

# Этапы обработки данных в системах, использующих машинное обучение

Сергей Владимирович Аксёнов,

Доцент отделения информационных технологий ИШИТР,

Томский политехнический университет

### Обработка данных в оффлане режиме

Yandex Search API

Kaggle.com API

> GitHub REST API

Подготовка к исследованию

Анализ данных, признаки, Схема данных,

Нормализация Соединение Агрегация Обработка естественного языка (NLP) Создание репозитория

JSON CSV Pickle MySQL MongoDB Cassandra Конструирован ие признаков (Feature Engineering)

Сокращение размерности Отбор признаков Обработка естественного языка (NLP) Токенизация Лемматизация Стемминг

Построение модели (Data Modelling)

Обучение с учителем: Классификация Регрессия,

Обучение без учителя: Кластеризация, Поиск аномалий Фильтрация Визуализация

Отчеты, Графики, Гистограммы Временные ряды, Геолокация

### Обработка данных в онлайн режиме

Обработка Потребитель событий потоковых данных Результат Сообщения Фильтрация, Трансформация, Telegram Разделение a a a Дашборды, Streaming данных, Бизнес процессы, Обнаружение Бизнес сервисы паттернов, Геолокация, Предобученные модели

### Задачи очистки данных

- ✓ Работа с отсутствующими данными
- ✓ Выявление нарушения взаимосвязи признаков
- ✓ Обработка некорректных данных
- ✓ Устранение дубликатов
- ✓ Преобразования данных
- ✓ Удаление выбросов

# Удаление дубликатов

Имя пациента	Дата исследования	С-реактивный белок
Эльвира Сафина	21/08/05	
Андрей Семёнов	21/08/05	-14.3
Юлия Келлер	21/08/05	0.5
Эльвира Сафина	21/08/05	-6.1
Андрей Семенов	21/08/05	20.6
Элвира Сафина	/08/05	8.77

Имя пациента	Дата исследования	С-реактивный белок
Эльвира Сафина	21/08/05	8.77
Андрей Семёнов	21/08/05	20.6
Юлия Келлер	21/08/05	0.5

#### Баланс классов

В машинном обучении балансировка классов означает изменение числа примеров в наборе, где наблюдается разные доли примеров для каждого класса.

Перед использованием алгоритма машинного обучения важно избежать дисбаланса классов, потому что наша конечная цель — обучить модель машинного обучения, которая хорошо обобщается для всех возможных классов.

Перед использованием алгоритма машинного обучения очень важно посмотреть на распределение классов, чтобы исправить проблему балансировки классов.

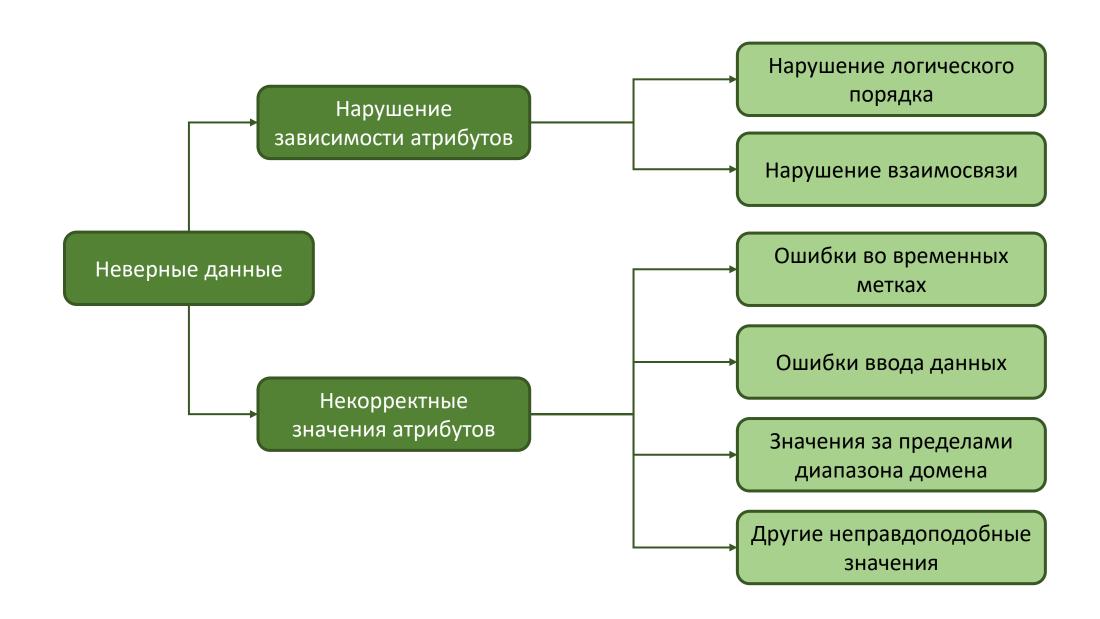
### Отсутствующие значения



# Методы обработки пропущенных значений



### Неверные данные



### Нарушение логического порядка

Запись из истории болезни:

