



Методы статистического анализа

Основы МО и примеры в ФЭЧ

Ильясов Айдар

Аспирант 2 курса НИЦ КИ, МНС ЛФРП ОФН КИ



План занятий

Занятие 1 (27.10.2022) - Теоретическое

- Основы машинного обучения
- Основные модели машинного обучения

Занятие 2 (03.11.2022) – Практическое, вариативное

- «Усложнённые» модели машинного обучения
- Использование машинного обучения в физике частиц
 - C++
 - Python
- Пример результата работы машинного обучения на C++
- Пример(ы) результата работы машинного обучения на Python
- Источники информации
- Другое? ilyasovaid@yandex.ru

Методы статистического анализа

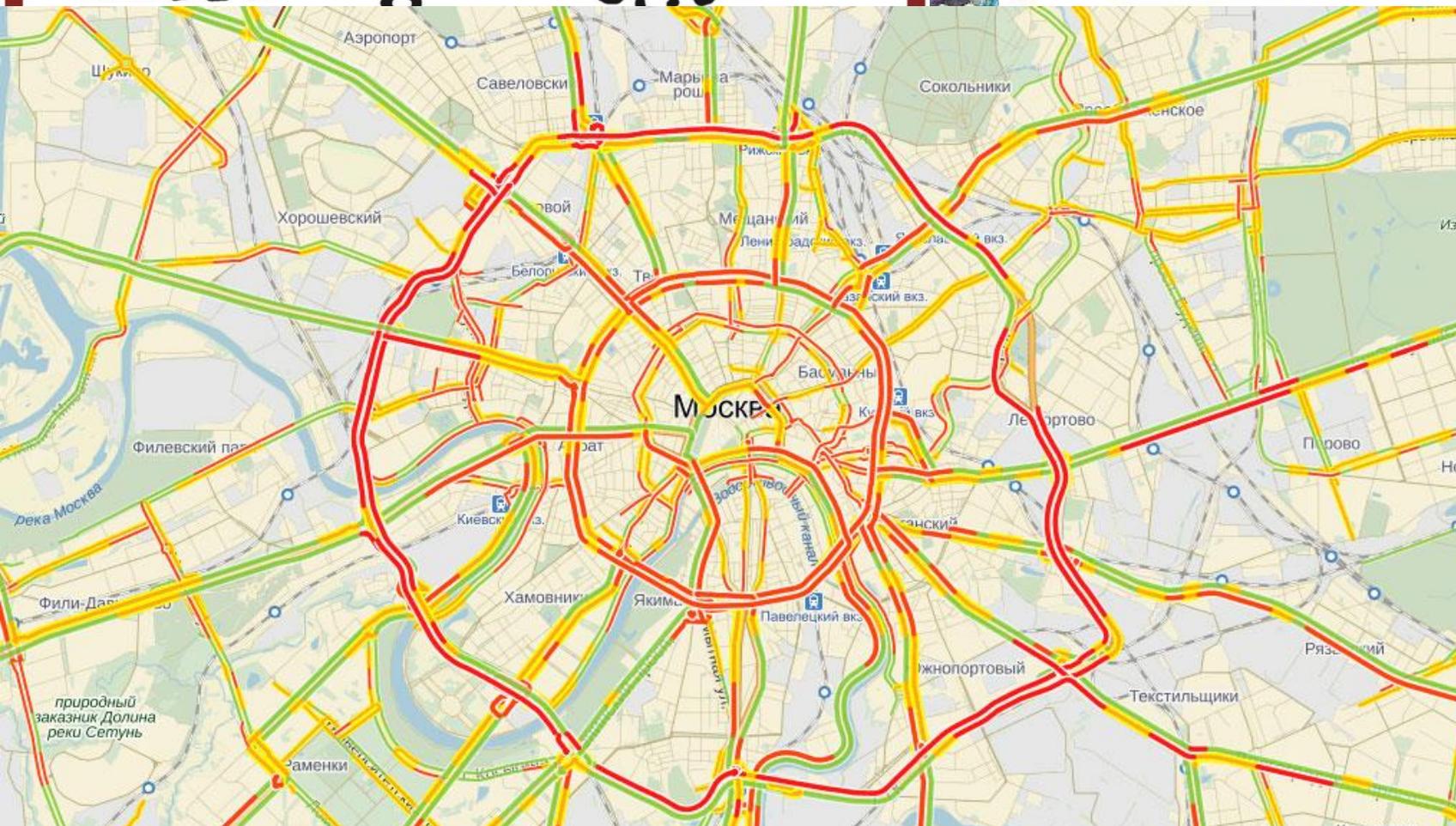
Часть 1: Основы



27.10.2022

Примеры из жизни

morning overtook



google.com

Сообщения Форум Помощь

Все задания Активные 4

gns
click skip

Я.Поллук

Приступить к заданию

Инструкция

Я.Щит

