает значительные изменения сразу на многих масштабах исследования и, соответственно, характеризуется возрастанием вейвлет-коэффициентов для многих уровней детализации, в то время как на стационарных участках фонем вейвлет-коэффициенты оказываются сгруппированными вблизи определённых масштабов. Таким образом, отыскание межфонемных границ сводится к отысканию моментов увеличения вейвлет-коэффициентов на значительном количестве уровней масштабирования.

На основе этих закономерностей формируется критерий определения межфонемных переходов, который в свою очередь также определяется экспериментальным путем, в виду чего не гарантируют, полученной в результате экспериментов точности на других данных из-за сложной структуры звукового сигнала.

Исходя из приведенных выше рассуждений в данной работе решено применить метод классификации машинного обучения.

### **1.3.1 Формальная постановка задачи классификации**

Пусть  $X$  — множество описаний объектов,

$Y$  - конечное множество меток классов.

Существует неизвестная целевая зависимость — отображение  $y^* : X \rightarrow Y$ , значения которой известны только на объектах конечной обучающей выборки  $X^m = \{(x_1, y_1), \dots, (x_m, y_m)\}$ .

Требуется построить алгоритм  $a : X \rightarrow Y$ , способный классифицировать произвольный объект  $x \in X$ .

Для решения этой задачи могут быть использованы следующие методы классификации [14]:

- вероятностные («наивный» Байес);
- метрические (метод  $k$  ближайших соседей);
- логические (решающее дерево);
- линейные (метод опорных векторов, логистическая регрессия);
- методы на основе искусственных нейронных сетей (FFBP, RNN, CNN).

Рассмотрим наиболее популярные подходы для каждого из методов.

### 1.3.2 Наивный байесовский классификатор

**Метод Байеса (Naïve Bayes, NB)** относится к вероятностным методам классификации.

Пусть  $P(c_i|d)$  - вероятность того, что объект, представленный вектором  $d = (t_1, \dots, t_n)$ , соответствует категории  $c_j$  для  $i = 1, \dots, |C|$ . Задача классификатора заключается в том, чтобы подобрать такие значения  $c_i$  и  $d$ , при которых значение вероятности  $P(c_i|d)$  будет максимальным:

$$CSV(d) = \underset{c_i \in C}{argmax} P(c_i|d) \quad (18)$$

Для вычисления значений  $P(c_i|d)$  используется теорема Байеса:

$$P(c_i|d) = \frac{P(c_i)P(d|c_i)}{P(d)}, \quad (19)$$

где

- $P(c_i)$  - априорная вероятность того, что вектор  $d$  отнесен к категории  $c_i$ ;
- $P(d|c_i)$  - вероятность найти вектор  $d = (t_1, \dots, t_n)$  в категории  $c_i$ ;
- $P(d)$  - вероятность того, что произвольно взятый объект можно представить в виде вектора признаков  $d = (t_1, \dots, t_n)$ .

$P(d)$  не зависит от категории  $c_i$ , а значения  $t_1, \dots, t_n$  заданы заранее, поэтому знаменатель не влияет на выбор наибольшего из значений  $P(c_i|d)$ .

Вычисление  $P(d|c_i)$  затруднительно из-за большого количества признаков  $t_1, \dots, t_n$  поэтому делают «наивное» предположение о том, что любые случайные величины не зависят друг от друга.

Исходя из этого допущения можно воспользоваться формулой:

$$P(d|c_i) = \prod_{k=1}^n P(t_k|c_i) \quad (20)$$

#### Преимущества метода:

- высокая скорость работы;
- поддержка инкрементного обучения;
- относительно простая реализация программная реализация алгоритма;

- легкая интерпретируемость результатов алгоритма.

### Недостатки метода:

- относительно низкое качество классификации по сравнению с другими методами;
- неспособность учитывать зависимость результата классификации от сочетания признаков из-за допущения о их независимости друг от друга.

#### 1.3.3 Метод $k$ -ближайших соседей

**Метод  $k$  ближайших соседей (k Nearest Neighbours, KNN)** относится к метрическим методам классификации.

Чтобы найти категорию, соответствующую объекту  $d$ , классификатор сравнивает  $d$  со всеми объектами из обучающей выборки  $L$ , то есть для каждого  $d_z \in L$  вычисляется расстояние  $(d_z, d)$ . Далее из обучающей выборки выбираются  $k$  объектов, ближайших к  $d$ . Согласно методу  $k$  ближайших соседей, объект  $d$  считается принадлежащим тому классу, который является наиболее распространенным среди соседей данного объекта, то есть для каждого класса  $c_i$  вычисляется функция ранжирования:

$$CSV(d) = \sum_{d_z \in L_k(d)} p(d_z, d) * \Phi(d_z, c_i), \quad (21)$$

где

- $L_k(d)$  - ближайшие  $k$  объектов из  $L$  к  $d$ ;
- $\Phi(d_z, c_i)$  - количество объектов  $d_z$  обучающей выборки относящиеся к категории  $c_i$ .

### Преимущества метода:

- возможность обновления обучающей выборки без переобучения классификатора;
- устойчивость алгоритма к аномальным выбросам в исходных данных;
- относительно простая программная реализация алгоритма;
- легкая интерпретируемость результатов работы алгоритма;
- хорошее обучение в случае с линейно неразделимыми выборками.

### **Недостатки метода:**

- репрезентативность набора данных, используемого для алгоритма;
- высокая зависимость результатов классификации от выбранной метрики;
- большая длительность работы из-за необходимости полного перебора обучающей выборки;
- невозможность решения задач большой размерности по количеству классов и объектов.

#### **1.3.4 Метод деревьев решений (Decision Trees, DT)**

**Метод деревьев решений (Decision Trees, DT)** относится к логическим методам классификации.

Деревом решений называют ациклический граф, по которому производится классификация объектов, описанных набором признаков. Каждый узел дерева содержит условие ветвления по одному из признаков. У каждого узла столько ветвлений, сколько значений имеет выбранный признак. В процессе классификации осуществляются последовательные переходы от одного узла к другому в соответствии со значениями признаков объекта.

Классификация считается завершенной, когда достигнут один из листьев (конечных узлов) дерева. Значение этого листа определит класс, которому принадлежит рассматриваемый объект. На практике обычно используют бинарные деревья решений, в которых принятие решения перехода по ребрам осуществляется простой проверкой наличия признака в объекте. Если значение признака меньше определенного значения, выбирается одна ветвь, если больше или равно, другая.

### **Преимущества метода:**

- относительно простая программная реализация алгоритма;
- легкая интерпретируемость результатов работы алгоритма.

### **Недостатки метода:**

- неустойчивость алгоритма к выбросам исходных данных;
- требуется большой объем данных для получения точных результатов.



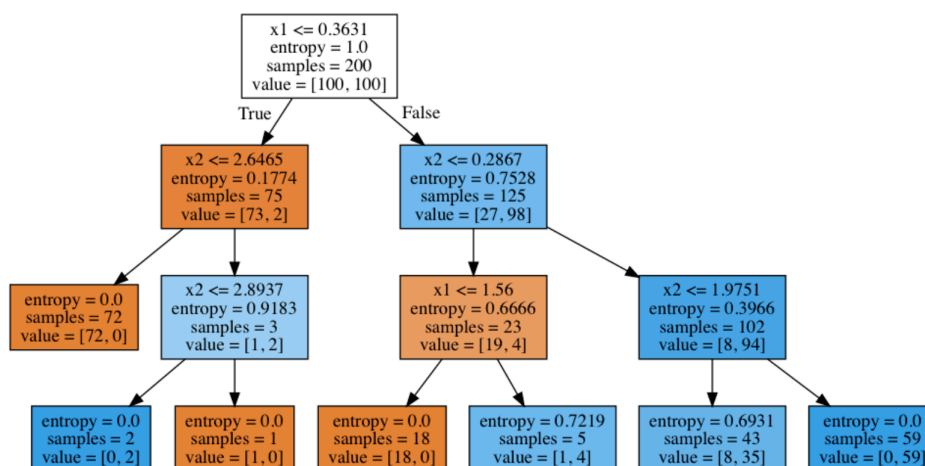


Рисунок 9 – Пример построенного бинарного дерева решений по двум признакам

### 1.3.5 Метод опорных векторов

**Метод опорных векторов (Support Vector Machine, SVM)** является линейным методом классификации. Рассмотрим множество объектов, которые необходимо классифицировать. Сопоставим ему множество точек в пространстве размерности  $|D|$ .

Выборку точек называют линейно разделимой, если принадлежащие разным классам точки можно разделить с помощью гиперплоскости (в двумерном случае - прямой). Необходимо провести прямую так, чтобы по одну сторону от нее лежали все точки одного класса, а по другую - все точки другого класса. Тогда для классификации неизвестных точек достаточно будет посмотреть, с какой стороны прямой они окажутся.

В методе опорных векторов расстоянием между прямой и множеством точек считается расстояние между прямой и ближайшей к ней точкой из множества. Именно такое расстояние и максимизируется в данном методе. Гиперплоскость, максимизирующая расстояние до двух параллельных гиперплоскостей, называется разделяющей (на рисунке 10 обозначена буквой  $L$ ). Ближайшие к параллельным гиперплоскостям точки называются опорными векторами, через них проходят пунктирные линии.

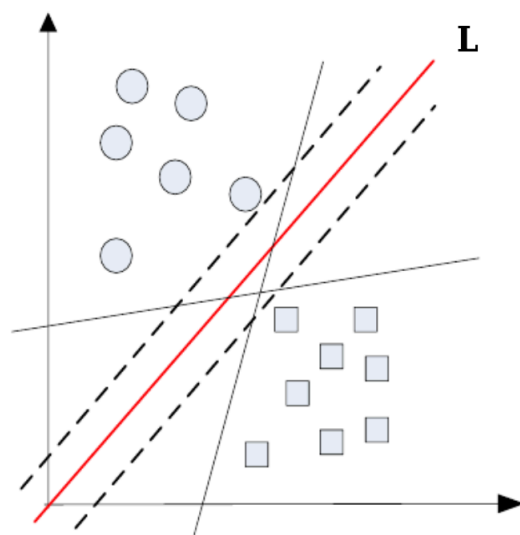


Рисунок 10 – Разделяющая гиперплоскость в методе опорных векторов

#### **Преимущества метода:**

- один из наиболее качественных методов;
- возможность работы с небольшим набором данных для обучения;
- сводимость к задаче выпуклой оптимизации, имеющей единственное решение.

#### **Недостатки метода:**

- сложная интерпретируемость параметров алгоритма;
- неустойчивость к выбросам в исходных данных.

#### **1.3.6 Методы на основе искусственных нейронных сетей**

Нейронная сеть — это система соединённых и взаимодействующих между собой простых процессоров (нейронов), соединённых между собой синапсами.

Существует большое количество разновидностей нейронных сетей, основные из них – сети прямого распространения, рекуррентные сети, радиально-базисные функции и самоорганизующиеся карты.

Рассмотрим структуру классической нейронной сети прямого распространения (Feed Forward Back Propagation, FFBP).

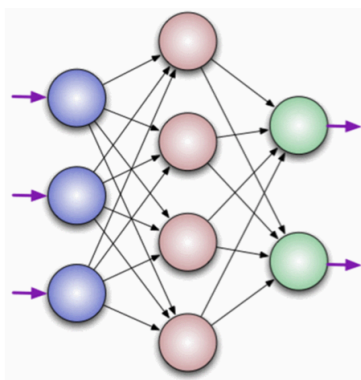


Рисунок 11 – Пример нейронных слоев (синий - входной, красный - скрытый, зеленый - выходной)

В нейронных сетях прямого распространения присутствуют входной слой, выходной слой и промежуточные (скрытые) слои: сигнал идет последовательно от входного слоя нейронов по промежуточным слоям к выходному.

Для классификации объекта  $d_i$  при помощи нейронной сети прямого распространения веса признаков документа подаются на соответствующие входы сети. Активация распространяется по сети; значения, получившиеся на выходах, есть результат классификации. Стандартный метод обучения такой сети - метод обратного распространения ошибки.

#### **Преимущества метода:**

- имеет очень высокое качество алгоритма при удачном подборе параметров;
- является универсальным аппроксиматором непрерывных функций;
- поддерживает инкрементное обучение.

#### **Недостатки метода:**

- вероятность возможной расходимости или медленной сходимости, поскольку для настройки сети используются градиентные методы;
- необходимость очень большого объема данных для обучения, чтобы достичь высокой точности;
- низкая скорость обучения;
- сложная интерпретируемость параметров алгоритма.

## 1.4 Постановка задачи

В результате проведения анализа предметной области, существующих методов выделения признаков речевого сигнала и алгоритмов классификации, постановка задачи, которую необходимо решить в рамках данной работы, может быть определена следующим образом.

Требуется спроектировать метод выделения звуков естественного языка в звучащей речи, а также разработать программный модуль для проверки работоспособности метода.

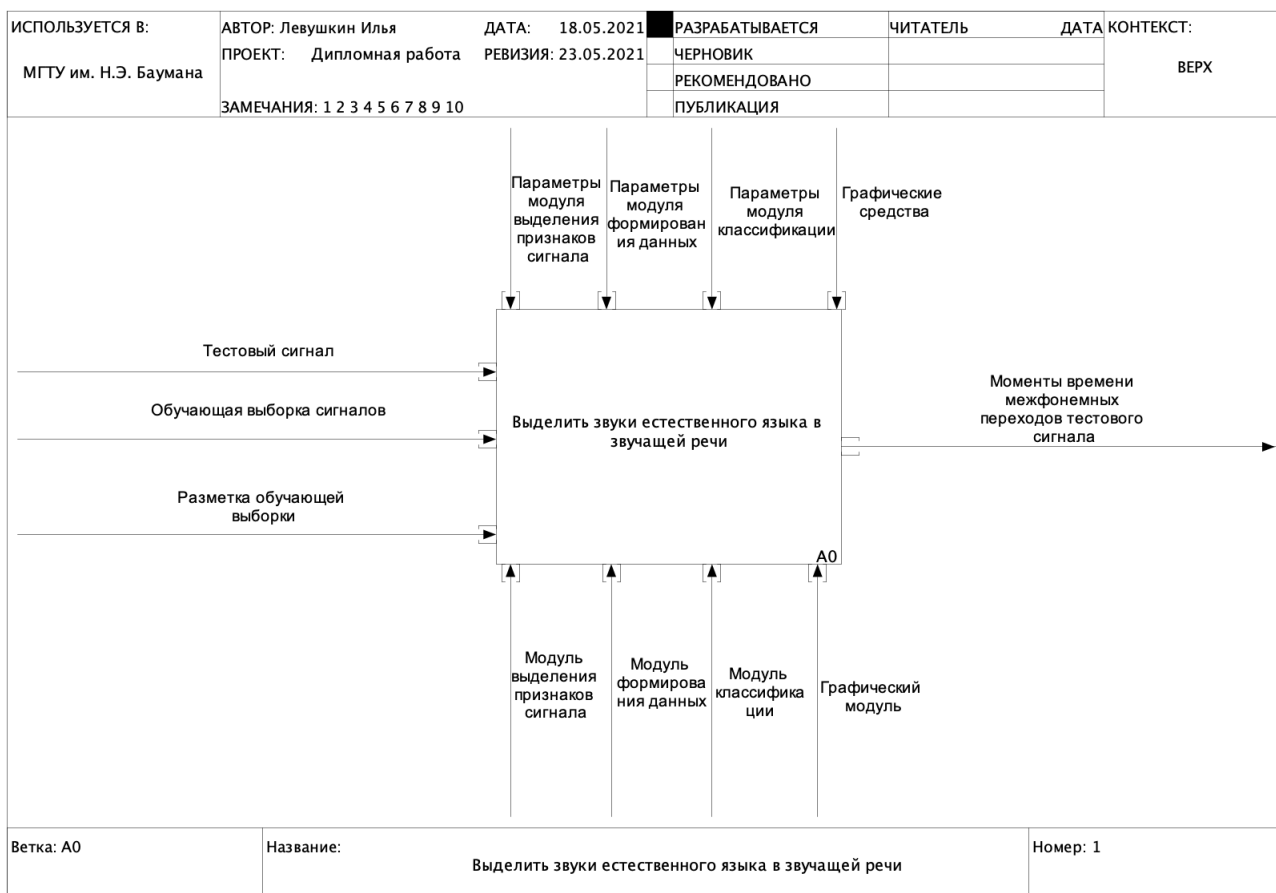


Рисунок 12 – IDEF0 диаграмма разрабатываемого модуля.