Заболела остро 17.09.2016г, утром, потрясающий озноб, повышение температуры до 39,5С, чувство разбитости, головная боль, тошнота, рвота, выраженная слабость. К утру 16.09.2016г появилась эритема в области правой стопы, обширная гематома на своде стопы.

# Нарушение взаимосвязи

	Имя пловца	Время, сек. Вольный стиль, 50 м.	Скорость движения, м/с
<b>/</b>	Рауза Садыкова	23	2.17
X	Мария Григорьева	24	2.8
<b>/</b>	Исмаил Ахметов	22	2.27

### Ошибки во временных метках

Проведения исследования: 15.01.22 10:14

Поступление образца в лабораторию: 15.01.22 14:47

Исследование образца: 15.01.22 17:25

Получения результата исследования: 16.01.22 9:30

# Ошибки набора данных

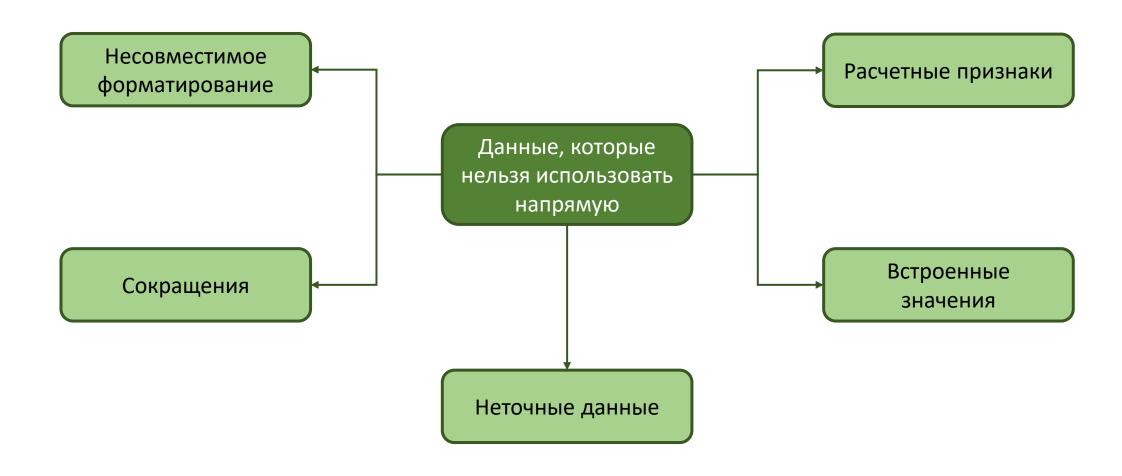
	Бренд	Фактура	Размер	Тип рукава
×	Маргарита	вльвет	35	Без руковов
X	Только ты	гладкий	399	Длинные
<b>/</b>	Savanna	кружевной	42	Короткие
X	BEAUTY	жаккард	38	Длиннные
<b>/</b>	Вирджиния	атлас	44	Короткие

## Значения за пределами диапазона домена



Имя пациента	Холестерин	Общий белок
Роксана Ковалевская	3.93	-1
Марат Денисов	5.72	80.3
Ольга Вальтер	6.24	38.5
Алия Сарымова	3.36	86.9
Серж Оганесян	-0.11	77.4

### Данные, которые нельзя использовать напрямую



### Несовместимое форматирование

Записи из осмотров разных пациентов:

Жалобы: На повышение температуры до 39.2С, распирающую боль в правом плече и предплечье, чувство жара, тошнота.

Жалобы: На высокую температуру до сорока градусов, отёк и гиперемия с гематомой в области свода стопы.

Жалобы: Слабость, повышение температуры до 38,9, жар, озноб, боль, гиперемия, отёк правой голени.

### Сокращения

Анамнез жизни: В детстве перенесла корь, ветряную оспу, ОРВИ, грипп, 1 раз пневмония, туберкулёз отрицает. С 2007г АГІІ, риск 2, ГЛЖ, ИБС, Стенокардия напряжения ФКІІ. Хр. холецистит, на описторхоз не обследована, хр. панкреатит.

План обследовния: OAK, OAM, БАК, коагулограмма, РМП на сифилис, ЭКГ

### Расчетные значения

Имя клиента	Рост	Вес	Индекс массы тела
София	1.75	65	?
Дарина	1.68	73	?
Полина	1.72	69	23.3
Анна	1.81	72	24.3
Камилла	1.62	58	?

### Неточные данные

#### Отзыв покупателя:

Сумка приехала ко мне в розовой версии, и хотя розоватый оттенок здесь действительно есть, что отмечали все мои подруги, кому я её показывала, по сути это, скорее, бежево-кремовый цвет. Впечатление от её специфическое: вроде и не белая, а какая — непонятно.

# Встроенные значения

Адрес дома	Общее потребление электроэнергии	Потребление электроэнергии на освещение
Ул. Пушкина, д.7	540	?
Ул. Пушкина, д.56	730	80
Ул. Светлая, д.102	849	?
Ул. Уральская, д.6	624	145
Ул. Ахматовой, д.23	428	?

### Преобразование категориальных признаков

- Порядковые (можно сравнить их и упорядочить)
- Номинальные (несравнимые)

Адрес Отопление		Наличие кондиционир ования	Общая площадь, кв.м.	Расход энергии	
Пушкина, 5	Электроэнергия	Да	120	Высокий	
Мира, 32	Уголь	Нет	104	Средний	
Королёва, 19	Электроэнергия	Да	145	Высокий	
Сибирская, 6	Природный газ	Нет	158	Средний	
Сибирская, 12	Природный газ	Нет	136	Низкий	

† Номинальный признак

. Порядковый признак

### Порядковые признаки

Задание соответствия вручную исходя из понимания признака.

Адрес Отопление		Наличие кондиционир ования	Общая площадь, кв.м.	Расход энергии
Пушкина, 5	Электроэнергия	Да	120	Высокий
Мира, 32	Уголь	Нет	104	Средний
Королёва, 19	Электроэнергия	Да	145	Высокий
Сибирская, 6	Природный газ	Нет	158	Средний
Сибирская, 12	Природный газ	Нет	136	Низкий

Расход энергии: «Высокий» $\rightarrow$ 2, «Средний» $\rightarrow$ 1, «Низкий» $\rightarrow$ 0

Адрес	Отопление	Наличие кондиционир ования	Общая площадь, кв.м.	Расход энергии
Пушкина, 5	Электроэнергия	Да	120	2
Мира, 32	Уголь	Нет	104	1
Королёва, 19	Электроэнергия	Да	145	2
Сибирская, 6	Природный газ	Нет	158	1
Сибирская, 12	Природный газ	Нет	136	0

### Номинальные признаки

#### Прямое кодирование (унитарное кодирование)

Адрес	Адрес Отопление		Общая площадь, кв.м.	Расход энергии	
Пушкина, 5	Электроэнергия	Да	120	Высокий	
Мира, 32	Уголь	Нет	104	Средний	
Королёва, 19	Электроэнергия	Да	145	Высокий	
Сибирская, 6	Природный газ	Нет	158	Средний	
Сибирская, 12	Природный газ	Нет	136	Низкий	

Адрес	Отопление: Электро- энергия	Отопление: Уголь	Отопление: Природный газ	Наличие кондицио- нирования	Общая площадь, кв.м.	Расход энергии
Пушкина, 5	1	0	0	1	120	Высокий
Мира, 32	0	1	0	0	104	Средний
Королёва, 19	1	0	0	1	145	Высокий
Сибирская, 6	0	0	1	0	158	Средний
Сибирская, 12	0	0	1	0	136	Низкий

### Приведение данных к одному масштабу

Кластеризация и **многие** алгоритмы обучения с учителем крайне нуждаются в приведении признаков к одной шкале, т.к. их вычисления завязаны на учет расстояния между объектами в выборке.