Пройти обучение

Инструкция

Я.Феникс

Пройти обучение

Инструкция

Я.Плутон

Пройти обучение

Инструкция

Я.Поллук

Приступить к заданию

Инструкция

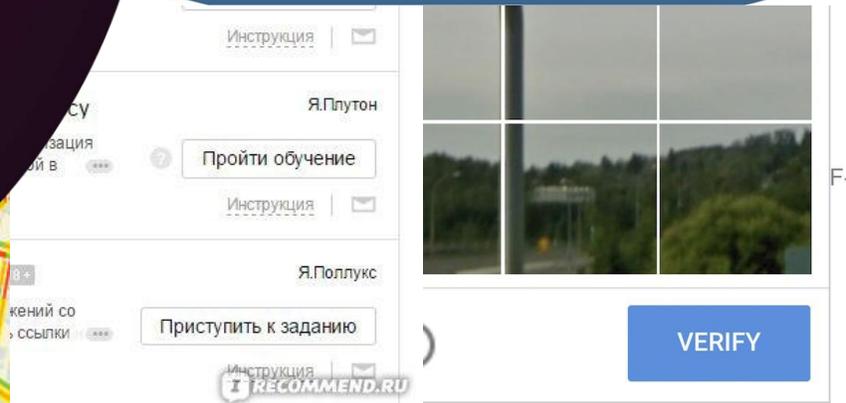
RECOMMEND.RU

VERIFY

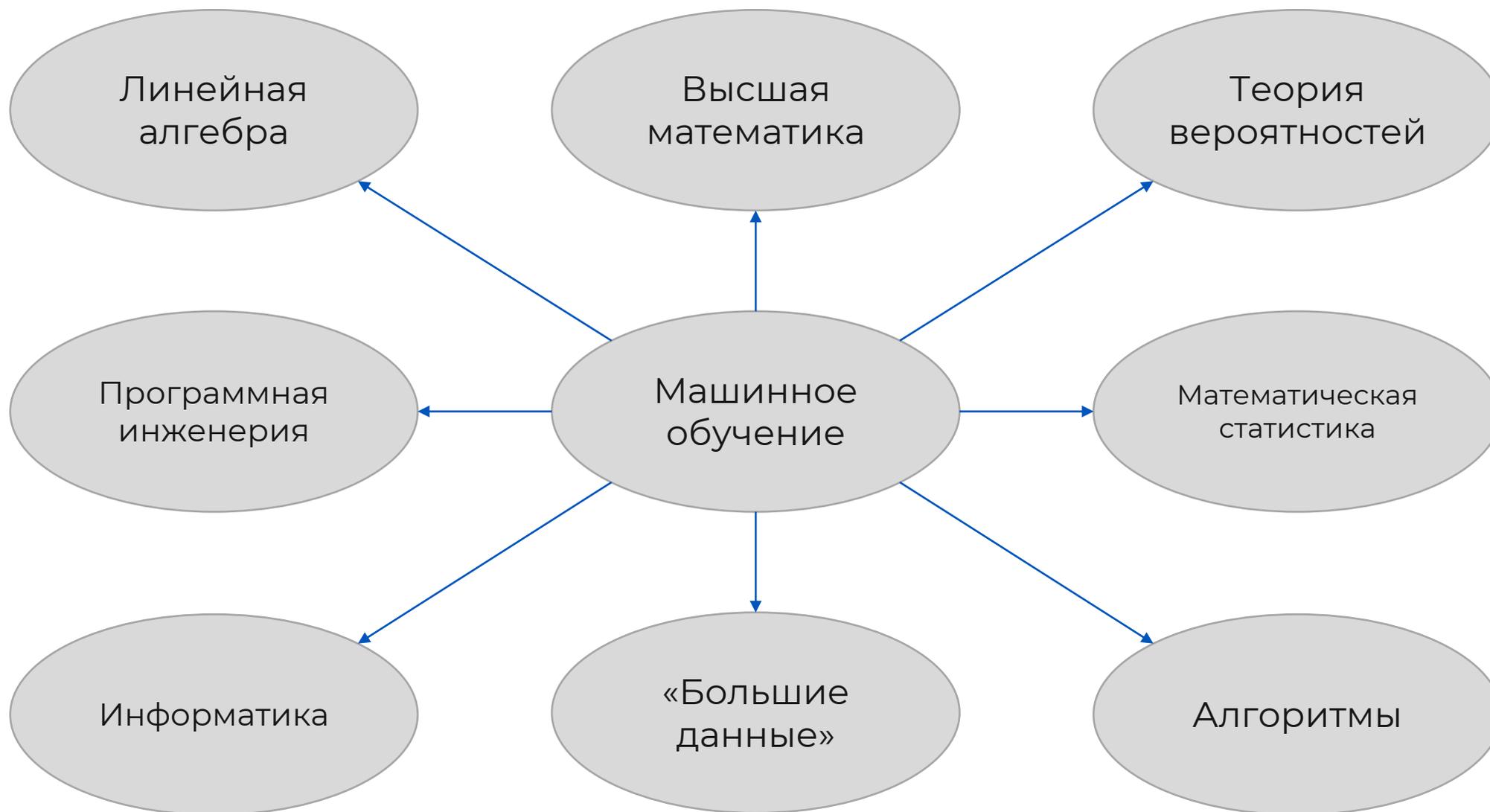
Примеры из жизни



Яндекс Д



Что такое машинное обучение?



Что такое машинное обучение?



Типы данных



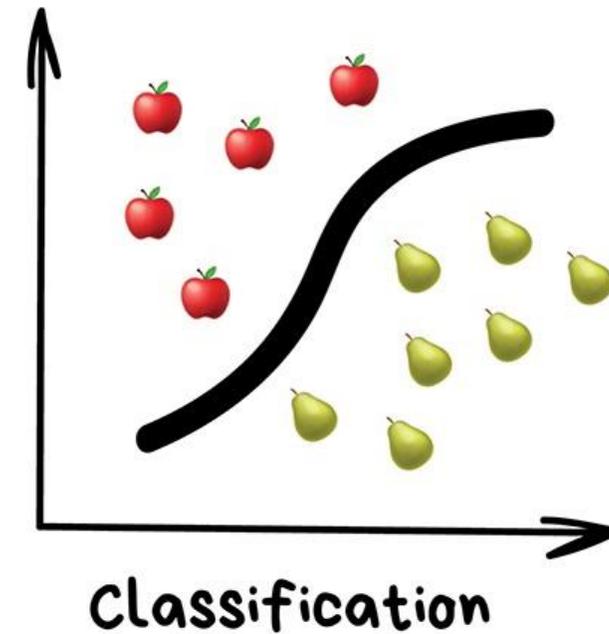
А на самом деле?



*Умный маленький человечек внутри компьютера

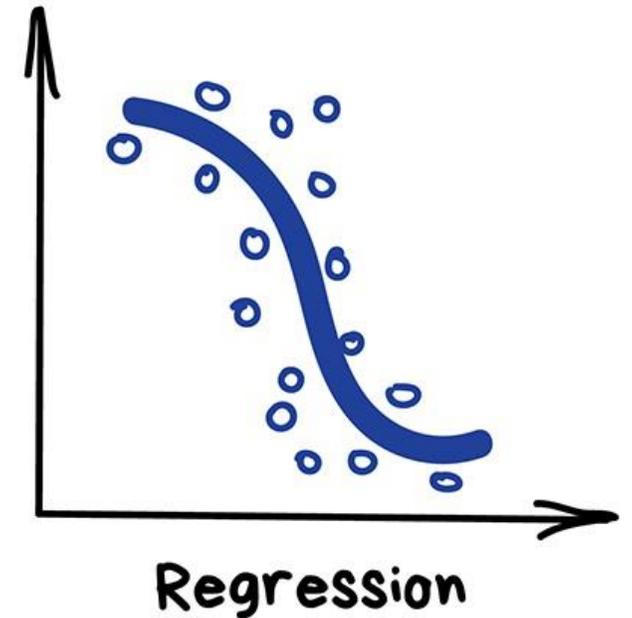
• С учителем:

- Классификация
 - Кошка-собака
 - 1-0
 - Птица-Самолёт-Вертолет
 - Мальчик-Девочка
 - Цветок-Животное



- **С учителем:**

- Классификация
- Регрессия
 - Какого цвета конфета?
 - Сколько будет стоить квартира?
 - Какая будет цена на нефть в 2025 году?



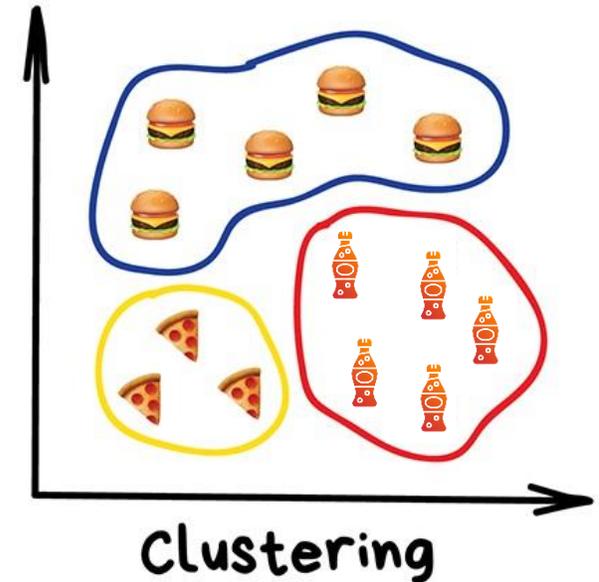
Основные типы задач машинного обучения

- **С учителем:**

- Классификация
- Регрессия

- **Без учителя:**

- Кластеризация
 - Машинное обучение - финансы - игры
 - Разложить похожие вещи по кучкам

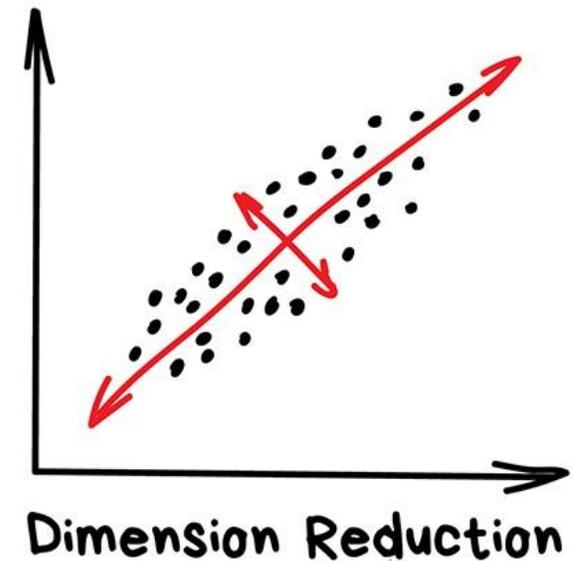


- **С учителем:**

- Классификация
- Регрессия

- **Без учителя:**

- Кластеризация
- Задача уменьшения размерности
 - Анализ параметров наборов данных

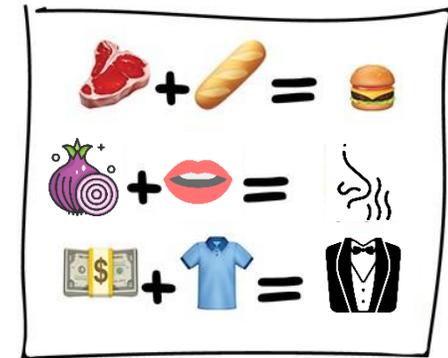


- **С учителем:**

- Классификация
- Регрессия

- **Без учителя:**

- Кластеризация
- Задача уменьшения размерности
- Задача поиска правил
 - Анализ товаров покупаемых вместе
 - Прогноз распродаж



**Association
Rule Learning**

Основные типы задач машинного обучения

- **С учителем:**

- Классификация
- Регрессия

- **Без учителя:**

- Кластеризация
- Задача уменьшения размерности
- Задача поиска правил



Классическое обучение

Основные типы задач машинного обучения

- **С учителем:**

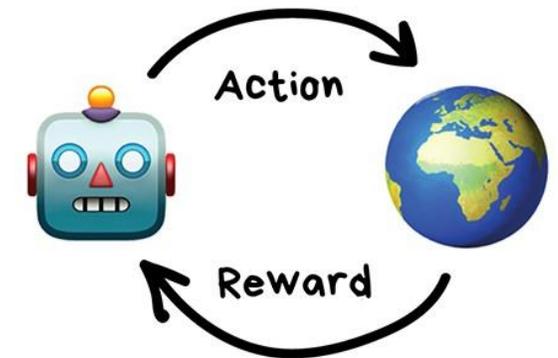
- Классификация
- Регрессия

- **Без учителя:**

- Кластеризация
- Задача уменьшения размерности
- Задача поиска правил

- **Обучение с подкреплением**

Классическое обучение



**Reinforcement
Learning**

Основные типы задач машинного обучения

- **С учителем:**

- Классификация
- Регрессия

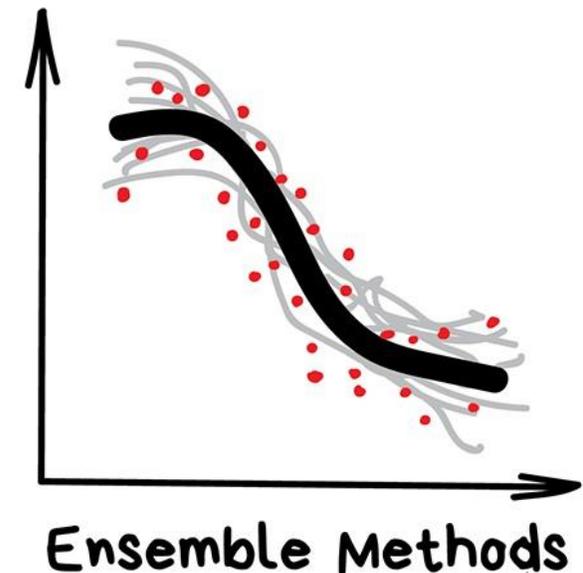
- **Без учителя:**

- Кластеризация
- Задача уменьшения размерности
- Задача поиска правил

Классическое обучение

- **Обучение с подкреплением**

- **Ансамбли**



Основные типы задач машинного обучения

- **С учителем:**

- