

**Нурматова Елена Вячеславовна,**  
*доцент кафедры аппаратного, программного и  
математического обеспечения вычислительных систем,  
кандидат технических наук,  
Института кибербезопасности и цифровых технологий  
РТУ МИРЭА  
nurmatova@mirea.ru*

**Пылаев Константин Витальевич,**  
*магистрант 2 курса  
специальности «Информационные системы и технологии»  
Института кибербезопасности и цифровых технологий  
РТУ МИРЭА*

## **ПОДХОД К РЕШЕНИЮ ЗАДАЧИ ДИАГНОСТИКИ СЕРДЕЧНО-СОСУДИСТЫХ ЗАБОЛЕВАНИЙ**

***Аннотация.** В статье описан выбор наилучшей модели для прогнозирования сердечно-сосудистых заболеваний и этапы работы, которую следует провести с данными перед их обучением.*

***Ключевые слова:** машинное обучение, статистические данные, прогнозирование, анализ данных.*

При рассмотрении темы, связанной с диагностикой сердечно-сосудистых заболеваний (ССЗ), возникает вопрос, как спрогнозировать заболевание? Ответом будет являться машинное обучение<sup>1</sup>. Алгоритм самостоятельного нахождения решений путем комплексного использования статистических данных, из которых выводятся закономерности и на основе которых будут производиться прогнозы.

Модели машинного обучения часто используются для прогнозирования многих распространённых заболеваний, включая диагностику ССЗ<sup>2</sup>. Часто существует множество факторов, которые способствуют выявлению пациентов, подверженных риску из этой группы заболеваний. Методы машинного обучения помогают выявлять скрытые закономерности в этих факторах, которые в противном случае могут быть попросту пропущены человеком<sup>3</sup>.

---

<sup>1</sup> Хабр // URL: <https://habr.com/ru/post/448892/> (дата обращения: 15.11.2022).

<sup>2</sup> Medium // URL: <https://medium.com/nuances-of-programming/все-модели-машинного-обучения-за-5-минут-9270566197e7> (дата обращения: 15.11.2022).

<sup>3</sup> BMC // URL: <https://bmcmmedinformdecismak.biomedcentral.com/articles/10.1186/s12911-019-0918-5> (дата обращения: 15.11.2022).

Существует достаточно моделей, которые могли помочь с решением поставленной задачи. Их характеристики приведены в табл. 1.

Получается, что каждая из моделей имеет, как и существенные достоинства в сравнении с другими, так и значительные недостатки. При этом нужно понимать, что у моделей как положительных, так и отрицательных сторон намного больше нежели, чем представлено в таблице. Для более детального выяснения и выбора наилучшего метода необходимо провести ряд тестов и ознакомиться с каждым поглубже<sup>4</sup>.

Таблица 1. Сравнение моделей машинного обучения

| <i>Название модели</i>                             | <i>Плюсы</i>   | <i>Минусы</i>   |
|--|--|---|
| 1. Logistic Regression (Логистическая регрессия)   | <ol style="list-style-type: none"> <li>1. Проще реализовать, интерпретировать и очень эффективно обучать.</li> <li>2. Очень быстро классифицирует неизвестные записи.</li> <li>3. Может интерпретировать коэффициенты модели как показатели важности функции.</li> </ol> | <ol style="list-style-type: none"> <li>1. Если количество наблюдений меньше количества признаков, логистическую регрессию не следует использовать, может привести к переобучению.</li> <li>2. Основным ограничением является предположение о линейности между зависимой переменной и независимыми переменными.</li> </ol> |
| 2. K Nearest Neighbours (k-ближайшие соседи)       | <ol style="list-style-type: none"> <li>1. Высокая точность.</li> <li>2. Нечувствительность к выбросам.</li> <li>3. Отсутствие допущений о вводе данных.</li> </ol>   | <ol style="list-style-type: none"> <li>1. Высокая временная сложность</li> <li>2. Большая пространственная сложность.</li> </ol>  |
| 3. Kernel SVM (Метод опорных векторов)             | <ol style="list-style-type: none"> <li>1. Эффективен в пространствах больших размеров.</li> <li>2. Универсальность.</li> </ol>   | <ol style="list-style-type: none"> <li>1. В случае превышения количество выборок, следует избегать чрезмерной подгонки при выборе функций ядра.</li> </ol>  |
| 4. Naïve Bayes (Наивный байесовский классификатор) | <ol style="list-style-type: none"> <li>1. Простота реализации</li> <li>2. Низкие вычислительные затраты при обучении и классификации.</li> </ol>   | <ol style="list-style-type: none"> <li>1. Низкое качество классификации в большинстве реальных задач.</li> </ol>  |
| 5. Decision Tree (Дерева решений)                  | <ol style="list-style-type: none"> <li>1. Простота.</li> <li>2. Не нужна нормализация.</li> <li>3. Встроенный отбор признаков работает одновременно с дискретными и непрерывными признаками.</li> </ol>  | <ol style="list-style-type: none"> <li>1. Для новых наблюдений требуется полная перестройка всего дерева.</li> </ol>  |

