

Прогнозирование риска развития сердечно-сосудистого заболевания пациента

Модель машинного обучения

Решение задачи



Лисовский Дмитрий Александрович

Коротко о себе

Лисовский Дмитрий Александрович, 29 лет

Образование: Врач (Лечебное дело, 2016 год)

Специальности:

- Хирургия
- Организация здравоохранения и общественное здоровье

Дополнительное образование:

Аналитик BI (Нетология)

Data Scientist (Skillfactory)

Предыдущие места работы:

НИИ ОЗИММ ДЗМ - Аналитик

Сеченовский Университет - Директор Департамента развития регионального здравоохранения



Актуальность

Умершие в РФ по классам причин смерти 2021 год

Причина смерти	на 100 000	%
Болезни системы кровообращения	638.9	38,3 %
COVID-19	318.4	19,0 %
Новообразования	193.7	11,6 %
Внешние причины смертности	95.1	5,7 %
Другое	424.0	25,4 %

Формулировка задачи

Задача — разработать модель машинного обучения, цель которой — предсказание наличия **сердечно-сосудистых заболеваний**.

Перечень сердечно-сосудистых заболеваний для предсказания:

1. Артериальная гипертензия
2. Ишемическая болезнь сердца
3. Сердечная недостаточность
4. Острое нарушение мозгового кровообращения
5. Другие заболевания сердца

Описание данных

Всего признаков (столбцов) - **39** (в том числе **5 целевых**):

Величина тренировочного датасета - **955 строк**

Величина тестового датасета - **638 строк**

Данные являются результатом заполнения опросника

Количество опрошенных у которых есть заболевания (всего 955 человек)				
АГ	ОНМК	ИБС	СН	Другое
509	41	117	96	86

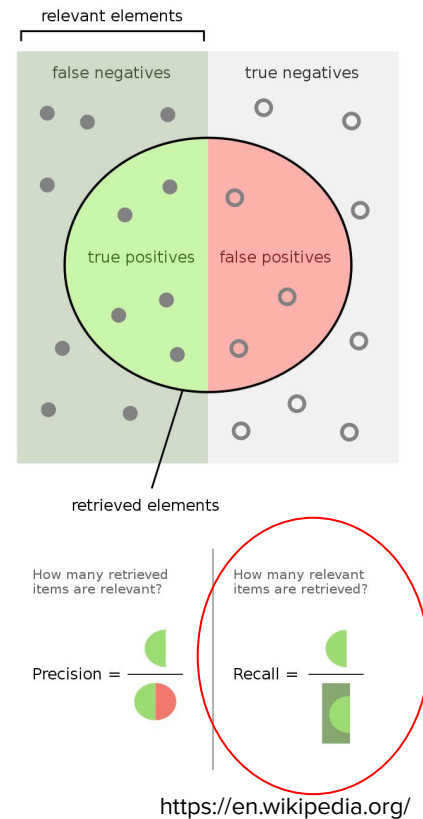
Уточнение задачи

В связи с малым количеством наблюдений и **малым** числом опрошенных, которые имеют заболевания такие как ИБС, СН, ОНМК и “Другие”:

Задача - выявление **высокого риска** наличия/развития конкретных заболеваний

Таким образом, результат работы модели следует трактовать как **“высокий риск развития заболевания”**, а не наличие самого заболевания.

В данном случае оправдано использование метрики среднего **recall**



Список признаков

'ID'	'id'	идентификатор опрошенного
'Пол'	'sex'	пол
'Семья'	'family'	семейное положение
'Этнос'	'ethnos'	этнос
'Национальность'	'nationality'	национальность
'Религия'	'religion'	отношение к религиозной группе
'Образование'	"education"	уровень образования
'Профессия'	"profession"	профессия
'Вы работаете?'	'job'	наличие работы в данный момент
'Выход на пенсию'	'retire'	является ли опрошенный пенсионером
'Прекращение работы по болезни'	'stop_work_due_disease'	связано ли прекращение работы с болезнью
'Сахарный диабет'	'diabetes'	наличие сахарного диабета (любого типа)
'Гепатит'	'hepatitis'	наличие гепатита
'Онкология'	'oncology'	наличие онкологических заболеваний
'Хроническое заболевание легких'	'chronic_lung_disease'	наличие хронического заболевания легких
'Бронхиальная астма'	'bronchial_asthma'	наличие бронхиальной астмы
'Туберкулез легких '	'tuberculosis'	наличие туберкулеза
'ВИЧ/СПИД'	'hiv/aids'	наличие ВИЧ/СПИД
'Регулярный прием лекарственных средств'	'intake_medicines'	факт регулярного приема лекарственных средств
'Травмы за год'	'trauma_last_year'	факт получения травм за последний год
'Переломы'	'fractures'	наличие переломов
'Статус Курения'	'smoking'	статус курения
'Возраст курения'	'smoking_duration'	длительность курения

'Сигарет в день'	'cigarettes_per_day'	количество сигарет в день
'Пассивное курение'	'passive_smoking'	наличие курящих в окружении
'Частота пасс кур'	'passive_smoking_frequency'	частота нахождения среди курящих
'Алкоголь'	'alcohol'	отношение к приему алкоголя
'Возраст алког'	'alcohol_duration'	длительность употребления алкоголя (годы)
'Время засыпания'	'time_sleep_onset'	время засыпания
'Время пробуждения'	'time_sleep_upset'	время пробуждения
'Сон после обеда'	'midday_sleep'	сон в середине дня (после обеда)
'Спорт, клубы'	'sport'	факт регулярного занятия спортом(фитнесом)
'Религия, клубы'	'religion_clubs'	факт регулярного посещения религиозных мероприятий, групп
'ID_y'	'id_y'	идентификатор (связь с таблицей диагнозов)
'Артериальная гипертензия'	'arterial_hypertension'	наличие артериальной гипертензии
'ОНМК'	'stroke'	наличие острого нарушения мозгового кровообращения (инсульта)
'Стенокардия, ИБС, инфаркт миокарда'	'IHD'	наличие ишемической болезни сердца (все формы)
'Сердечная недостаточность'	'heart_failure'	наличие сердечной недостаточности
'Прочие заболевания сердца'	'other_cardio_diseases'	наличие других сердечных заболеваний

Дополнительные признаки

<i>unknown_id_feature</i>	Последняя цифра id (всего возможных цифр 4) представленная как порядковый признак
<i>'id_1', 'id_2', 'id_3', 'id_4'</i>	Последняя цифры id представленная как 4 отдельных колонки (закодированы как 1/0)
<i>'comorbid_count'</i>	Количество болезней (сумма по столбцам diabetes, chronic_lung_diseases и др.)
<i>'is_comorbid'</i>	Наличие хотя бы одного заболевания (в столбцах diabetes, chronic_lung_diseases и др.)
<i>'fracture_last_year'</i>	Признак о переломе кости за последний год (предположительно)
<i>'trauma_on_retire'</i>	Признак о наличии травмы за последний год будучи в пенсионном возрасте
<i>'fracture_on_retire'</i>	Признак о переломе кости за последний год (предположительно) будучи на пенсии
<i>'smoker_score'</i>	Рейтинг курильщика (0 - отсутствие курение, 1 - курение в прошлом, 2 - курение в настоящее время)
<i>'smoking_duration_score'</i>	Рейтинг длительности курения (каждые 5 лет курения - 1 балл, максимальный балл - 9)
<i>'cigarettes_per_day_log'</i>	Прологарифмированный признак количества сигарет в день
<i>'cigarettes_per_day_score'</i>	Рейтинг количества сигарет в день (каждые 10 сигарет - 1 балл, максимальный балл - 5)
<i>'smoking_score_int_log'</i>	Варианты рейтинга курильщика (комбинации рейтинга курильщика, рейтинга длительности курения и количества сигарет в день)
<i>'smoking_score_int_score'</i>	
<i>'smoking_score_int_score_2'</i>	
<i>'smoking_score_int_score_3'</i>	
<i>'smoking_score_int_score_log'</i>	

<i>'alcohol_score'</i>	Рейтинг употребления алкоголя (0 - никогда не употреблял, 1 - употреблял в прошлом, 2 - употребляю на стоящее время)
<i>'alcohol_duration_score'</i>	Рейтинг длительности приема алкоголя, где каждые 5 лет - 1 балл, максимальный балл - 9
<i>'alcohol_duration_log'</i>	Прологарифмированный признак длительности употребления алкоголя
<i>'alcohol_int_score_1'</i>	Варианты рейтинга употребления алкоголя (комбинации рейтинга употребления и длительности)
<i>'alcohol_int_score_2'</i>	
<i>'alcohol_int_score_3'</i>	
<i>'sleep_onset_early_22'</i>	Признаки указывающие на конкретный час засыпания
<i>'sleep_onset_22'</i>	
<i>'sleep_onset_23'</i>	
<i>'sleep_onset_later_after_0'</i>	
<i>'early_onset'</i>	Раннее засыпание (ранее 23 часов)
<i>'lately_onset'</i>	Позднее засыпание (после 23 часов)
<i>'sleep_upset_early_6'</i>	Признаки указывающие на конкретный час просыпания
<i>'sleep_upset_6'</i>	
<i>'sleep_upset_7'</i>	
<i>'sleep_upset_8'</i>	
<i>'sleep_upset_after_9'</i>	
<i>'early_upset'</i>	Ранний подъем (6 - 7 часов утра)
<i>'lately_upset'</i>	Поздний подъем (после 8 утра)
<i>'sleep_time'</i>	Длительность сна
<i>'high_amount_sleep'</i>	Признак указывающий на длинную продолжительность сна (более 9 часов)
<i>'low_amount_sleep'</i>	Признак указывающий на короткую продолжительность сна (менее 6 часов)
<i>'cardio_score'</i>	Рейтинг сердечно сосудистого риска (Модифицированный Jakarta Score)

Описание работы модели

Использованные библиотеки: **pandas, sklearn, imblearn**

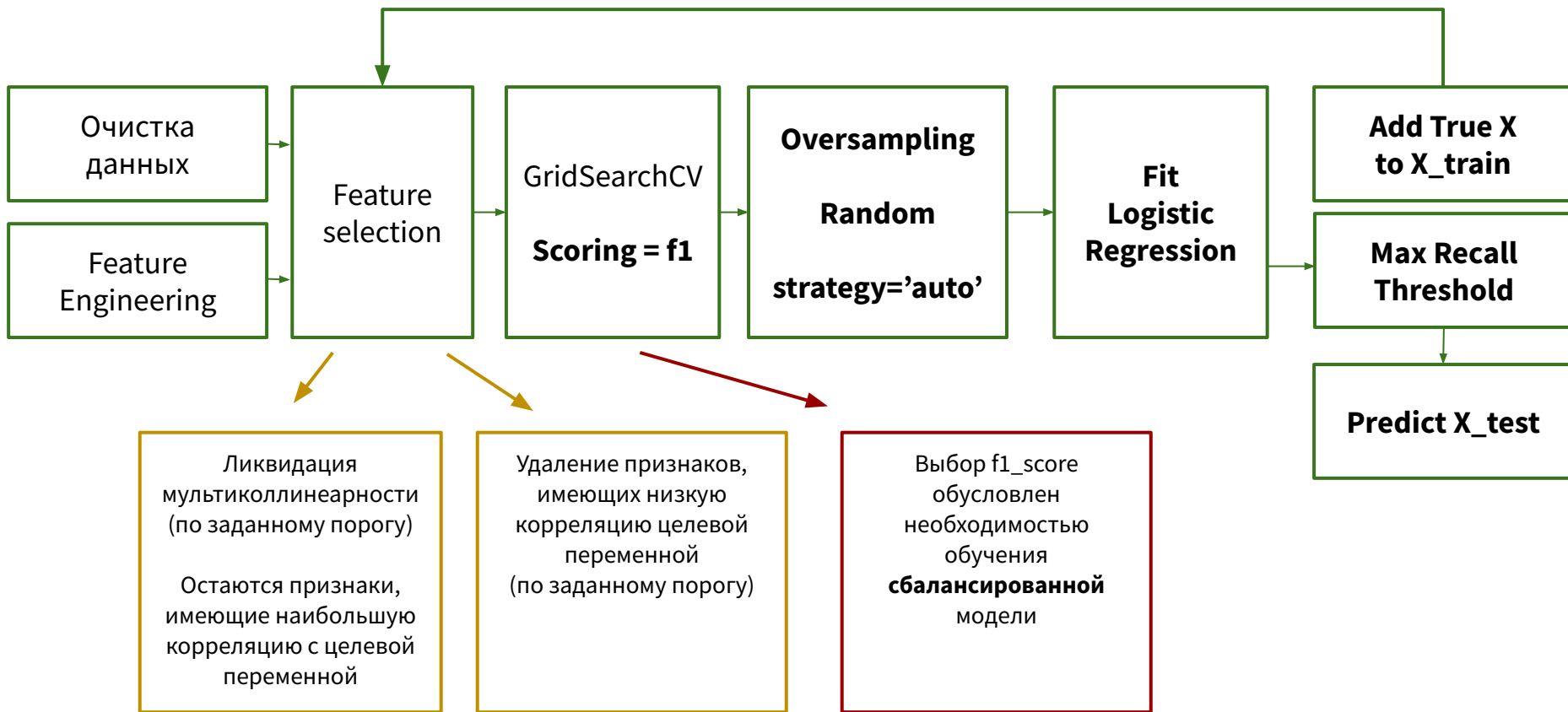
Использованные модули: **LogisticRegression, GridsearchCV, RandomOversampler**

Основные этапы работы **прототипа**:

1. Очистка данных, перевод признаков формат, воспринимаемый моделью (числа)
2. Генерация дополнительных признаков, таких как: “рейтинг курильщика”, “рейтинг сердечно сосудистого риска” и других из имеющихся данных
3. Выбор признаков на основе наибольшей корреляции с каждым наличием каждого заболевания в отдельности
4. Обучение модели логистической регрессии для определения каждого заболевания в отдельности (всего 5 моделей, которые в последующем могут быть объединены в одну последовательно выполняющуюся программу)
5. Каждая модель вычисляет число от 0 до 1 для каждого нового наблюдения (опроса), что является условной “вероятностью” наличия болезни.
6. Модель автоматически выбирает пороговое значение (от 0 до 1) для определения наличия или отсутствия болезни для максимизации целевой метрики

Схема работы модели

Следующее заболевание



Результаты работы модели

Результаты на hacks-ai.ru
(публичный лидерборд)
без установленного random_state

submit_stable_figuration_10_no_
seed_5.csv 0.733723

submit_stable_figuration_10_no_
seed_1.csv 0.725988

submit_stable_figuration_10_no_
seed_2.csv 0.711430

submit_stable_figuration_10_no_
seed_4.csv 0.719489

submit_stable_figuration_10_no_
seed_3.csv 0.722951

Из 638 наблюдений, модель
предсказывает наличие болезней
в количестве:

Артериальная гипертензия	336	
ОНМК		203
Стенокардия, ИБС, инфаркт миокарда	293	
Сердечная недостаточность	299	
Прочие заболевания сердца	298	
dtype: int64		

Модель явно предсказывает болезней **больше, чем в реальности** есть у опрошенных, но это хорошо, так как таким образом **формируются группы риска** (во многом это благодаря выбора метрики recall)

Возможность применения модели

Поскольку задача модели - определить группу риска, целью работы модели может являться **побудить человека обратиться в медицинское учреждение для комплексного обследования**

Где такой опросник может быть заполнен?

1. В рамках более широкого опроса, вроде переписи населения и др.
2. С помощью терминала с сенсорным экраном
 - а. В социальном учреждении
 - б. В парке
 - с. др.

Требуется разработка пилотного проекта с исследованием результатов.
Возможно использование модели будет неэффективно



<https://kiosks.ru/index.php/product/terminal-dlya-polikliniki-atm-med/>

Дополнительные возможности использования модели

Благодаря тому, что модель является интерпретируемой, есть возможность в результате определять не только группу риска, но и выводить некий совет по дальнейшему снижению риска, например:

1. Отказаться от курения (рекомендации + предложение посетить центр здоровья*)
2. Заняться спортом (рекомендации + предложения)
3. Другие рекомендации

* в каждой поликлинике функционирует центр здоровья, в рамках которого работает кабинет по отказу от курения

Приказ Минздравсоцразвития России от 15.05.2012 № 543н

“Об утверждении Положения об организации оказания первичной медико-санитарной помощи взрослому населению”

Контактные данные

Лисовский Дмитрий Александрович

тел. +7-985-999-65-91

e-mail: Lisikux@gmail.com

telegram: @Lisikux