**Входными данными** являются аудиофайл и обучающая выборка аудиофайлов формата mp3, содержащие в себе информацию об амплитуде и частоте звуков в дискретные моменты времени, а также разметка обучающей выборки аудиофайлов в csv формате.

**Выходным параметром** является список моментов времени, на которых сигнал содержит межфонемные переходы.

Спроектированный метод должен обеспечить возможность

- выбрать тестовый аудиофайл;
- вычислить список моментов времени, на которых сигнал содержит межфонемные переходы и визуализировать их графически на осциллограмме (амплитудно-временном графике сигнала) [16].

### 1.5 Выводы

В ходе анализа предметной области были описаны основные этапы распознавания речевого сигнала, а также было дано обоснование актуальности задачи сегментации речевого сигнала.

Был проведен анализ существующих решений сегментации речевого сигнала на фрагменты, на основе которых были выделены два основных этапа в задаче сегментации:

- выделение признаков сигнала;
- определение межфонемных переходов на основе выделенных признаков сигнала.

Был произведен анализ существующих методов выделения признаков речевого сигнала. Ниже приведена сводная таблица особенностей каждого из рассмотренных методов относительно друг друга.

Метод	Базис	Представление	Нелинейность	Нестационарность	Вычислительная сложность
Фурье	Априорный	Энергия-частота	Нет	Нет	Низкая
Вейвлет	Априорный	Энергия-время-частота	Нет	Да	Средняя
Гильберт-Хуанг	Адаптивный	Энергия-время-частота	Да	Да	Высокая

Таблица 1 – Сравнение существующих методов выделения признаков речевого сигнала.

В данной работе было решено использовать метод, основанный на Вейвлет преобразовании, а именно, многомасштабный Вейвлет-анализ, поскольку он учитывает нестационарность человеческой речи, а также имеет относительно невысокую вычислительную сложность по сравнению с преобразованием Гильберта-Хуанга, что имеет большое значение в системах распознавания речи.

Также, было обнаружено, что большинство работ используют закономерности, эвристики, основанные на экспериментальном опыте других работ и не имеющие под собой существенной теоретической базы, в виду чего эти решения не гарантируют, полученной в результате экспериментов, точности на других данных.

Поэтому, в качестве альтернативы было решено использовать наиболее оптимальный для задачи определения межфонемных интервалов сигнала алгоритм классификации на основе выделенных признаков сигнала.

Ниже приведена сводная таблица наиболее важных особенностей для решения поставленной задачи каждого из методов классификации относительно друг друга для малого объема обучающей выборки.

Метод	Качество обучения	Скорость работы	Поддержка инкрементного обучения
Наивный Байес	Низкое	Высокая	Есть
k-ближайших соседей (KNN)	Среднее	Низкая	Есть
Деревья решений (DT)	Среднее	Средняя	Нет
Опорные вектора (SVM)	Высокое	Средняя	Нет
Искусственные нейронные сети	Среднее	Низкая	Есть

Таблица 2 – Сравнение существующих методов классификации объектов.

В данной работе было решено использовать метод опорных векторов из-за недостаточно большого объема обучающей выборки. Метод SVM дает наилучшее качество обучения на небольшом наборе данных, а также не требует значительных ресурсных затрат для обучения и предсказания результатов по сравнению с другими методами.

## **2 Конструкторская часть**

Данный раздел содержит в себе следующие задачи:

- изложение принципов спроектированного метода выделения звуков естественного языка в звучащей речи;
- подробное описание всех этапов разработанного алгоритма, реализующего данный метод;
- проектирование архитектуры системы, проверяющей работоспособность данного метода.
- описание формата входных и выходных данных и структуры ПО;

### **2.1 Метод выделения звуков естественного языка в звучащей речи**

Данный метод состоит из двух основных этапов:

- формирования признаков из поступающего на вход речевого сигнала одним из существующих методов выделения признаков сигнала;
- получения границ межфонемных переходов из сформированных признаков сигнала при помощи применения алгоритмов классификации.

На вход данный метод получает цифровой сигнал, содержащий в себе значения амплитуд (уровней сигнала) в дискретные моменты времени, а на выходе выдает список границ (моментов времени в секундах) межфонемных переходов (переходов с одного лингвистического элемента языка [1] к другому) естественного языка в звучащей речи.

### **2.2 Описание алгоритма метода выделения звуков естественного языка в звучащей речи**

Алгоритм, реализующий данный метод, состоит из следующих этапов:

- 1) выделение признаков речевых сигналов из проверяемого сигнала и обучающей выборки (рисунок 13 блок А1);
- 2) формирование данных для применения на них (обучения и предсказания результатов) модели классификации (рисунок 13 блок А2);
- 3) обучение модели классификации на сформированной обучающей выборке (рисунок 13 блок А3);
- 4) применение модели классификации на тестовых данных (проверяемом сигнале) для получения предсказанных значений - межфонемных переходов (рисунок 13 блок А4);



5) графическое представление результатов предсказания (рисунок 13 блок A5).

Ниже приведена IDEF0-диаграмма иллюстрирующая основные этапы алгоритма.

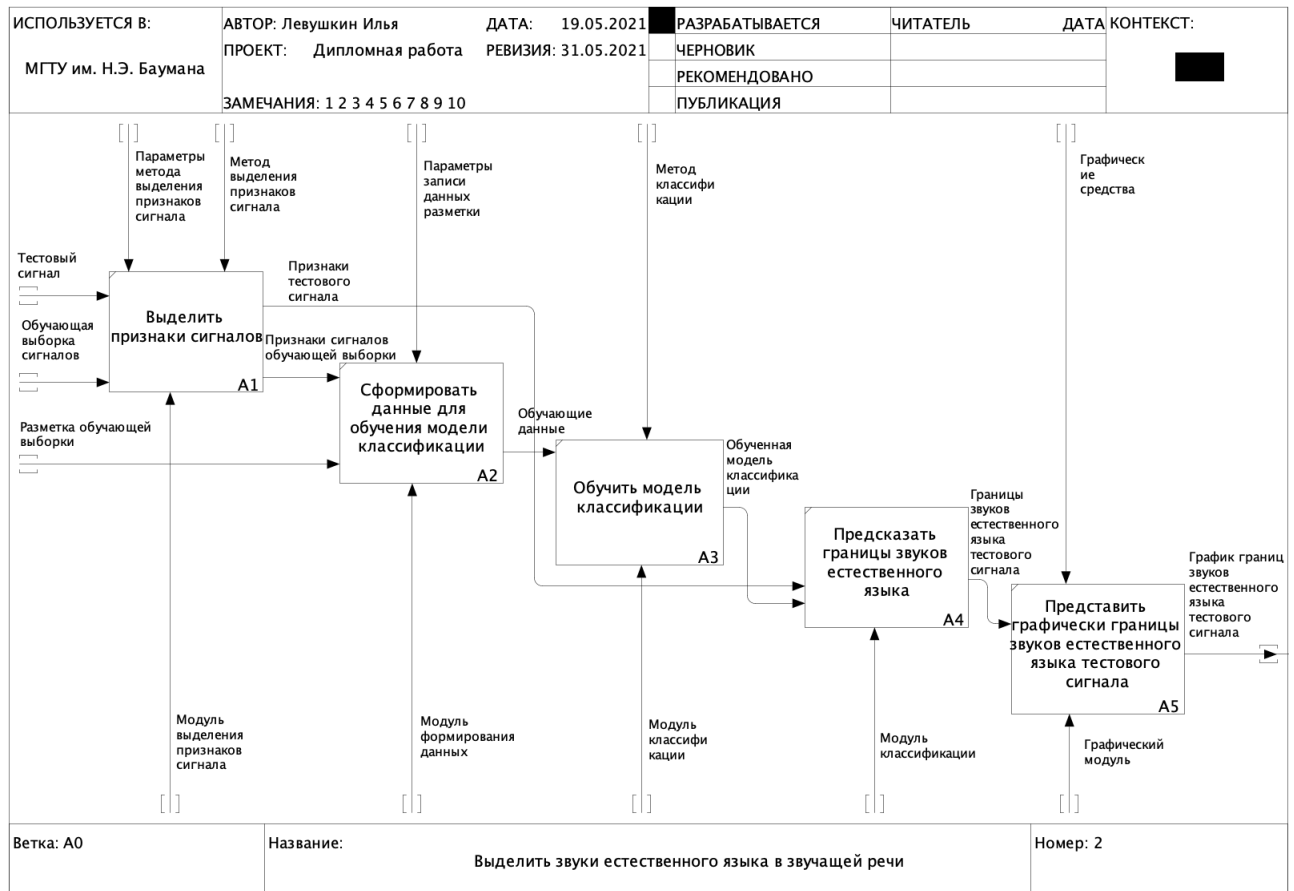


Рисунок 13 – IDEF0 диаграмма разрабатываемого алгоритма.

### 2.3 Выделение признаков речевого сигнала

В данной работе было решено использовать метод многомасштабного Вейвлет-анализа сигнала (раздел 1.5), в результате чего схема алгоритма выделения признаков речевого сигнала выглядит следующим образом [7]:

- 1) В качестве предобработки сигнал нормализуется: все отсчеты делятся на максимальное значение, для установки единых пороговых значений для любых входных сигналов;
- 2) Входной сигнал разбивается на фреймы по  $n_{frame}$  отсчётов при частоте дискретизации  $sr$  с перекрытием  $hop$ .
- 3) Каждый фрейм накрывается окном Хэмминга для устранения дефектов на краях;

- 4) К каждому фрейму применяется вейвлет-преобразование в базисе  $base$ . Используется разложение до  $k$ -го уровня декомпозиции;
- 5) Для каждого уровня декомпозиции определяется энергия, как сумма квадратов значений коэффициентов детализации  $E$  [6].

Ниже приведены варьируемые параметры, поступающие на вход модулю выделения признаков сигнала:

- $n_{frame}$  - количество отсчетов в одном фрейме (по умолчанию - 2500);
- $sr$  - частота дискретизации сигнала, Гц (по умолчанию - 16000);
- $hop$  - размер перекрытия фрейма другим фреймом, % (по умолчанию - 40);
- $base$  - базисная функция вейвлет-преобразования (по умолчанию - Вейвлет Добеши 8);
- $k$  - уровень декомпозиции сигнала (по умолчанию - 6).

## 2.4 Формирование данных для обучения модели классификации

Из формализованной постановки задачи классификации (раздел 1.3.1) следует, что для обучения модели классификации необходимо сформировать обучающую выборку таким образом, чтобы каждому вектору признаков  $(x^1, x^2, \dots, x^d)$  объекта  $x$  множества  $X^m$  соответствовало единственное значение ответа  $y = y(x)$  множества  $Y^m$ :

$$X^m = (x_1, y_1), \dots, (x_m, y_m). \quad (22)$$

В терминах данной задачи это означает, что каждому вектору признаков речевого сигнала за определенный промежуток времени  $t$  требуется сопоставить ответ на вопрос - *является ли данный интервал  $t$  границей межфонемного перехода?*

Говоря формально, требуется сформировать список ответов  $Y^m$ , элементы  $y_i$  которого принимают значения  $\{0, 1\}$ :

- 1 - в случае если объект  $x_i$  (вектор признаков речевого сигнала на интервале времени  $t_i$ ) является границей межфонемного перехода;
- 0 - в противном случае.

Исходя из приведенных выше рассуждений, схема формирования данных для модели классификации выглядит следующим образом:

- 1) Извлечение границ (моментов времени)  $T = (t_1, \dots, t_s)$  межфонемных переходов из файла с обучающей выборкой в соответствии с ее параметрами

записи данных;

- 2) Избавление от зашумленных данных путем усреднения значений  $t_i$  и  $t_{i+1}$  при условии, что  $t_{i+1} - t_i < \lim_{phonem}$ , где  $\lim_{phonem}$  - минимальный порог длительности фонемы (если расстояние между соседними значениями  $t$  меньше порогового, то заменяем их средним значением);
- 3) Преобразование границ межфонемных переходов  $T = (t_1, \dots, t_s)$  в список ответов  $Y = (y_1, \dots, y_m)$ , где  $y_i \in 0, 1$ .

## 2.5 Обучение модели классификации

В данной работе было решено использовать метод опорных векторов SVM (раздел 1.5) в качестве алгоритма классификации.

Ниже приведены настраиваемые параметры метода:

- kernel - ядро SVM, задает функцию разделяющей гиперплоскости. Наиболее популярные ядра - linear (линейная функция), rbf (радиально-базисная функция), poly (полиномиальная функция). В данной реализации метода выбрано по умолчанию ядро rbf;
- C - параметр регуляризации. Варьирование этого параметра позволяет понизить уровень переобучения модели, поскольку он контролирует соотношение между гладкой границей и корректной классификацией рассматриваемых точек. По умолчанию - 1.0;
- gamma - «ширина» ядра. Варьирование этого параметра также позволяет понизить уровень переобучения модели для нелинейных ядер. По умолчанию равен  $\frac{1}{n_{features} * X_{var}}$ ;
- class\_weight - параметр, позволяющий варьировать размер штрафа при неудачном выборе класса. Используется при несбалансированной выборке классов. По умолчанию равен balanced, то есть устанавливает штраф за ошибку выбора класса пропорционально размеру этого класса в обучающей выборке. Применяется, поскольку в данной задаче классы имеют значительный дисбаланс (интервалов времени, на которых присутствуют границы межфонемных переходов, значительно меньше чем тех, на которых их нет).