<sup>4</sup> BMC // URL: <https://bmcmmedinformdecismak.biomedcentral.com/articles/10.1186/s12911-019-0918-5> (дата обращения: 15.11.2022).

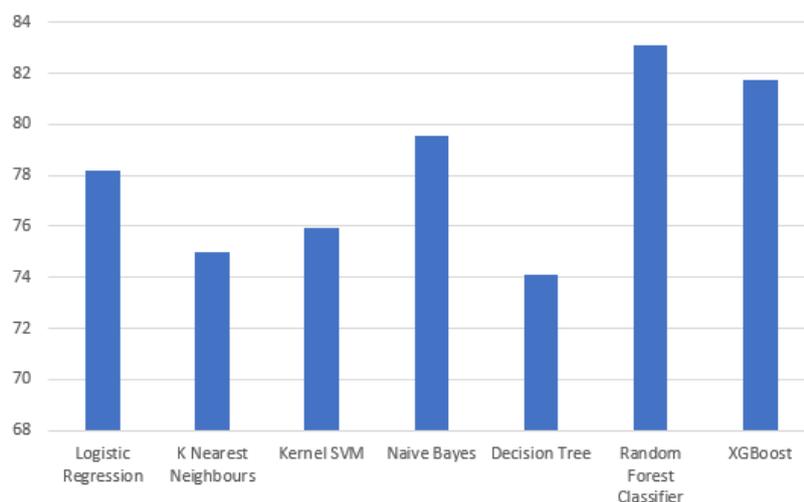
| <i>Название модели</i>                                      | <i>Плюсы</i>  | <i>Минусы</i>   |
|---|---|---|
| 6. Random Forest Classifier (Классификация случайного леса) | 1. Имеет высокую точность предсказания.<br>2. Не чувствителен к масштабированию.<br>3. Хорошо работает с пропущенными данными | 1. Нет формальных выводов, доступных для оценки важности переменных.<br>2. Алгоритм склонен к переобучению на некоторых задачах, особенно на зашумленных данных |

Основываясь на исследованиях, проведенных по диагностированию риска сердечной недостаточности лучший результат показал: Random Forest Classifier<sup>5</sup>.

Таблица 2. Точность и стандартное отклонение моделей

| <i>Название модели</i>                                      | <i>Точность</i> | <i>Стандартное отклонение</i> |
|---|-----------------|-------------------------------|
| 1. Logistic Regression (Логистическая регрессия)            | 78.20 %         | 8.53 %                        |
| 2. K Nearest Neighbours (к-ближайшие соседи)                | 75.02 %         | 6.22 %                        |
| 3. Kernel SVM (Метод опорных векторов)                      | 75.93 %         | 7.84 %                        |
| 4. Naïve Bayes (Наивный байесовский классификатор)          | 79.55 %         | 9.67 %                        |
| 5. Decision Tree (Дерева решений)                           | 74.13 %         | 5.70 %                        |
| 6. Random Forest Classifier (Классификация случайного леса) | 83.08 %         | 7.02 %                        |

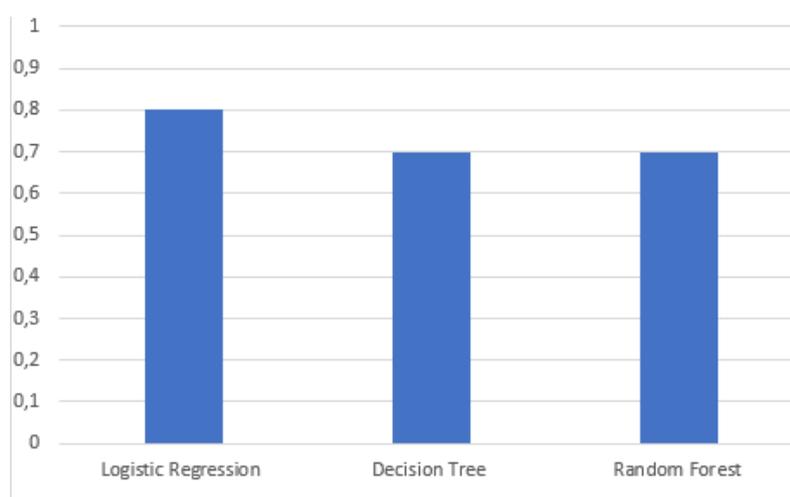
<sup>5</sup> Хабр // URL: <https://habr.com/ru/post/448892/> (дата обращения: 15.11.2022).



*Рисунок 1. Отображение точности моделей*

По данным результатам (рис. 1) предпочтительнее взять за основу модель Классификация случайного леса, так как точность по сравнению с другими немалого, но все же выше и достигает 83%, при этом стандартное отклонение в пределах 7%, что также является хорошим результатом в сравнении с моделью Дерево решений, у которого 5%, но это нивелирует высокая точность первого.

А уже по результатам другого исследования, ориентированного на выявление ССЗ лучшим себя, показал метод Логистической регрессии, но стоит заметить, что тестирование проводилось с использованием лишь трех моделей машинного обучения (рис. 2).



*Рисунок 2. Точность 3-х моделей*

Лучший результат среди всех показала, как ранее было сказано, логическая регрессия с точностью 80%. В связи с такими результатами стоит рассмотреть в

научной работе два вида моделей: Логистическую регрессию, классификатор случайного леса.

Также на точность результата можно будет повлиять, используя k-кратную перекрестную проверку, а после чего настроить гиперпараметры.

Стоит отметить, что перед тем, как отправлять данные на обучение с ними необходимо будет «поработать», используя выборку данных, очистку данных, генерацию признаков, интеграцию, форматирование и прочее.

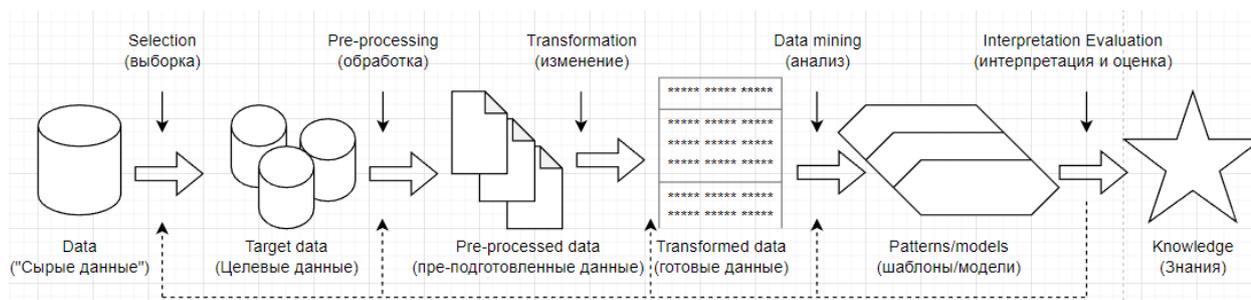


Рисунок 3. Работа с «сырыми» данными

По итогу можно сказать, что, решая задачу диагностики стоит много подзадач, исходя из которых сложится результат решения проблемы, но проведя уже такой небольшой анализ методов и описав подцели работы с “сырыми” данными:

—во-первых — это сильно упростит реализацию поставленной задачи;

—во-вторых — данный анализ, описанный выше, дал представление, как должна выглядеть блок-схема, визуализирующая процесс обработки данных (рис. 3).

### Список источников и литературы

1. Хабр // URL: <https://habr.com/ru/post/448892/> (дата обращения: 15.11.2022).

2. Medium // URL: <https://medium.com/nuances-of-programming/все-модели-машинного-обучения-за-5-минут-9270566197e7> (дата обращения: 15.11.2022).

3. BMC // URL: <https://bmcmedinformdecismak.biomedcentral.com/articles/10.1186/s12911-019-0918-5> (дата обращения: 15.11.2022).

© Нурматова Е.В., Пылаев К.В., 2022