#### Нормализация:

$$x_{norm}^{(i)} = \frac{x^{(i)} - x_{min}}{x_{max} - x_{min}}$$

 $arkappa_{min}$  - наименьшее значение

 $x_{max}$  - наибольшее значение

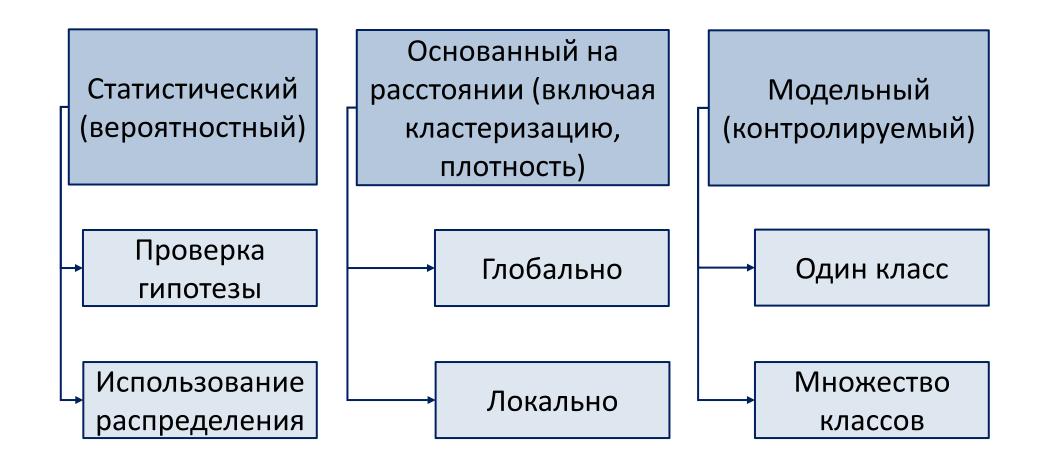
#### Стандартизация:

$$x_{std}^{(i)} = \frac{x^{(i)} - \mu_x}{\sigma_x}$$

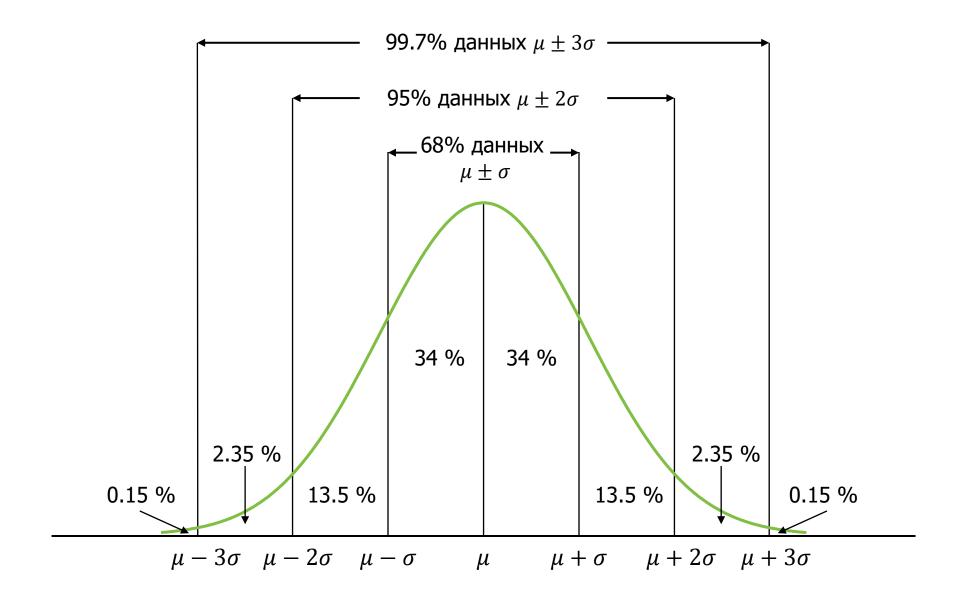
 $\mu_{\chi}$  - эмпирическое среднее

 $\sigma_{\chi}$  - стандартное отклонение

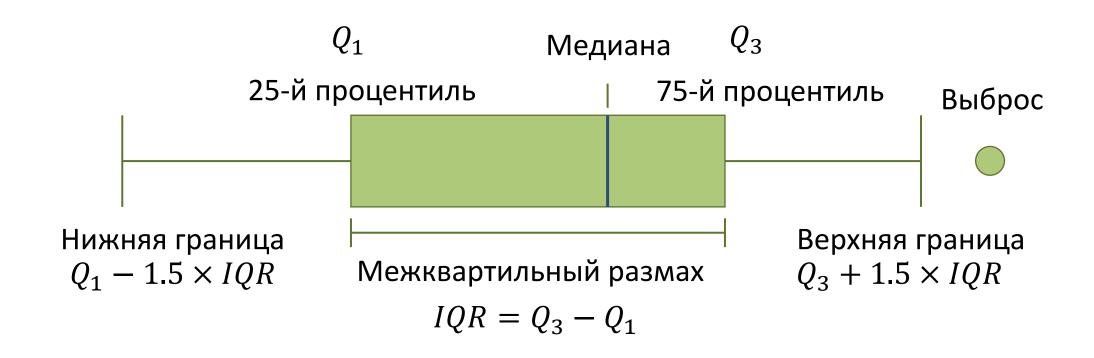
### Методы обнаружения выбросов



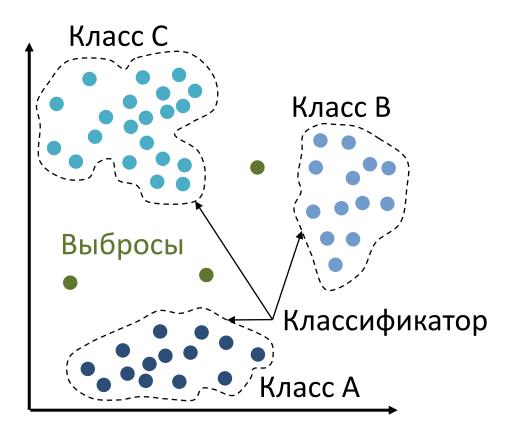
## Статистический подход



### Диаграмма размаха



### Методы обнаружения выбросов на основе моделей



Выделение выбросов для многих классов

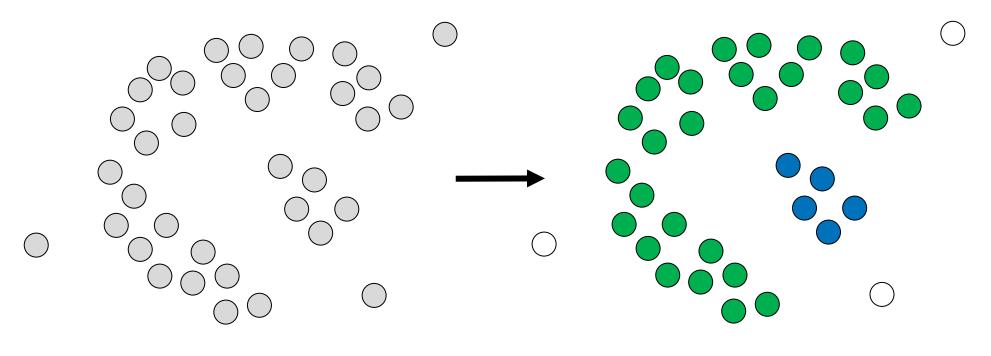


Выделение выбросов для одного класса

### Использование кластеризации

#### Исходные данные

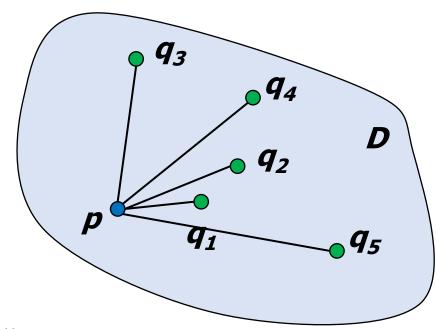
#### Кластеризация DBSCAN



- 🤍 Кластер А
- Кластер В
- - Выброс

### К-расстояние от объекта р

distance(p,q) или d(p,q):



Расстояние между объектом p и его k-ближайшими соседями.

$$d(p, q_1) \le d(p, q_2) \le d(p, q_3) \le d(p, q_4) \le d(p, q_5)$$

D — множество объектов

### Фактор локального выброса

Расстояние достижимости:

$$reach - dist_k(p, q) = \max(k - dist(p), dist(p, q))$$

Плотность локальной достижимости:

$$Lrd_{MinPts}(p) = 1 / \left[ \frac{\sum_{q \in N_{MinPts}(p)} reach - dist_k(p, q)}{|N_{MinPts}(p)|} \right]$$

Фактор локального выброс:

$$LOF_{MinPts}(p) = \frac{\sum_{q \in N_{MinPts}(p)} \frac{Lrd_{MinPts}(q)}{Lrd_{MinPts}(p)}}{|N_{MinPts}(q)|}$$

 $LOF_{MinPts}(p){\sim}1$  такая же плотность, как у соседей

 $LOF_{MinPts}(p) < 1$  более высокая плотность, чем у соседей (невыброс)

 $LOF_{MinPts}(p) > 1$  более низкая плотность, чем у соседей (выброс)

### Эллиптическая оболочка

Элиптическая Выбросы оболочка Невыбросы