## 2.6 Структура ПО

Архитектура ПО использует внутри себя принципы архитектурного паттерна MVC (Model-View-Controller), который разделяет общую структуру кода

на три отдельных компонента:

- модель - предоставление данных конкретным элементам системы;
- представление - реагирует на изменение данных в системе и обеспечивает их отображение пользователю;
- контроллер - точка входа для доступа к модели и представлению. Контроллер обрабатывает действия пользователя, после чего отдает сигнал модели о необходимости каким-либо образом измениться.

При использовании MVC достигается возможность разделения компонентов работы с данными, пользовательским интерфейсом и логикой взаимодействия пользователя с приложением, в результате чего модификация одного из компонентов оказывает минимальное воздействие на остальные.

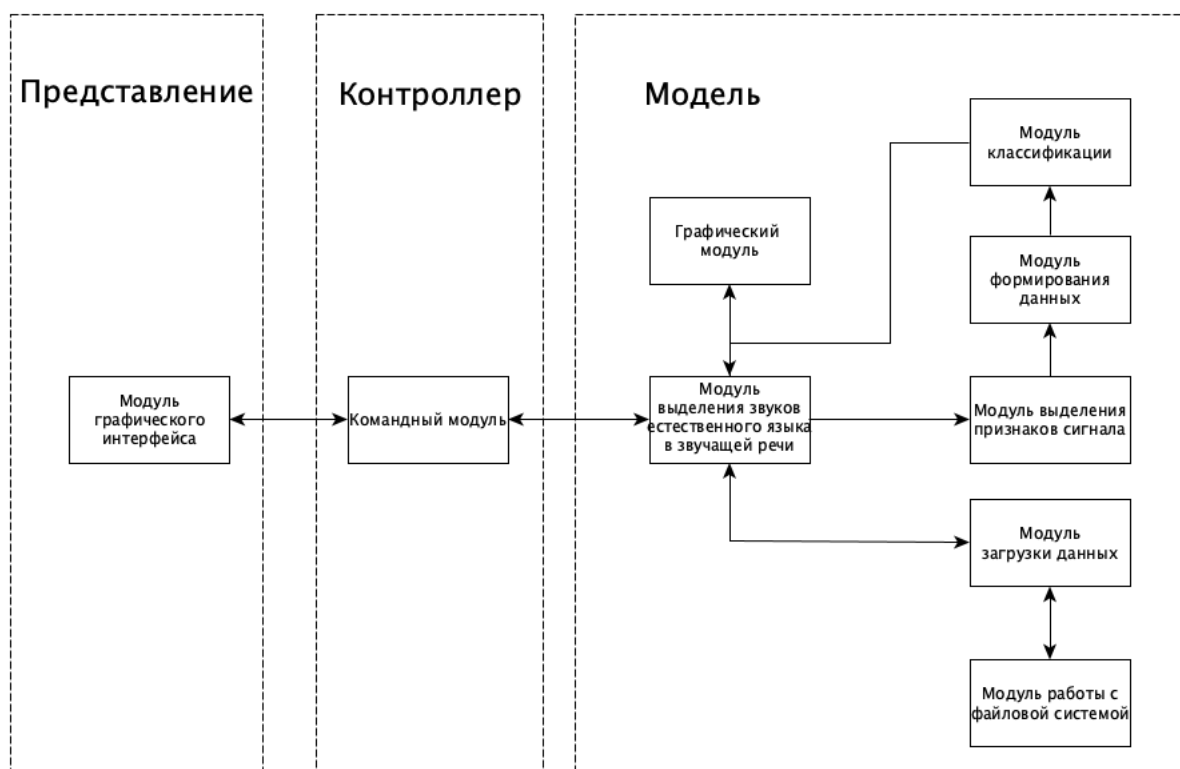


Рисунок 14 – Архитектура программного обеспечения.

Графический интерфейс предоставляет пользователю возможность взаимодействия с разработанным ПО.

Доступный функционал:

- возможность выбрать аудио-файл (формат mp3) для тестирования метода, аудио-файлы (формат mp3) и разметку (формат csv) для обучения модели

- классификации;
- возможность обучить модель классификации;
- возможность сохранить обученную модель классификации в файл (формат pkl);
- возможность загрузить обученную модель классификации из файла (формат pkl);
- возможность определения границ межфонемных переходов выбранного сигнала для тестирования;
- сохранение результатов в файл (формат csv);
- графическое отображение результатов в диалоговом окне.

Модули, изображенные на рисунке 14, содержат необходимые классы и методы, выполняющие основную логику приложения.

Ниже приведено описание основных классов, представленных на рисунке 15.

- Классы *MainView*, *ResultView*, *TeachModelView* - обрабатывают логику запросов пользователя к графическому интерфейсу приложения;
- Класс *MainController* агрегирует в себе классы *MainView*, *ResultController*, *TeachModelController*, *SoundsExtractionModule*. Является точкой входа в приложение и отвечает за распределение задач от *MainView* к *SoundsExtractionModule* и по своим контроллерам - *ResultController*, *TeachModelController*;
- *ResultController* и *TeachModelController* отвечают за распределение задач от *ResultView* и *TeachModelView* к *SoundsExtractionModule*;
- *SoundsExtractionModule* агрегирует в себе все модули: *GraphModule*, *FeatureExtractionModule*, *DataPreprocessingModule*, *ClassificationModule* и *LoadDataModule*. Отвечает за распределение задач от контроллеров по модулям;
- *LoadDataModule* отвечает за работу с файловой системой. Осуществляет загрузку и выгрузку данных: *test\_signal*, *train\_signals*, *train\_markup*, *model*;
- *ClassificationModule* отвечает за обучение модели и осуществление ею предсказаний;
- *DataPreprocessingModule* отвечает за предобработку данных. Формирует обучающую и тестовую выборки для *ClassificationModule*;

- *FeatureExtractionModule* отвечает за извлечение признаков из сигналов.

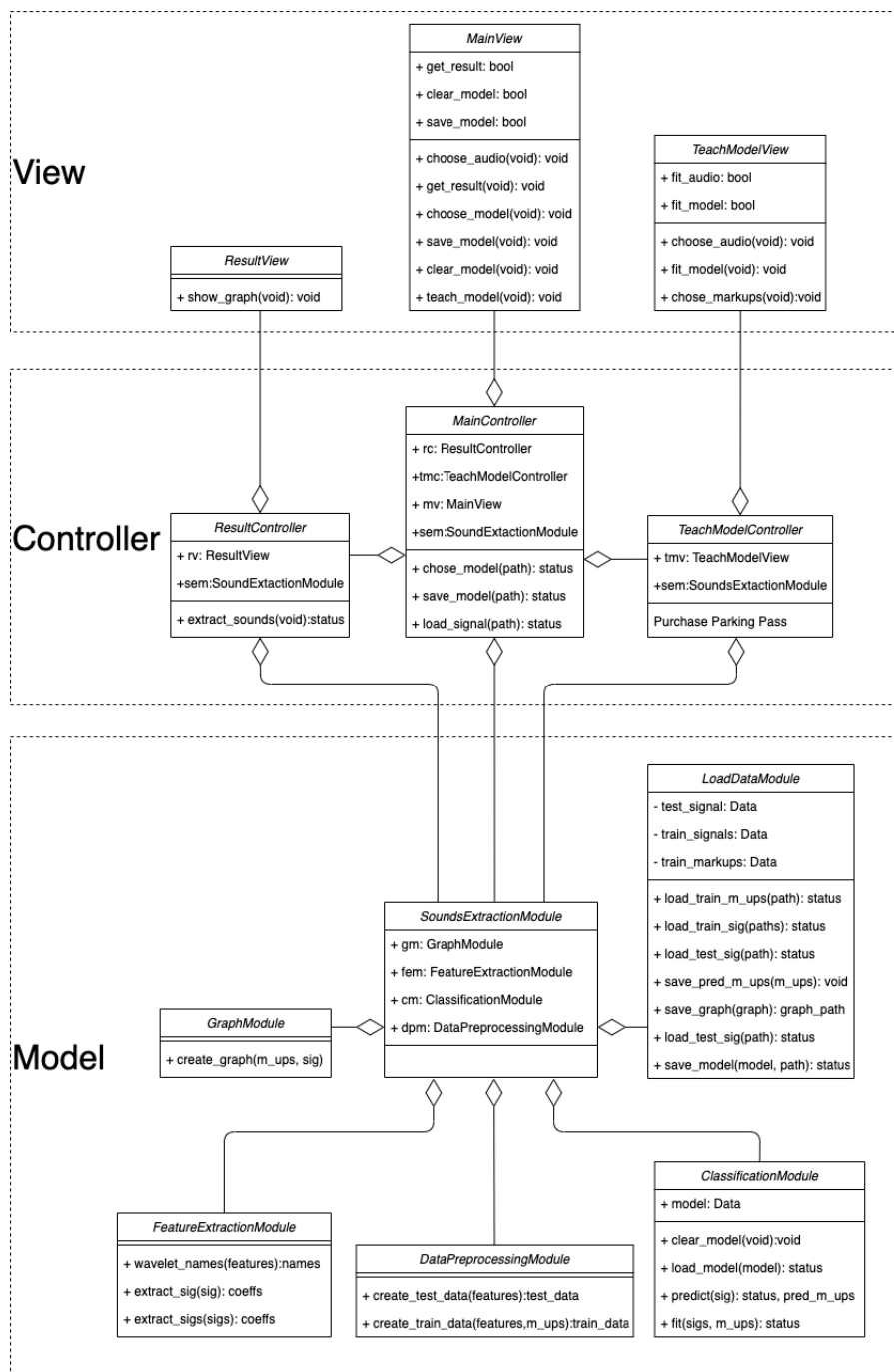


Рисунок 15 – Схематичное представление архитектуры программного обеспечения.

## 2.7 Описание формата входных и выходных данных

Для обучения модели пользователь выбирает файлы аудио-сигналов (формат mp3) и разметку этих аудио-файлов (формат csv). Разметка аудио-файлов представляет из себя таблицу со столбцами *filename*, *start*, *end*, *intervals*, где

- *filename* - имя аудио-файла;

- start - начало фонемы;
- end - конец фонемы;
- intervals - список границ, на которых был размечен аудио-файл.

Вместо обучения пользователь может загрузить файл обученной модели (формат pkl).

Для тестирования модели пользователь выбирает файл с аудио-сигналом (формат mp3).

**На выходе** пользователь получает график границ межфонемных переходов, выведенный на экран и сохраненный в файл (формат png), список выделенных границ межфонемных переходов, сохраненных в файл (формат csv).

Также, пользователь имеет возможность сохранить обученную модель в файл (в формате pkl).

## 2.8 Выводы

В результате написания конструкторского раздела были решены следующие задачи:

- спроектирован метод выделения звуков естественного языка в звучащей речи;
- разработан алгоритм, реализующий данный метод, и приведено подробное его описание;
- спроектирована архитектура системы для проверки работоспособности метода;
- описаны формат входных и выходных данных и структура разрабатываемого программного обеспечения.

### 3 Технологическая часть

#### 3.1 Выбор языка программирования

В качестве языка программирования был выбран язык программирования *Python3* - объектно-ориентированный интерпретируемый язык программирования высокого уровня [17].

Преимуществами данного языка являются:

- встроенные типы данных высокого уровня: списки, словари, строки и другие;
- объектно-ориентированность - язык предоставляет удобный интерфейс создания и использования объектов;
- широко используется в научных исследованиях, так как обладает большим объемом различных библиотек для решения прикладных задач из различных областей, включая такие области как цифровая обработка сигналов, машинное обучение, работа со сложными математическими вычислениями и большим объемом данных.

#### 3.2 Выбор используемых библиотек

Спроектированный метод выделения звуков естественного языка звучащей речи состоит из двух этапов (раздел 1.1.2):

- 1) выделение признаков сигнала;
- 2) определение межфонемных переходов на основе выделенных признаков сигнала.

**Первый этап** подразумевает работу с аудио-сигналами, а именно - загрузка аудио-файлов, дискретизация сигнала, разделение сигнала на перекрывающиеся участки - фреймы.

Для решения этих задач было решено использовать библиотеку *Librosa* [19] - модуль python, предоставляющий удобный интерфейс для анализа звуковых сигналов.

В качестве модуля выделения признаков сигнала был выбран *Многомасштабный Вейвлет-анализ* (раздел 1.5). Для решения этой задачи была выбрана библиотека *pywt* [20] - модуль python, предоставляющий необходимый инструментарий для применения вейвлет-преобразования на дискретном сигнале.

На **втором этапе** для определения границ между основными лингвисти-



ческими элементами языка было решено применить метод классификации машинного обучения (раздел 1.3.1) - метод опорных векторов SVM [15]. Наиболее популярным подходом к решению задач классификации является использование библиотеки *Scikit-Learn* [21] в связке с такими инструментами как *Pandas* [22], *NumPy* [24] и *Matplotlib* [23]. Такой подход позволяет наиболее эффективно использовать объем доступной памяти и процессорное время, предоставляемое системой.

- *Pandas* - библиотека python для обработки и анализа больших данных. Использует объекты *NumPy* и предоставляет специальные структуры данных и операции для манипулирования числовыми таблицами и временными рядами;
- *NumPy* - библиотека python, добавляющая поддержку больших многомерных массивов и матриц, вместе с большой библиотекой высокоуровневых и быстрых математических функций для операций с этими массивами;
- *Matplotlib* - библиотека python, используемая для построения графиков и диаграмм. Использует объекты *NumPy* и предоставляет объектно-ориентированный интерфейс для данных, отображая получившиеся результаты с помощью встроенных графических библиотек.

Графический интерфейс, используемый пользователем был написан при помощи библиотеки *PyQt5* [25] - графической библиотеки python, предоставляющей возможность создавать графические интерфейсы для пользователя. Данная библиотека использует понятную структуру наследования и предлагает объектно-ориентированные решения с логической иерархией между объектами.

### 3.3 Выбор среды разработки

В качестве среды разработки было решено использовать *PyCharm* [18], имеющий доступный интерфейс, а также большое количество полезных функций и плагинов, упрощающих процесс разработки.

Поскольку *PyCharm* нативно поддерживает Python, данная среда разработки выполняет большое количество рутинных действий за разработчика. Среда снабжена отладчиком, необходимым для разработки ПО. Также, среда предоставляет возможность удобного рефакторинга кода, меняя все зависимости во всем проекте.

### 3.4 Пользовательский интерфейс

Пользователь взаимодействует с графическим интерфейсом ПО курсором мыши и клавиатурой.

Ниже приведено описание взаимодействия пользователя с графическим интерфейсом ПО.

- **Обучение модели.**

Для обучения модели классификации, пользователю необходимо нажать на кнопку *Обучить модель классификации*.

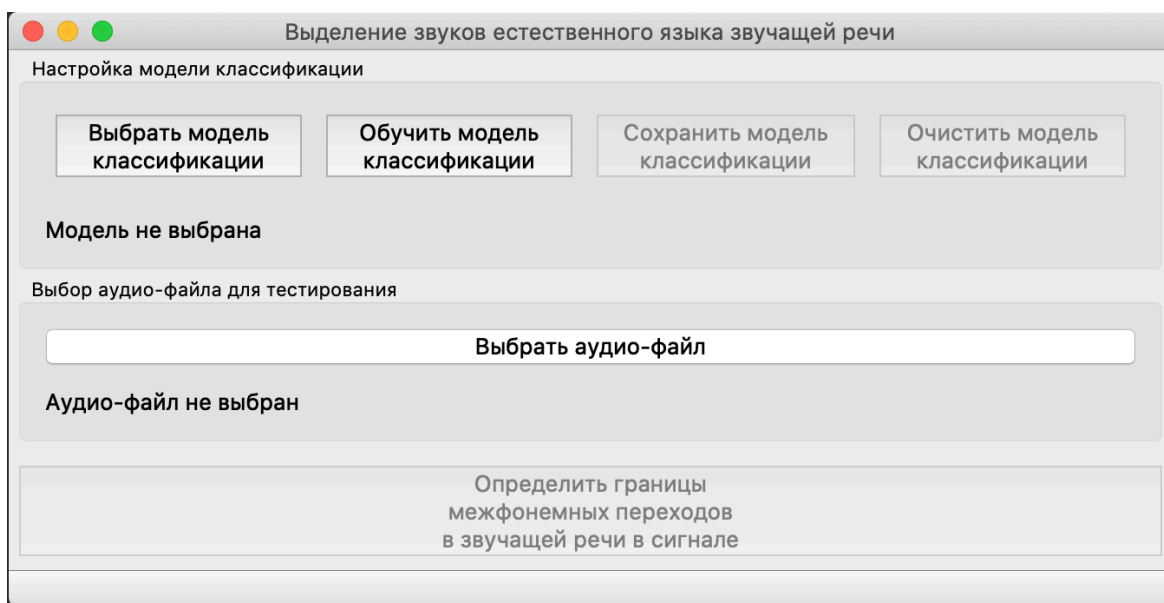


Рисунок 16 – Главное меню.

В открывшемся диалоговом окне необходимо выбрать файл-разметки аудио-сигналов для обучения модели, нажав на кнопку *Выбрать разметку для обучения*. Затем выбрать аудио-файлы для обучения модели, нажав на кнопку *Выбрать аудио-файлы для обучения*. После загрузки файлов, нажать кнопку *Обучить модель*, дождаться конца обучения и нажать кнопку *OK*.

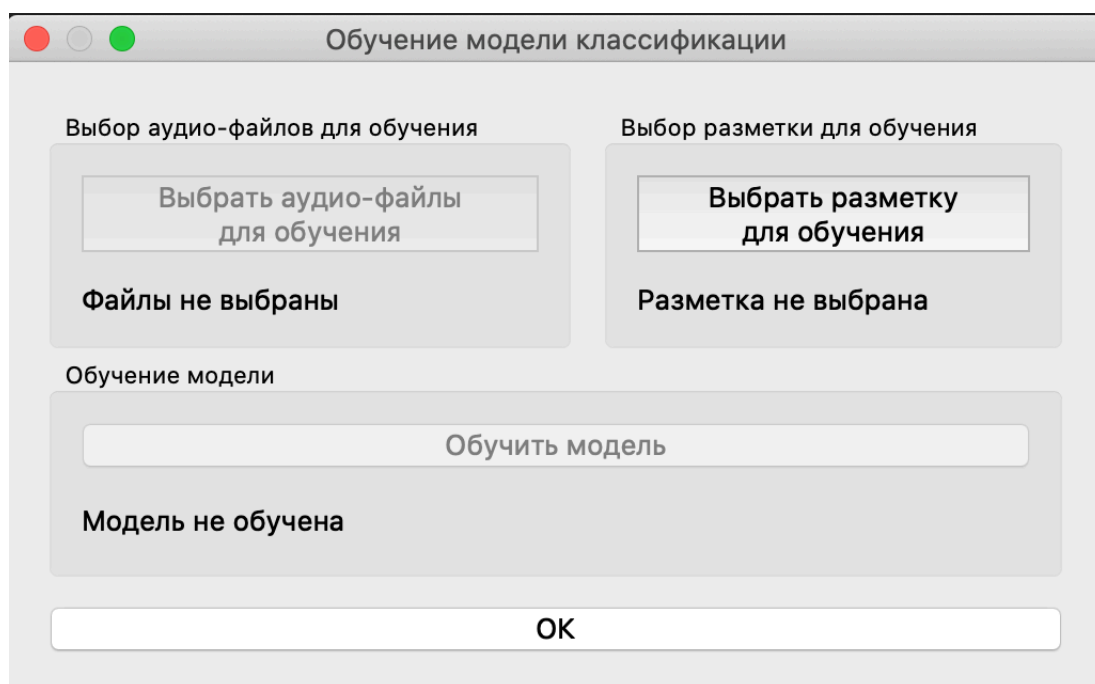


Рисунок 17 – Обучение модели.

- **Сохранение модели.**

После обучения модели появляется возможность сохранить модель. Для этого необходимо нажать на кнопку *Сохранить модель классификации*.

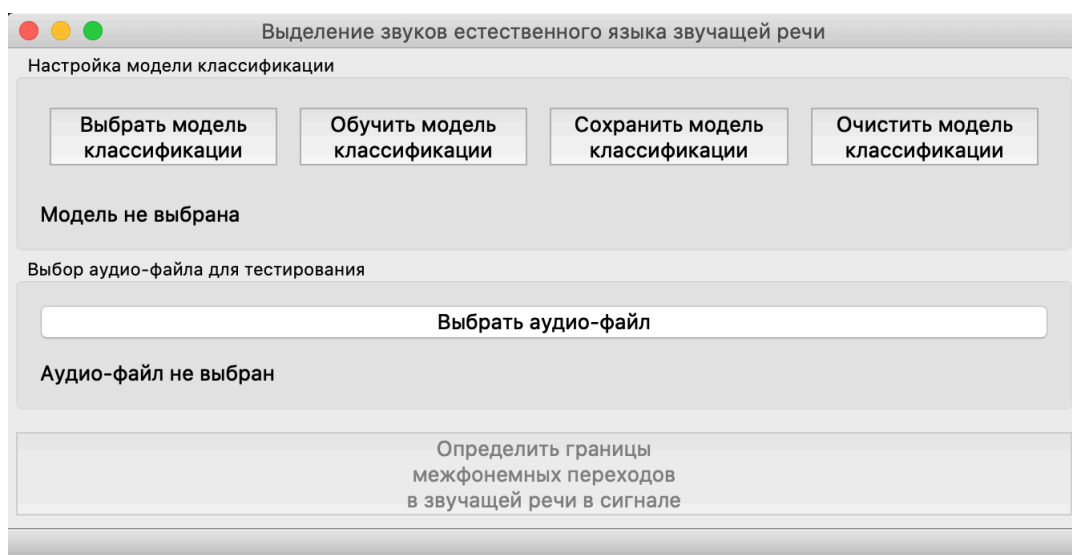


Рисунок 18 – Сохранение модели.

- **Загрузка модели.**

Вместо обучения модели пользователь имеет возможность загрузить файл обученной модели. Для этого необходимо нажать на кнопку *Выбрать модель классификации*.

- **Сброс модели.**

После загрузки или обучения модели классификации у пользователя появляется возможность сбросить модель. Для этого необходимо нажать на кнопку *Очистить модель классификации*.

- **Выбор аудио-файла.**

Пользователь имеет возможность загрузить аудио-файл для тестирования ПО. Для этого необходимо нажать на кнопку *Выбрать аудио-файл*.

- **Выделение звуков естественного языка звучащей речи**

Для получения результатов предсказания модели пользователь должен нажать на кнопку **Определить границы межфонемных переходов звучащей речи в сигнале** после выбора тестового аудио-файла и обучения/загрузки модели.

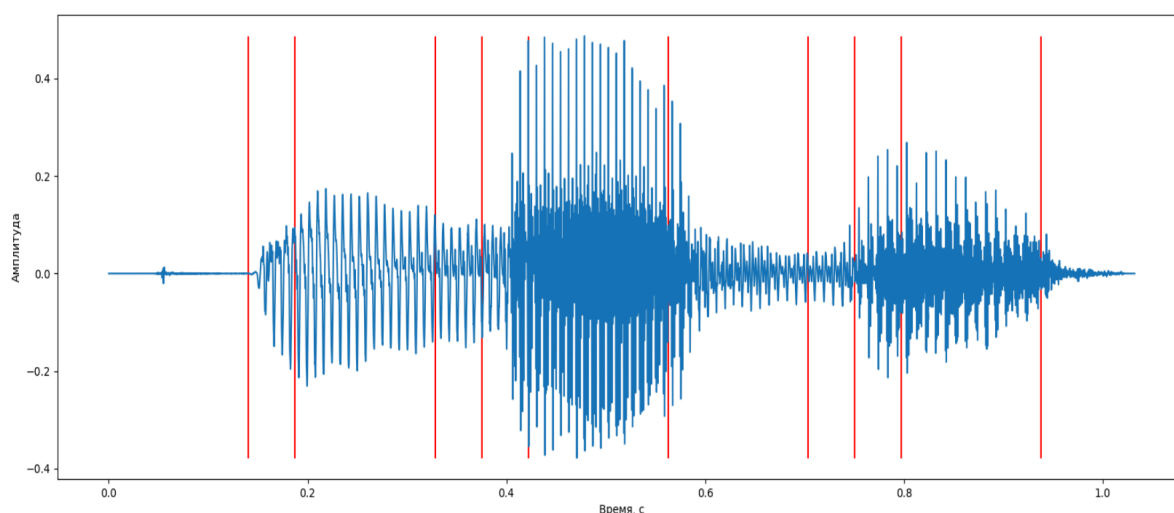


Рисунок 19 – Выделение границ межфонемных переходов звучащей речи в сигнале.

### 3.5 Реализация

Ниже приведены листинги кода реализаций некоторых компонентов ПО: блоков работы с данными - Model, блока работы с пользовательским интерфейсом - View, блока логики взаимодействия пользователя с приложением - Controller (полный листинг кода разработанного ПО в разделе *ПРИЛОЖЕНИЕ А*).

**MainView** - логика запросов пользователя к графическому интерфейсу главного экрана приложения.

## Листинг 1 – MainView

```
1 from PyQt5 import uic
2 from PyQt5.QtCore import pyqtSlot, Qt
3 from PyQt5.QtWidgets import QMessageBox, QMainWindow, QFileDialog
4
5
6 class MainWindow(QMainWindow):
7     def __init__(self, in_controller):
8         super(MainWindow, self).__init__()
9         self.ui = uic.loadUi("view/main_view/main_window.ui", self)
10        self.controller = in_controller
11        self.ui.get_result_button.setEnabled(False)
12        self.ui.clear_model_button.setEnabled(False)
13        self.ui.save_model_button.setEnabled(False)
14
15        @pyqtSlot(name='on_choose_model_button_clicked')
16        def choose_model(self):
17            filename, _ = QFileDialog.getOpenFileName(self, Выберите' модель для
загрузки',
18            './Data/models/', 'Model File (*.pkl)')
19            if len(filename) > 0:
20                status = self.controller.choose_model(filename)
21                if status == Модель' успешно загружена':
22                    self.ui.clear_model_button.setEnabled(True)
23                    self.try_enable_get_result_btn()
24            else:
25                msgBox = QMessageBox()
26                msgBox.setStandardButtons(QMessageBox.Ok)
27                msgBox.setIcon(QMessageBox.Critical)
28                msgBox.setText(status)
29                msgBox.exec()
30
31        @pyqtSlot(name='on_teach_model_button_clicked')
32        def teach_model(self):
33            self.controller.launch_teach_model_view()
34            self.ui.clear_model_button.setEnabled(True)
35            self.ui.save_model_button.setEnabled(True)
36
```

## Листинг 2 – MainView Продолжение

```
1  @pyqtSlot(name='on_save_model_button_clicked')
2  def save_model(self):
3      filename, _ = QFileDialog.getSaveFileName(self, Сохранить' модель',
4      './Data/models/', 'Model File (*.pkl)')
5      if len(filename) > 0:
6          status = self.controller.save_model(filename)
7          msgBox = QMessageBox()
8          msgBox.setStandardButtons(QMessageBox.Ok)
9          if status == Модель' успешно сохранена':
10             msgBox.setIcon(QMessageBox.Information)
11             msgBox.setText(status)
12             self.ui.save_model_button.setEnabled(False)
13         else:
14             msgBox.setIcon(QMessageBox.Critical)
15             msgBox.setText(status)
16             msgBox.exec()
17
18  @pyqtSlot(name='on_clear_model_button_clicked')
19  def clear_model(self):
20      self.controller.clear_model()
21      self.ui.clear_model_button.setEnabled(False)
22      self.ui.get_result_button.setEnabled(False)
23
24
25  @pyqtSlot(name='on_choose_test_audio_button_clicked')
26  def choose_test_audio(self):
27      filename, _ = QFileDialog.getOpenFileName(self, Выберите' mp3файл- для
28      тестирования',
29      './Data/test/', 'Audio File (*.mp3)')
30      if len(filename) > 0:
31          status = self.controller.load_test_signal(filename)
32          self.ui.choose_test_audio_label.setText(status)
33          self.try_enable_get_result_btn()
34
35  @pyqtSlot(name='on_get_result_button_clicked')
36  def get_result(self):
37      self.controller.launch_result_view()
```

**MainController** - точка входа в приложение. Отвечает за распределение задач от *MainView* к *SoundsExtractionModel* и по агрегируемым в себе контроллерам - *ResultController*, *TeachModelController*.

### Листинг 3 – MainController

```
1 class MainController:
2     def __init__(self):
3         self.main_window = MainWindow(self)
4         self.main_window.show()
5         self.sounds_extraction_module = SoundsExtractionModule()
6
7         # инициализация teach_model_controller
8         self.teach_model_controller = TeachModelController(self,
9             self.sounds_extraction_module)
10
11        # инициализация result_controller
12        self.result_controller = ResultController(self,
13            self.sounds_extraction_module)
14
15    def launch_teach_model_view(self):
16        self.teach_model_controller.show_window()
17
18    def teach_model_window_closed(self):
19        self.main_window.try_enable_get_result_btn()
20
21    def launch_result_view(self):
22        self.result_controller.show_window()
23
24    def clear_model(self):
25        self.sounds_extraction_module.clear_classification_model()
26
27    def load_test_signal(self, filename):
28        return self.sounds_extraction_module.load_test_signal(filename)
29
30    def choose_model(self, path):
31        return self.sounds_extraction_module.choose_model(path)
32
33    def save_model(self, path):
34        return self.sounds_extraction_module.save_classification_model(path)
```

**FeatureExtractionModule** - модуль, отвечающий за извлечение признаков ИЗ СИГНАЛОВ.

#### Листинг 4 – FeatureExtractionModule

```

1 from scipy import signal
2 from pywt import wavedec
3
4 from model.params.audio_params import AudioParams
5
6 class FeatureExtractionModule:
7
8     @staticmethod
9     def get_t_between_frames(frames_coeffs, offset):
10         t_values = []
11         cur_time = offset
12         dt = AudioParams.frame_sz() * AudioParams.hop_part() / AudioParams.sr()
13         for i in range(len(frames_coeffs)):
14             t_values.append(cur_time)
15             cur_time += dt
16         return np.array(t_values)
17
18     # получает массив сигналов: [[path, signal],...], возвращает словарь
    признаков: {path: feature}
19     def extract_from_signals(self, sigs):
20         res = []
21         for path_sig in sigs:
22             res.append([path_sig[0], self.extract_from_signal(path_sig[1])])
23         return res
24
25     # получает на вход signal, возвращает признаки signal
26     def extract_from_signal(self, sig):
27         # разбиение сигнала на фреймы
28         frames_sig = self._splitting_signal_on_frames(sig,
29             AudioParams.frame_sz(), AudioParams.hop_part())
30
31         # устранение деффекта на краях с помощью оконной функции Хамминга
32         frames_hamming_sig = self._hamming_func(frames_sig,
33             AudioParams.frame_sz())
34
35         # применение вейвлетпреобразования-
36         frames_coeffs = self._frames_wavelet(frames_hamming_sig,
37             AudioParams.wavelet(), AudioParams.wavelet_level())
38
39         return frames_coeffs

```



#### Листинг 5 – FeatureExtractionModule Продолжение

```
1 def _splitting_signal_on_frames(self, descrete_sig, frame_sz, hop_part):
2     frames = librosa.util.frame(descrete_sig, frame_length=frame_sz,
3     hop_length=int(hop_part * frame_sz), axis=0)
4     return frames
5
6 def _hamming_func(self, frames_signal, frame_sz):
7     w = signal.windows.hamming(frame_sz, sym=True)
8     return np.array([w * frame for frame in frames_signal])
9
10 def _wavelet_conv(self, sig, method, level):
11     return np.array(wavedec(sig, method, level=level))
12
13 def _frames_wavelet(self, frames_signal, method, level):
14     return np.array([self._wavelet_conv(sig, method, level) for sig in
15     frames_signal], dtype=object)
```

**DataPreprocessingModule** - модуль, формирующий обучающую и тестовую выборки для *ClassificationModule*.

#### Листинг 6 – DataPreprocessingModule

```
1 import numpy as np
2 from itertools import chain
3 import pandas as pd
4 import ast
5
6 from model.feature_extraction_module.feature_extraction_module import
7     FeatureExtractionModule
8 from model.params.markups_params import MarkupsParams
9 from model.params.audio_params import AudioParams
10
11 class DataPreprocessingModule:
12     # получает на вход signal, возвращает признаки signal
13     def create_test_data(self, features):
14         features_names =
15             FeatureExtractionModule.get_wavelet_columns_names(features[0])
16         res_data = []
17         for feature in features:
18             res_data.append(list(chain(*feature)))
19         return pd.DataFrame(res_data, columns=features_names)
```

```

1  # получает массив признаков [[path, feature],...], разметку в pandas с
    колонками: ['recordname', 'start', 'end']
2  # возвращает обучающую выборку в pandas формате: [feature, feature, ...,
    y], y = {0,1}
3  def create_train_data(self, features, df_markup):
4      # получаем названия признаков
5      features_names =
        FeatureExtractionModule.get_wavelet_columns_names(features[0][1][0])
6      features_names.append('y')
7      res_data = []
8      for path, frames_coeffs in features:
9          filename = self._get_filename_from_path(path)
10         for interval in
            ast.literal_eval(df_markup[df_markup[MarkupsParams.filename()] ==
                filename]
11             .intervals.iloc[0]):
12             # получаем временные промежутки между фреймами
13             t_positions_between_frames =
                FeatureExtractionModule.get_t_between_frames(frames_coeffs, interval[0])
14             # формируем список Вейвлеткоэффициентов- (n) - один фрейм
15             # длина( списка - t_positions_between_frames) и добавляем к нему
            y={0,1}
16             y = np.zeros(frames_coeffs.shape[0])
17             borders = self._borders_preprocessing(
18                 df_markup[df_markup[MarkupsParams.filename()] == filename],
19                 MarkupsParams.start_end()[0], MarkupsParams.start_end()[1])
20             # оставляем только те границы, что лежат в интервале interval
21             borders = self._filter_borders(borders, interval)
22             for border in borders:
23                 for t_pos in range(len(t_positions_between_frames)):
24                     if t_positions_between_frames[t_pos] > border:
25                         y[t_pos] = 1.
26                         break
27             # преобразуем frames_coeffs
28             res_data_i = []
29             for frame, y_i in zip(frames_coeffs, y):
30                 lst_i = list(chain(*frame))
31                 lst_i.append(y_i)
32                 res_data_i.append(lst_i)
33             res_data.extend(res_data_i)
34             return pd.DataFrame(res_data, columns=features_names)
35

```

**ClassificationModule** - модуль, отвечающий за обучение модели и осуществление ею предсказаний.

## Листинг 8 – ClassificationModule

```
1 from model.data import Data
2 from model.params.classification_params import ClassificationParams
3
4 class ClassificationModule:
5     def __init__(self, in_data_preprocessing_module,
6                 in_feature_extraction_module):
7         self.data_preprocessing_module = in_data_preprocessing_module
8         self.feature_extraction_module = in_feature_extraction_module
9         self.params = ClassificationParams()
10        self.model = DataМодель(' не выбрана')
11
12    def clear_model(self):
13        self.model = DataМодель(' не выбрана')
14
15    def fit(self, signals, markups):
16        if self.model.status == Модель' обучена':
17            return self.model.status
18        try:
19            # получает массив сигналов: [[path, audio],...], возвращает массив
20            признаков: [[path, feature],...]
21            features = self.feature_extraction_module.extract_from_signals(signals)
22            # получает массив признаков [[path, feature],...],
23            # разметку в pandas с колонками: ['filename', 'start', 'end']
24            # возвращает обучающую выборку в pandas формате: [feature, feature,
25            ..., y], y = {0,1}
26            train_data =
27            self.data_preprocessing_module.create_train_data(features, markups)
28
29            self.params.model.fit(train_data.drop(axis=0, columns=['y']),
30                                train_data['y'].astype('int64'))
31
32            self.model.data = self.params.model
33            self.model.status = Модель' обучена'
34        except:
35            return Ошибка' обучения модели'
36        return self.model.status
37
```

#### Листинг 9 – ClassificationModule Продолжение

```
1  # возвращает статус и список predicted_markupс временных промежутков
2  def predict(self, signal):
3      if self.model.status != Модель' обучена':
4          return self.model.status
5      try:
6          features = self.feature_extraction_module.extract_from_signal(signal)
7          # получает на вход признаки signal, возвращает pandas признаки с их
            именами
8          test_data = self.data_preprocessing_module.create_test_data(features)
9          predictions = self.model.data.predict(test_data)
10         return Звуки' успешно выделены',
            self._convert_predictions_to_time_list(predictions)
11     except:
12         return Ошибка'! Не удалось выделить звуки', None
13
```

### 3.6 Выводы

В данном разделе были выполнены следующие задачи:

- была реализована и протестировано программное обеспечение, спроектированное в разделе 2, для проверки работоспособности разработанного метода выделения звуков естественного языка в звучащей речи;
- было дано обоснование выбора языка программирования, среды разработки и используемых библиотек;
- было дано описание пользовательского интерфейса;
- были приведены листинги кода реализаций основных компонентов программного обеспечения.

## 4 Исследовательская часть

### 4.1 Формирование обучающей и тестовой выборок данных

В качестве данных для обучения и тестирования разработанного метода был задействован, размеченный студентами с кафедры Л4 и университета МГЛУ, речевой корпус. Данные представляли из себя аудио-файлы в mp3 формате и разметку к ним в формате json.

В результате выполнения предобработки данных было получено 5806 объектов по 7 признаков в каждом. Предобработка данных выполнялась при следующих начальных параметрах (раздел 2.3):

- $n_{frame} = 2500$ ;
- $sr = 16000$  Гц;
- $hop = 40\%$ ;
- $base = db4$  (Вейвлет Добеши 8);
- $k = 6$ .

Ниже приведен вывод первых пяти сформированных объектов выборки.

	A	D6	D5	D4	D3	D2	D1	y
0	0.486720	0.036786	0.010464	0.006818	0.004979	0.001184	0.000364	0
1	0.223234	0.024111	0.077241	0.060741	0.253940	0.503590	0.466350	0
2	0.704831	0.394256	0.155101	0.106348	0.133811	0.140230	0.078296	0
3	0.177642	0.124712	0.056273	0.017455	0.006439	0.003420	0.001974	0
4	0.213507	0.071184	0.017467	0.006277	0.003426	0.000703	0.000307	0

Рисунок 20 – Формат сформированной выборки.

Под тестовую выборку было решено выделить всего 20% объектов данных. Такой малый объем тестовых данных объясняется тем, что на этапе поиска оптимальных гиперпараметров модели планируется использовать кросс-валидацию [26] из 5 блоков, чтобы нейтрализовать эффект неравномерного распределения классов в тестовой выборке.

### 4.2 Метрики оценки качества моделей классификации

В сформированном датасете наблюдался дисбаланс классов - отношение количества объектов класса 1 (данный участок содержит в себе границу меж-фонемного перехода) к классу 0 равнялось 1 : 3. В виду чего применение такой

метрики оценки качества модели как *доля правильных ответов алгоритма* - ассигасу, будет давать неинформативные результаты.

В виду чего было решено использовать другие метрики, независящие в той или иной степени от баланса классов в выборке.

#### 4.2.1 precision и recall

Базовыми метриками для оценки качества модели классификации с несбалансированными классами являются *precision* и *recall*:

$$precision = \frac{TP}{TP + FP} \quad (23)$$

$$recall = \frac{TP}{TP + FN} \quad (24)$$

Значения  $TP, FP, FN$  получаются из таблицы сопряженности (таблица 4.2.1), строящейся на основе результатов классификации моделью и фактической принадлежности объектов к классам.

	$y = 1$	$y = 0$
$\hat{y} = 1$	True Positive (TP)	False Positive (FP)
$\hat{y} = 0$	False Negative (FN)	True Negative (TN)

Таблица 3 – Таблица сопряженности.

- TP (True Positive) - верно классифицированные положительные объекты;
- TN (True Negative) - верно классифицированные отрицательные объекты;
- FN (False Negative) - неверно классифицированные отрицательные объекты;
- FP (False Positive) - неверно классифицированные положительные объекты.

В данной работе за положительный объект принимается класс 1 (объект содержит в себе границу межфонемного перехода), а за отрицательный - класс 0.

Метрику *precision* можно интерпретировать как долю объектов, названных классификатором положительными и при этом действительно являющимися положительными, а *recall* показывает, какую долю объектов положитель-

ного класса из всех объектов положительного класса нашел алгоритм.

#### 4.2.2 Коэффициент корреляции Мэтьюса

При оптимизации гиперпараметров метода как правило используют одну метрику, улучшение результатов которой планируется увидеть на тестовой выборке.

Одной из наиболее популярных метрик является F-мера - среднее гармоническое *precision* и *recall*:

$$F_1 = \frac{precision * recall}{precision + recall} \quad (25)$$

Однако, данная метрика все еще может дать смещенную оценку предсказания, поскольку она не учитывает показатель *TN*.

Наиболее оптимальным вариантом является использование *Коэффициента Корреляции Мэтьюса (KKM)*, учитывающего все базовые показатели из таблицы сопряженности:

$$KKM = \frac{TP * TN - FP * FN}{\sqrt{(TP + FP) * (TP + FN) * (TN + FP) * (TN + FN)}} \quad (26)$$

где  $KKM \in [-1; +1]$ .

Коэффициент +1 представляет собой идеальное предсказание, 0 не лучше, чем случайное предсказание, а -1 указывает на полное несоответствие между предсказанием и наблюдением.

#### 4.2.3 AUC-ROC

При конвертации вероятности принадлежности объекта к классу в бинарную метку алгоритм классификации выбирает порог при котором класс 0 становится 1.

По умолчанию он равен 0.5, но такое решение далеко не всегда является оптимальным при отсутствии баланса классов в данных.

Наиболее популярным подходом оценить модель, не привязываясь к конкретному порогу, является *AUC-ROC*, где

- *ROC* (Receiver Operating Characteristic curve) - кривая ошибок, равная отношению доли истинно положительных объектов *TPR* (True Positive Rate) к доли ложных положительных объектов *FPR* (False Positive Rate);

- *AUC* (Area Under Curve) - площадь под графиком кривой ошибок *ROC*.

При этом *TPR* и *FPR* задаются следующими формулами:

$$TPR = \frac{TP}{TP + FN}, FPR = \frac{FP}{FP + TN} \quad (27)$$

Критерий *AUC-ROC* устойчив к несбалансированным классам, и интерпретируется как вероятность того, что случайно выбранный *positive* объект будет иметь более высокую вероятность идентифицироваться классификатором как *positive*, чем случайно выбранный *negative* объект.

### 4.3 Поиск оптимальных гиперпараметров метода

В качестве варьируемых гиперпараметров разработанного метода были выбраны следующие характеристики (разделы 2.3 и 2.5):

- *k* - количество отсчетов (замеров величины сигнала) в одном фрейме;
- *C* - параметр регуляризации модели классификации;
- *class\_weight* - параметр, позволяющий варьировать размер штрафа при неудачном выборе класса в модели классификации.

При этом *kernel* - ядро SVM, задающее функцию разделяющей гиперплоскости, принимает значение *rbf*, *gamma* - ширина ядра, принимает значение *scale* по умолчанию.

#### 4.3.1 GridSearchCV с применением ККМ

Подбор гиперпараметров *C* и *class\_weight* осуществлялся путем применения инструмента *GridSearchCV* из библиотеки *sklearn*, реализующий внутри себя кросс-валидацию. В качестве метрики качества был использован *Коэффициент Корреляции Мэтьюса* (раздел 4.2.2).

Интервал варьирования гиперпараметра *C* равнялся  $[10, 10^6]$ , а гиперпараметр *class\_weight* принимал значение *None* либо *balanced*.



Ниже приведена таблица усредненных на всех блоках результатов применения *GridSearchCV* для гиперпараметров  $C$  и  $class\_weight$ :

	$class\_weight = None$	$class\_weight = balanced$
$C = 10$	0.290	0.274
$C = 10^2$	0.302	0.293
$C = 10^3$	0.313	0.307
$C = 7 * 10^3$	0.318	0.318
$C = 8 * 10^3$	0.316	0.314
$C = 9 * 10^3$	0.318	0.317
$C = 10^4$	0.317	0.313
$C = 2 * 10^4$	0.320	0.307
$C = 3 * 10^4$	0.318	0.306
$C = 4 * 10^4$	0.318	0.303
$C = 10^5$	0.324	0.304
$C = 10^6$	0.307	0.293

Таблица 4 – Результаты *GridSearchCV*. Метрика качества - Коэффициент Корреляции Мэтьюса.

Из приведенной выше таблицы следует, что наиболее оптимальными значениями гиперпараметров  $C$  и  $class\_weight$  являются  $10^5$  и  $None$ , соответственно.

Подбор гиперпараметра  $k$  осуществлялся вне *GridSearchCV* из-за необходимости проведения слишком большого количества экспериментов.

Для исследования влияния параметра  $k$  на разработанный метод, было проведено 7 экспериментов обучения и тестирования метода при  $k = 4, 6, 8, 10, 12, 14, 16$ . Для проведения исследования были выбраны наиболее оптимальные значения гиперпараметров  $C$  и  $class\_weight$ , полученных в предыдущем шаге.

Ниже приведен график зависимости *Коэффициента Корреляции Мэтьюса ККМ* от параметра  $k$ .

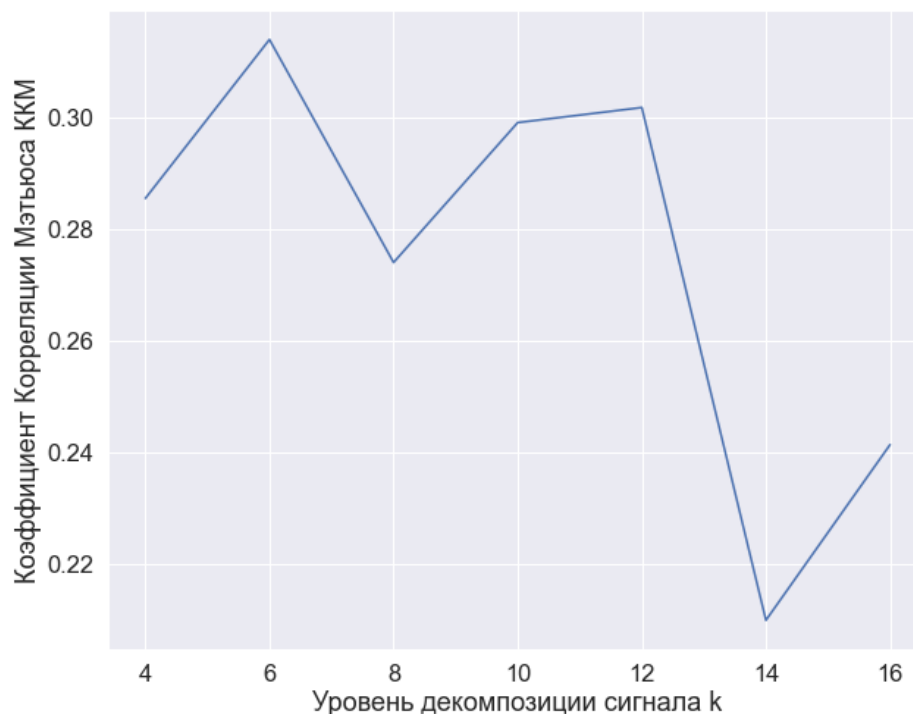


Рисунок 21 – График зависимости *Коэффициента Корреляции Мэтьюса ККМ* от уровня декомпозиции сигнала  $k$ .

Из приведенного выше графика видно, что наибольшее значение параметр *ММК* принимает при  $k = 6$ .

#### 4.4 Оценка качества модели классификации

Ниже приведены результаты применения описанных выше метрик на обученной модели классификации при выборе оптимальных параметров:

	precision	recall
positive	0.72	0.56
negative	0.63	0.77

Таблица 5 – Точность (precision) и полнота (recall) предсказания модели.

Коэффициент Корреляции Мэтьюса (ККМ) равен 0.33.

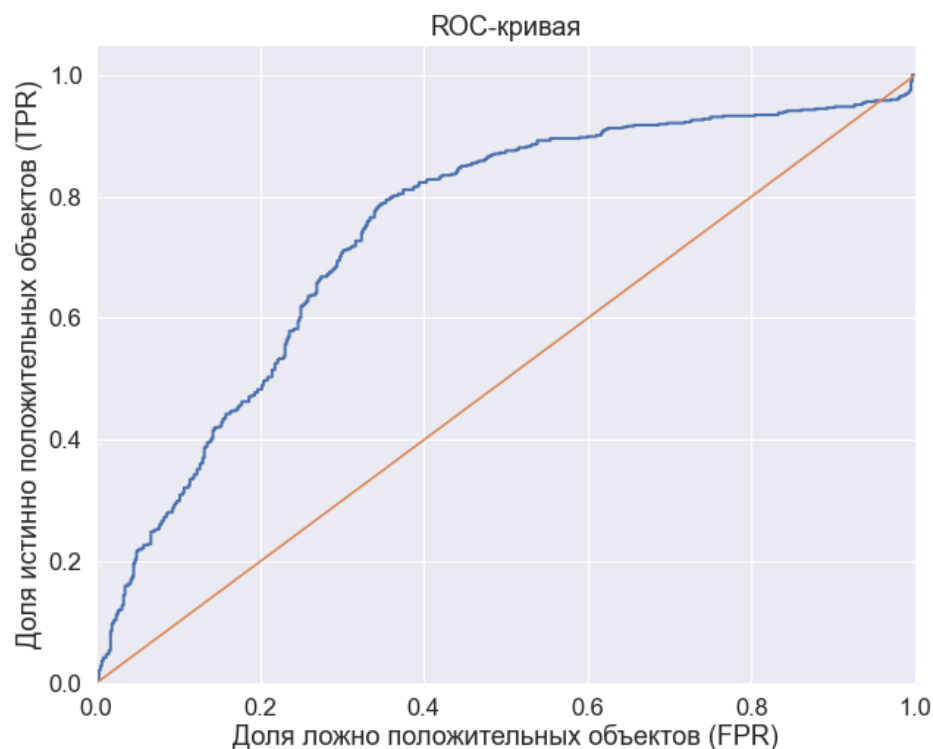


Рисунок 22 – ROC-кривая.

При этом  $AUC = 0.7325$ .

В литературе [27] приводится следующая экспертная шкала для значений  $AUC$ , по которой можно судить о качестве модели:

Интервал AUC	Качество модели
0.9-1.0	отличное
0.8-0.9	очень хорошее
0.7-0.8	хорошее
0.6-0.7	среднее
0.5-0.6	удовлетворительное

Таблица 6 – Экспертная шкала значений AUC.

Откуда следует, что полученная модель имеет хорошее качество предска-

зания.

#### **4.5 Выводы**

В результате проведения исследовательской работы над разработанным методом были выполнены следующие задачи:

- были сформированы обучающие и тестовые выборки данных для обучения модели классификации обнаружения межфонемных переходов;
- были рассмотрены основные метрики оценки качества модели классификации применительно к данной задаче, а именно, к задаче с несбалансированными классами;
- был произведен поиск оптимальных гиперпараметров метода при помощи Коэффициента корреляции Мэтьюса (ККМ);
- была произведена оценка качества модели на оптимальных гиперпараметрах при помощи рассмотренных метрик.

Таким образом, полученная модель имеет хорошее качество предсказания, что означает, что разработанный метод может быть применим для решения задачи сегментации речевого сигнала.

Такое относительно невысокое качество распознавания обусловлено малым размером обучающей выборки и сильной погрешностью измерений при разметке сформированных данных. Это связано с субъективностью человеческого слухового восприятия - все звуки размечались вручную разными людьми в результате чего возникла большая погрешность измерений.

## ЗАКЛЮЧЕНИЕ

В результате выполнения выпускной квалификационной работы были успешно достигнуты все поставленные задачи:

- проанализирована предметная область:
  - описаны основные этапы распознавания речевого сигнала;
  - обоснована актуальность задачи сегментации речевого сигнала;
  - проведен анализ существующих решений сегментации речевого сигнала.
- проведен анализ существующих методов выделения признаков речевого сигнала и алгоритмов классификации;
- выбраны методы наиболее оптимальные для поставленной цели;
- на их основе спроектирован метод выделения звуков естественного языка в звучащей речи и разработан алгоритм, реализующий данный метод;
- спроектирована система для проверки работоспособности метода;
- сформирована выборка для проведения эксперимента;
- проведено исследование разработанного метода на данной выборке;

Проведя исследование разработанного метода, выяснилось, что разработанный метод может быть применим для решения задачи сегментации речевого сигнала. Однако, реализованная модель имеет относительно невысокое качество предсказания, что обусловлено малым размером обучающей выборки и сильной погрешностью измерений при разметке сформированных данных.

Это связано с субъективностью человеческого слухового восприятия - все звуки размечались вручную разными людьми в результате чего возникла большая погрешность измерений.

Можно обозначить следующие пути развития данной работы:

- подготовка большего количества обучающей выборки, имеющей меньшую погрешность измерений для повышения точности метода;
- применение альтернативных алгоритмов классификации, позволяющими повысить качество распознавания с применением большего количества данных.

## СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. Аграновский А.В., Леднов Д.А., Телеснин Б.А. Сегментация речи (математическая модель) // Информационные технологии. 1998. № 9. С. 24-28
2. Казачкин, А. Е. Методы распознавания речи, современные речевые технологии [Электронный ресурс]. - Режим доступа: <https://moluch.ru/archive/277/62675/>, свободный - (17.03.2021)
3. Шпаков, Д. В. Распознавание голоса в сфере информационных технологий [Электронный ресурс]. - Режим доступа: <https://moluch.ru/archive/163/45163/>, свободный - (17.03.2021)
4. Федосин С.А., Еремин А. Ю. Классификация систем распознавания речи. [Электронный ресурс]. - Режим доступа: <http://fetmag.mrsu.ru/2010-2/pdf/SpeechRecognition.pdf>, свободный - (19.04.2021)
5. Тампель И.Б, Карпов А.А. Автоматическое распознавание речи. — СПб. : Университет ИТМО, 2016. — С. 113. [Электронный ресурс]. - Режим доступа: <https://books.ifmo.ru/file/pdf/1921.pdf>, свободный - (19.04.2021)
6. Рамишвили Г. С. Автоматическое опознавание говорящего по голосу. М. : Радио и связь, 1981. 224 с
7. О.А. Вишнякова, Д.Н. Лавров АВТОМАТИЧЕСКАЯ СЕГМЕНТАЦИЯ РЕЧЕВОГО СИГНАЛА НА БАЗЕ ДИСКРЕТНОГО ВЕЙВЛЕТ-ПРЕОБРАЗОВАНИЯ [Электронный ресурс]. - Режим доступа: <http://msm.omsu.ru/jrns/jrn23/VishnyakovaLavrov.pdf>, свободный - (22.04.2021)
8. Т. Ермоленко, В. Шевчук АЛГОРИТМЫ СЕГМЕНТАЦИИ С ПРИМЕНЕНИЕМ БЫСТРОГО ВЕЙВЛЕТ-ПРЕОБРАЗОВАНИЯ [Электронный ресурс]. - Режим доступа: <http://www.dialog-21.ru/media/2715/ermolenko.pdf>, свободный - (22.04.2021)
9. О.А. Вишнякова, Д.Н. Лавров ПРИМЕНЕНИЕ ПРЕОБРАЗОВАНИЯ ГИЛЬБЕРТА-ХУАНГА К ЗАДАЧЕ СЕГМЕНТАЦИИ РЕЧИ [Электронный ресурс]. - Режим доступа:

- [http://msm.omsu.ru/jrns/jrn24/Vishnakova\\_n24\\_2011.pdf](http://msm.omsu.ru/jrns/jrn24/Vishnakova_n24_2011.pdf), свободный - (10.05.2021)
10. И.В. Дашкевич. МАГИСТЕРСКАЯ ДИССЕРТАЦИЯ Модуль голосовой идентификации диктора [Электронный ресурс]. - Режим доступа: [http://elib.sfu-kras.ru/bitstream/handle/2311/74791/diplom\\_magistratura.pdf?sequence=1](http://elib.sfu-kras.ru/bitstream/handle/2311/74791/diplom_magistratura.pdf?sequence=1), свободный - (23.04.2021)
11. Музычук, Д. С. Использование преобразования Гильберта-Хуанга для формирования моделей фонем русского языка в задаче распознавания речи [Электронный ресурс]. - Режим доступа: <https://moluch.ru/archive/53/7041/>, свободный - (15.05.2021).
12. Давыдов А. В. Цифровая обработка сигналов: Тематические лекции. — Екатеринбург: УГГУ, ИГиГ, ГИН, Фонд электронных документов, 2005.
13. Томчук Кирилл Константинович, Сегментация речевых сигналов для задач автоматической обработки речи [Электронный ресурс]. - Режим доступа: [https://fs.guap.ru/disssov/tomchuk\\_kk/full.pdf](https://fs.guap.ru/disssov/tomchuk_kk/full.pdf), свободный - (10.05.2021)
14. Т.В. Батура, МЕТОДЫ АВТОМАТИЧЕСКОЙ КЛАССИФИКАЦИИ ТЕКСТОВ [Электронный ресурс]. - Режим доступа: [https://www.researchgate.net/publication/315328102\\_Metody\\_avtomaticheskoy\\_klassifikacii\\_tekstov](https://www.researchgate.net/publication/315328102_Metody_avtomaticheskoy_klassifikacii_tekstov), свободный - (11.05.2021)
15. К.В. Воронцов Лекции по методу опорных векторов [Электронный ресурс]. - Режим доступа: <http://www.ccas.ru/voron/download/SVM.pdf>, свободный - (11.05.2021)
16. Осциллографирование электрических сигналов цифровым осциллографом [Электронный ресурс]. - Режим доступа: [https://portal.tpu.ru/SHARED/k/KOZHEMYAK/Teaching/Tab1/Electronics/Lab\\_1\\_Digital\\_oscilloscope\\_2018.pdf](https://portal.tpu.ru/SHARED/k/KOZHEMYAK/Teaching/Tab1/Electronics/Lab_1_Digital_oscilloscope_2018.pdf), свободный - (10.05.2021)
17. Python documentation [Электронный ресурс]. - Режим доступа: <https://www.python.org/doc/>, свободный - (11.05.2021)

18. Pycharm documentation [Электронный ресурс]. - Режим доступа: <https://www.jetbrains.com/help/pycharm/inline-documentation.html>, свободный - (11.05.2021)
19. Librosa documentation [Электронный ресурс]. - Режим доступа: <https://librosa.org/doc/latest/index.html>, свободный - (11.05.2021)
20. PyWavelets documentation [Электронный ресурс]. - Режим доступа: [https://pywavelets.readthedocs.io/\\_/downloads/en/v1.0.0/pdf/](https://pywavelets.readthedocs.io/_/downloads/en/v1.0.0/pdf/), свободный - (11.05.2021)
21. Scikit-learn documentation [Электронный ресурс]. - Режим доступа: <https://www.sklearn.org/documentation.html>, свободный - (11.05.2021)
22. Pandas documentation [Электронный ресурс]. - Режим доступа: <https://pandas.pydata.org/pandas-docs/stable/index.html>, свободный - (11.05.2021)
23. Matplotlib documentation [Электронный ресурс]. - Режим доступа: <https://matplotlib.org/>, свободный - (11.05.2021)
24. NumPy documentation [Электронный ресурс]. - Режим доступа: <https://numpy.org/doc/>, свободный - (11.05.2021)
25. PyQt5 documentation [Электронный ресурс]. - Режим доступа: <https://doc.qt.io/qtforpython/>, свободный - (11.05.2021)
26. Кросс-валидация [Электронный ресурс]. - Режим доступа: <http://datascientist.one/cross-validation/>, свободный - (06.04.2021)
27. Zweig, Mark H.; Campbell, Gregory. Receiver-operating characteristic (ROC) plots: a fundamental evaluation tool in clinical medicine. Clinical Chemistry. Нью-Йорк: Изд-во Oxford, 577 с.



## ПРИЛОЖЕНИЕ А

Ниже приведены листинги разработанного ПО.

### MainView

```
1 from PyQt5 import uic
2 from PyQt5.QtCore import pyqtSlot, Qt
3 from PyQt5.QtWidgets import QMessageBox, QMainWindow, QFileDialog
4
5
6 class MainWindow(QMainWindow):
7     def __init__(self, in_controller):
8         super(MainWindow, self).__init__()
9         self.ui = uic.loadUi("view/main_view/main_window.ui", self)
10        self.controller = in_controller
11        self.ui.get_result_button.setEnabled(False)
12        self.ui.clear_model_button.setEnabled(False)
13        self.ui.save_model_button.setEnabled(False)
14
15        @pyqtSlot(name='on_choose_model_button_clicked')
16        def choose_model(self):
17            filename, _ = QFileDialog.getOpenFileName(self, Выберите' модель для
загрузки',
18                                                    './Data/models/', 'Model
File(*.pkl)')
19            if len(filename) > 0:
20                status = self.controller.choose_model(filename)
21                if status == Модель' успешно загружена':
22                    self.ui.teach_model_button.setEnabled(False)
23                    self.ui.choose_model_button.setEnabled(False)
24                    self.ui.clear_model_button.setEnabled(True)
25                    self.ui.choose_model_label.setTextМодель(' выбрана')
26                    self.try_enable_get_result_btn()
27                else:
28                    msgBox = QMessageBox()
29                    msgBox.setStandardButtons(QMessageBox.Ok)
30                    msgBox.setIcon(QMessageBox.Critical)
31                    msgBox.setText(status)
32                    msgBox.exec()
33
34        @pyqtSlot(name='on_teach_model_button_clicked')
35        def teach_model(self):
36            self.controller.launch_teach_model_view()
37            self.ui.clear_model_button.setEnabled(True)
38            self.ui.save_model_button.setEnabled(True)
```

## MainView

```
1  Выберите модель для загрузки Модель успешно загружена Модель выбрана
2
3  @pyqtSlot(name='on_save_model_button_clicked')
4  def save_model(self):
5      filename, _ = QFileDialog.getSaveFileName(self, Сохранить ' модель',
6                                              './Data/models/', 'Model
7  File (*.pkl)')
8      if len(filename) > 0:
9          status = self.controller.save_model(filename)
10         msgBox = QMessageBox()
11         msgBox.setStandardButtons(QMessageBox.Ok)
12         if status == Модель ' успешно сохранена':
13             msgBox.setIcon(QMessageBox.Information)
14             msgBox.setText(status)
15             self.ui.save_model_button.setEnabled(False)
16         else:
17             msgBox.setIcon(QMessageBox.Critical)
18             msgBox.setText(status)
19             msgBox.exec()
20
21 @pyqtSlot(name='on_clear_model_button_clicked')
22 def clear_model(self):
23     self.controller.clear_model()
24     self.ui.clear_model_button.setEnabled(False)
25     self.ui.save_model_button.setEnabled(False)
26     self.ui.get_result_button.setEnabled(False)
27     self.ui.teach_model_button.setEnabled(True)
28     self.ui.choose_model_button.setEnabled(True)
29
30 @pyqtSlot(name='on_choose_test_audio_button_clicked')
31 def choose_test_audio(self):
32     filename, _ = QFileDialog.getOpenFileName(self, Выберите ' mp3 файл-
33     для тестирования',
34                                             './Data/test/', 'Audio
35     File (*.mp3)')
36     if len(filename) > 0:
37         status = self.controller.load_test_signal(filename)
38         self.ui.choose_test_audio_label.setText(status)
39         self.try_enable_get_result_btn()
40
41 @pyqtSlot(name='on_get_result_button_clicked')
42 def get_result(self):
43     self.controller.launch_result_view()
```

## ResultView

```
1 from PyQt5 import uic
2 from PyQt5.QtCore import pyqtSlot, Qt
3 from PyQt5.QtGui import QPixmap
4 from PyQt5.QtWidgets import QMessageBox, QWidget
5
6
7 class ResultWindow(QWidget):
8     def __init__(self, in_controller):
9         super(ResultWindow, self).__init__()
10        self.ui = uic.loadUi("view/result_view/result_window.ui", self)
11        self.controller = in_controller
12
13        @pyqtSlot(name='on_close_button_clicked')
14        def close_btn(self):
15            self.close()
16
17        def paint_graph(self):
18            status, result_graph_path = self.controller.extract_sounds()
19            if len(result_graph_path) > 0:
20                self.ui.graph_label.setStyleSheet("QLabel{{border-image:
url({});}}")
21
22
23                .format(result_graph_path))
24            else:
25                msgBox = QMessageBox()
26                msgBox.setStandardButtons(QMessageBox.Ok)
27                msgBox.setIcon(QMessageBox.Critical)
28                msgBox.setText(status)
29                msgBox.exec()
```

## TeachModelView

```
1 from PyQt5 import uic
2 from PyQt5.QtCore import pyqtSlot
3 from PyQt5.QtWidgets import QMessageBox, QDialog, QFileDialog
4
5
6 class TeachModelWindow(QDialog):
7     def __init__(self, in_controller):
8         super(TeachModelWindow, self).__init__()
9         self.ui = uic.loadUi("view/teach_model_view/teach_model_window.ui",
10 self)
11         self.controller = in_controller
12         self.ui.fit_model_button.setEnabled(False)
13         self.ui.choose_fit_audio_button.setEnabled(False)
14
15         @pyqtSlot(name='on_close_button_clicked')
16         def close_btn(self):
17             self.controller.window_closed()
18             self.close()
19
20         @pyqtSlot(name='on_choose_fit_audio_button_clicked')
21         def choose_fit_audio(self):
22             filenames, _ = QFileDialog.getOpenFileNames(self, 'Choose mp3-files
23 for learning',
24                                                         './Data/train/', 'Audio
25 Files (*.mp3)')
26             if len(filenames) > 0:
27                 status = self.controller.load_train_signals(filenames)
28                 self.ui.choose_fit_audio_label.setText(status)
29
30                 markups_status = self.controller.get_status('train_markups')
31                 if status == 'Audio-files have chosen' and markups_status ==
32 'Markups have chosen':
33                     model_status = self.controller.get_status('model')
34                     if model_status == 'Model havent chosen':
35                         self.ui.fit_model_button.setEnabled(True)
36
37         @pyqtSlot(name='on_fit_model_button_clicked')
38         def fit_model_btn(self):
39             status = self.controller.fit_model()
40             msgBox = QMessageBox()
41             msgBox.setStandardButtons(QMessageBox.Ok)
42             if status == 'Model have chosen':
43                 msgBox.setIcon(QMessageBox.Information)
44                 msgBox.setText(status)
```

## TeachModelView

```
1         self.ui.fit_model_label.setText(status)
2     else:
3         msgBox.setIcon(QMessageBox.Critical)
4         msgBox.setText(status)
5         msgBox.exec()
6
7     @pyqtSlot(name='on_choose_fit_markup_button_clicked')
8     def choose_fit_markup(self):
9         filename, _ = QFileDialog.getOpenFileName(self, 'Choose csv-markups
for learning',
10                                                     './Data/train/', 'CSV
File(*.csv)')
11         if len(filename) > 0:
12             status = self.controller.load_train_markup(filename)
13             self.ui.choose_fit_markup_label.setText(status)
14
15             if status == 'Markups have chosen':
16                 self.ui.choose_fit_audio_button.setEnabled(True)
17                 audios_status = self.controller.get_status('train_signals')
18                 if audios_status == 'Audio-files have chosen':
19                     model_status = self.controller.get_status('model')
20                     if model_status == 'Model havent chosen':
21                         self.ui.fit_model_button.setEnabled(True)
```

## MainController

```
1
2 from view.main_view.main_window import MainWindow
3 from view.result_view.result_window import ResultWindow
4 from view.teach_model_view.teach_model_window import TeachModelWindow
5
6 from ..result_controller.result_controller import ResultController
7 from ..teach_model_controller.teach_model_controller import
    TeachModelController
8
9 from model.sounds_extraction_module.sounds_extraction_module import
    SoundsExtractionModule
```

## MainController

```
1
2 class MainController:
3     def __init__(self):
4         self.main_window = MainWindow(self)
5         self.main_window.show()
6         self.sounds_extraction_module = SoundsExtractionModule()
7
8         # инициализация teach_model_controller
9         self.teach_model_controller = TeachModelController(self,
10 self.sounds_extraction_module)
11
12         # инициализация result_controller
13         self.result_controller = ResultController(self,
14 self.sounds_extraction_module)
15
16     def launch_teach_model_view(self):
17         self.teach_model_controller.show_window()
18
19     def teach_model_window_closed(self):
20         self.main_window.try_enable_get_result_btn()
21
22     def launch_result_view(self):
23         self.result_controller.show_window()
24
25     def clear_model(self):
26         self.sounds_extraction_module.clear_classification_model()
27
28     def load_test_signal(self, filename):
29         return self.sounds_extraction_module.load_test_signal(filename)
30
31     def get_status(self, data):
32         if data == 'test_signal':
33             return self.sounds_extraction_module.get_status_test_signal()
34         elif data == 'model':
35             return
36         self.sounds_extraction_module.get_status_classification_module()
37
38     def choose_model(self, path):
39         return self.sounds_extraction_module.choose_model(path)
40
41     def save_model(self, path):
42         return self.sounds_extraction_module.save_classification_model(path)
```

## ResultController

```
1
2 from view.result_view.result_window import ResultWindow
3
4
5 class ResultController:
6     def __init__(self, in_main_controller, in_sounds_extraction_module):
7         self.main_controller = in_main_controller
8         self.sounds_extraction_module = in_sounds_extraction_module
9         self.window = ResultWindow(self)
10
11     def show_window(self):
12         self.window.show()
13         self.window.paint_graph()
14
15     def extract_sounds(self):
16         return self.sounds_extraction_module.predict_classification_model()
```

## TeachModelController

```
1
2 from view.teach_model_view.teach_model_window import TeachModelWindow
3
4 class TeachModelController:
5     def __init__(self, in_main_controller, in_sounds_extraction_module):
6         self.main_controller = in_main_controller
7         self.window = TeachModelWindow(self)
8         self.sounds_extraction_module = in_sounds_extraction_module
9
10    def show_window(self):
11        self.window.show()
12
13    def window_closed(self):
14        self.main_controller.teach_model_window_closed()
15
16    def load_train_signals(self, filenames):
17        return self.sounds_extraction_module.load_train_signals(filenames)
18
19    def load_train_markup(self, filename):
20        return self.sounds_extraction_module.load_train_markup(filename)
21
22    def get_status(self, data):
23        if data == 'train_signals':
24            return self.sounds_extraction_module.get_status_train_signals()
```

## ClassificationModule

```
1
2 from model.data import Data
3 from model.params.classification_params import ClassificationParams
4
5 class ClassificationModule:
6     def __init__(self, in_data_preprocessing_module,
7         in_feature_extraction_module):
8         self.data_preprocessing_module = in_data_preprocessing_module
9         self.feature_extraction_module = in_feature_extraction_module
10        self.params = ClassificationParams()
11        self.model = DataМодель(' не выбрана')
12
13    def clear_model(self):
14        self.model = DataМодель(' не выбрана')
15
16    def fit(self, signals, markups):
17        if self.model.status == Модель' обучена':
18            return self.model.status
19        try:
20            features =
21            self.feature_extraction_module.extract_from_signals(signals)
22
23            train_data =
24            self.data_preprocessing_module.create_train_data(features, markups)
25
26            train_data = train_data.sample(frac=1)
27
28            y_train = train_data['y'].astype('int64')
29            X_train = train_data.drop(axis=0, columns=['y'])
30
31            self.params.model.fit(X_train, y_train)
32
33            self.model.data = self.params.model
34            self.model.status = Модель' обучена'
35        except:
36            return Ошибка' обучения модели'
37        return self.model.status
38
39    def predict(self, signal):
40        if self.model.status != 'Model have teached':
```



## ClassificationModule

```
1 МодельневыбранаМодельневыбранаМодельобученаМодельобученаОшибкаобучениямодели
2         return self.model.status
3     try:
4         features =
5         self.feature_extraction_module.extract_from_signal(signal)
6
7         test_data =
8         self.data_preprocessing_module.create_test_data(features)
9
10        test_data = test_data.sample(frac=1)
11        predictions = self.model.data.predict(test_data)
12        return 'Sounds have got',
13        self._convert_predictions_to_time_list(predictions)
14    except:
15        return 'Error! Cant get sounds.', None
16
17    def _convert_predictions_to_time_list(self, predictions):
18
19        time_list =
20        self.feature_extraction_module.get_t_between_frames(predictions, 0.)
21
22        res = []
23        dt = time_list[1] - time_list[0]
24        for t, prediction in zip(time_list, predictions):
25            if prediction:
26                res.append(t + dt)
27        return res
28
29    def load_model(self, model):
30        self.model.data = model
31        self.model.status = 'Model have chosen'
32        return 'Model loaded'
```

## DataPreprocessingModule

```
1 import numpy as np
2 from itertools import chain
3 import pandas as pd
4 import ast
5
6 from model.feature_extraction_module.feature_extraction_module import
    FeatureExtractionModule
7 from model.params.markups_params import MarkupsParams
8 from model.params.audio_params import AudioParams
9
10 class DataPreprocessingModule:
11
12
13     def create_test_data(self, features):
14         features_names =
            FeatureExtractionModule.get_wavelet_columns_names(features[0])
15         res_data = []
16         #for feature in features:
17             #    res_data.append(list(chain(*feature)))
18         #return pd.DataFrame(res_data, columns=features_names)
19         return pd.DataFrame(features, columns=features_names)
20
21     def create_train_data(self, features, df_markups):
22
23         features_names =
            FeatureExtractionModule.get_wavelet_columns_names(features[0][1][0])
24         features_names.append('y')
25         res_data = []
26         for path, frames_coeffs in features:
27             filename = self._get_filename_from_path(path)
28             for interval in
                ast.literal_eval(df_markups[df_markups[MarkupsParams.filename()] ==
                    filename]
29                                     .intervals.iloc[0]):
30
31                 t_positions_between_frames =
                    FeatureExtractionModule.get_t_between_frames(frames_coeffs, interval[0])
32
33                 y = np.zeros(frames_coeffs.shape[0])
34                 borders = self._borders_preprocessing(
35                     df_markups[df_markups[MarkupsParams.filename()] ==
36                         filename],
37                         MarkupsParams.start_end()[0],
38                         MarkupsParams.start_end()[1])
```

## DataPreprocessingModule

```
1
2         borders = self._filter_borders(borders, interval)
3         for border in borders:
4             for t_pos in range(len(t_positions_between_frames)):
5                 if t_positions_between_frames[t_pos] > border:
6                     y[t_pos] = 1.
7                     break
8
9         res_data_i = []
10        for frame, y_i in zip(frames_coeffs, y):
11            #lst_i = list(chain(*frame))
12            lst_i = frame
13            lst_i.append(y_i)
14            res_data_i.append(lst_i)
15        res_data.extend(res_data_i)
16        return pd.DataFrame(res_data, columns=features_names)
17
18
19    def _filter_borders(self, borders, interval):
20        res = []
21        for border in borders:
22            if border >= interval[0] and border <= interval[1]:
23                res.append(border)
24        return res
25
26    def _get_filename_from_path(self, path):
27        pos = path.rfind('/')
28        if pos != -1:
29            return path[pos + 1:]
30        return path
31
32    def _borders_preprocessing(self, df, column_name_1, column_name_2):
33
34        borders = np.sort(np.array(self._merge_2_columns_on_row(df,
35        column_name_1, column_name_2)))
36
37        return self._check_on_min_phonem_interval(borders,
38        AudioParams.min_phonem_interval())
39
40    def _merge_2_columns_on_row(self, df, column_name_1, column_name_2):
41        lst = []
```

## FeatureExtractionModule

```
1 import numpy as np
2 import librosa
3 from scipy import signal
4 from pywt import wavedec
5
6 from model.params.audio_params import AudioParams
7
8 class FeatureExtractionModule:
9
10     @staticmethod
11     def get_wavelet_columns_names(features):
12         '''columns = []
13         n = len(features) - 1
14         for pos, A_coeffs in enumerate(features[0]):
15             columns.append('A_' + str(n) + str(pos))
16         for i, D_coeffs in enumerate(features[1:]):
17             for j, D_i_coeffs in enumerate(D_coeffs):
18                 columns.append('D_' + str(n - i) + '_' + str(j))
19         return columns'''
20         #return ['A', 'D6', 'D5', 'D4', 'D3', 'D2', 'D1']
21         res = ['D' + str(i) for i in range(len(features) - 1)]
22         res.insert(0, 'A')
23         return res
24         #return ['D' + str(i) for i in range(len(features))]
25
26     @staticmethod
27     def get_t_between_frames(frames_coeffs, offset):
28         t_values = []
29         cur_time = offset
30         dt = AudioParams.frame_sz() * AudioParams.hop_part() /
AudioParams.sr()
31         for i in range(len(frames_coeffs)):
32             t_values.append(cur_time)
33             cur_time += dt
34         return np.array(t_values)
35
36     @staticmethod
37     def get_t_between_pos(sig, offset):
38         t_values = []
39         cur_time = offset
40         dt = 1 / AudioParams.sr()
41         for i in range(len(sig)):
42             t_values.append(cur_time)
43             cur_time += dt
44         return np.array(t_values)
```

## FeatureExtractionModule

```
1         return np.array(t_values)
2
3         # получает массив сигналов: [[path, signal],...], возвращает словарь
признаков: {path: feature}
4     def extract_from_signals(self, sigs):
5         res = []
6         for path_sig in sigs:
7             res.append([path_sig[0], self.extract_from_signal(path_sig[1])])
8         return res
9
10        # получает на вход signal, возвращает признаки signal
11    def extract_from_signal(self, sig):
12        # разбиение сигнала на фреймы
13        frames_sig = self._splitting_signal_on_frames(sig,
AudioParams.frame_sz(), AudioParams.hop_part())
14
15        # устранение деффекта на краях с помощью оконной функции Хамминга
16        frames_hamming_sig = self._hamming_func(frames_sig,
AudioParams.frame_sz())
17
18        # применение вейвлетпреобразования-
19        frames_coeffs = self._frames_wavelet(frames_hamming_sig,
AudioParams.wavelet(), AudioParams.wavelet_level())
20
21        return frames_coeffs
22
23    def _splitting_signal_on_frames(self, desccrete_sig, frame_sz, hop_part):
24        frames = librosa.util.frame(desccrete_sig, frame_length=frame_sz,
hop_length=int(hop_part * frame_sz), axis=0)
25        return frames
26
27    def _hamming_func(self, frames_signal, frame_sz):
28        w = signal.windows.hamming(frame_sz, sym=True)
29        return np.array([w * frame for frame in frames_signal])
30
31    def _wavelet_conv(self, sig, method, level):
32        levels = np.array(wavedec(sig, method, level=level))
33        return [np.square(level).sum() for level in levels]
34
35    def _frames_wavelet(self, frames_signal, method, level):
36        return np.array([self._wavelet_conv(sig, method, level) for sig in
frames_signal], dtype=object)
```

## GraphModule

```
1 import matplotlib.pyplot as plt
2
3 class GraphModule:
4     # получает список временных промежутков и дискретный сигнал. Возвращает
    график
5     def create_predicted_graph(self, predicted_markup, test_signal, t_lst):
6         fig = plt.figure(figsize=(20, 7))
7         ax = plt.plot(t_lst, test_signal)
8         plt.vlines(predicted_markup, test_signal.min(), test_signal.max(),
9                     color='r')
10        plt.xlabel('Время', c='r')
11        plt.ylabel('Амплитуда')
12        #plt.legend([name])
13        #plt.title(name)
14        return fig
```

## LoadDataModule

```
1 import pandas as pd
2 import librosa
3 from itertools import chain
4 import ast
5 import matplotlib.pyplot as plt
6 import pickle
7
8 from model.data import Data
9 from model.params.audio_params import AudioParams
10 from model.params.markups_params import MarkupsParams
11
12
13 class LoadDataModule:
14     def __init__(self):
15         self.test_signal = Data('Audio-file havent chosen')
16         self.test_signal_filename = 'undefined_filename'
17         self.train_signals = Data('Audio-files havent chosen')
18         self.train_markup = Data('Markups havent chosen')
19
20
21     def load_train_signals(self, paths):
22         res = []
23         for path in paths:
24             df = self.train_markup.data
25             filename = self._get_filename_from_path(path)
```

## LoadDataModule

```
1         intervals = ast.literal_eval(df[df.recordname ==
2 filename].intervals.iloc[0])
3         for interval in intervals:
4             res.append([path, librosa.load(path, sr=AudioParams.sr(),
5 offset=interval[0],
6 duration=interval[1]-interval[0])[0]])
7         self.train_signals.data = res
8         self.train_signals.status = 'Audio-files have chosen'
9         return self.train_signals.status
10
11 def _get_filename_from_path(self, path):
12     pos = path.rfind('/')
13     if pos != -1:
14         return path[pos + 1:]
15     return path
16
17 def load_train_markup(self, path):
18     df = pd.read_csv(path)
19     self.train_markup.data = df
20     self.train_markup.status = 'Markups have chosen'
21     return self.train_markup.status
22
23 def load_test_signal(self, path):
24     self.test_signal_filename = self._get_filename_from_path(path)[-4:]
25     self.test_signal.data, _ = librosa.load(path, sr=AudioParams.sr(),
26 offset=0.)
27     self.test_signal.status = 'Audio-file have chosen'
28     return self.test_signal.status
29
30 def save_predicted_markup(self, predicted_markup):
31     df = pd.DataFrame({'y': predicted_markup})
32     df.to_csv('./Data/result_markup/'+self.test_signal_filename+'.csv')
33
34
35 def save_predicted_graph(self, graph):
36     graph_path = './Data/result_photo/'+self.test_signal_filename+'.png'
37     plt.savefig(graph_path)
38     return graph_path
```

## LoadDataModule

```
1
2     def clear_model_data(self):
3         self.train_signals = Data('Audio-files havent chosen')
4         self.train_markup = Data('Markups havent chosen')
5
6     def load_model(self, path):
7         with open(path, 'rb') as file:
8             return pickle.load(file)
9
10    def save_model(self, model, path):
11        with open(path, 'wb') as file:
12            pickle.dump(model, file)
13        return 'Model saved'
```

## SoundsExtractionModule

```
1
2 from model.load_data_module.load_data_module import LoadDataModule
3 from model.classification_module.classification_module import
4     ClassificationModule
5 from model.feature_extraction_module.feature_extraction_module import
6     FeatureExtractionModule
7 from model.data_preprocessing_module.data_preprocessing_module import
8     DataPreprocessingModule
9 from model.graph_module.graph_module import GraphModule
10
11 class SoundsExtractionModule:
12     def __init__(self):
13         self.load_data_module = LoadDataModule()
14         self.data_preprocessing_module = DataPreprocessingModule()
15         self.feature_extraction_module = FeatureExtractionModule()
16         self.graph_module = GraphModule()
17         self.classification_module =
18         ClassificationModule(self.data_preprocessing_module,
19
20         self.feature_extraction_module)
21
22     def get_status_classification_module(self):
23         return self.classification_module.model.status
24
25     def get_status_test_signal(self):
```



## SoundsExtractionModule

```
1         return self.load_data_module.test_signal.status
2
3
4     def get_status_train_signals(self):
5         return self.load_data_module.train_signals.status
6
7
8     def get_status_train_markup(self):
9         return self.load_data_module.train_markup.status
10
11
12     def clear_classification_model(self):
13         self.load_data_module.clear_model_data()
14         self.classification_module.clear_model()
15
16
17     def fit_classification_model(self):
18         if self.get_status_classification_module() != 'Model havent chosen':
19             return 'Model chosen'
20         train_signals_status = self.get_status_train_signals()
21         if train_signals_status == 'Audio-files havent chosen':
22             return train_signals_status
23         train_markup_status = self.get_status_train_markup()
24         if train_markup_status == 'Markups havent chosen':
25             return train_markup_status
26         status =
27         self.classification_module.fit(self.load_data_module.train_signals.data,
28
29         self.load_data_module.train_markup.data)
30         return status
31
32
33     def predict_classification_model(self):
34         if self.get_status_classification_module() != 'Model teached':
35             return 'Model didnt chosen', ''
36         test_signal_status = self.get_status_test_signal()
37         if test_signal_status == 'Audio-files havent chosen':
38             return test_signal_status, ''
39         # возвращает статус и список predicted_markup временных промежутков
40         status, predicted_markup =
41         self.classification_module.predict(self.load_data_module.test_signal.data)
42         if predicted_markup:
43             self.load_data_module.save_predicted_markup(predicted_markup)
44         graph =
45         self.graph_module.create_predicted_graph(predicted_markup,
46         self.load_data_module.test_signal.data,
```

## SoundsExtractionModule

```
1  возвращает статус и список временных промежутков
2
3  self.feature_extraction_module.get_t_between_pos(
4
5  self.load_data_module.test_signal.data, 0.))
6      graph_path = self.load_data_module.save_predicted_graph(graph)
7      return status, graph_path
8
9  else:
10     return status, ''
11
12 def save_classification_model(self, path):
13     classification_status = self.get_status_classification_module()
14     if classification_status != 'Model teached':
15         return 'Model havent chosen'
16     return
17 self.load_data_module.save_model(self.classification_module.model.data,
18 path)
19
20
21 def load_train_signals(self, filenames):
22     status = self.get_status_train_signals()
23     if status == 'Audio-files have chosen':
24         return status
25     return self.load_data_module.load_train_signals(filenames)
26
27
28 def load_train_markup(self, filename):
29     status = self.get_status_train_markup()
30     if status == 'Markups have chosen':
31         return status
32     return self.load_data_module.load_train_markup(filename)
33
34
35 def load_test_signal(self, filename):
36     return self.load_data_module.load_test_signal(filename)
37
38
39 def choose_model(self, path):
40     try:
41         model = self.load_data_module.load_model(path)
42         if model:
43             return self.classification_module.load_model(model)
```

## Parameters

```
1
2 class Data:
3     def __init__(self, status, data=None):
4         self.data = data
5         self.status = status
```

```
1
2 class AudioParams:
3
4     @staticmethod
5     def sr():
6         return 16000
7
8     @staticmethod
9     def frame_sz():
10        return 2500
11
12    @staticmethod
13    def hop_part():
14        return 0.6
15
16    @staticmethod
17    def wavelet():
18        return 'db4'
19
20    @staticmethod
21    def wavelet_level():
22        return 6
23
24    @staticmethod
25    def min_phonem_interval():
26        return 0.025
```

```
1 from sklearn.svm import SVC
2
3 class ClassificationParams:
4     def __init__(self):
5         self.model = SVC(kernel='rbf')
```

## Parameters

```
1
2 class MarkupsParams:
3
4     @staticmethod
5     def filename():
6         return 'recordname'
7
8     @staticmethod
9     def intervals():
10        return 'intervals'
11
12    @staticmethod
13    def start_end():
14        return ['start', 'end']
```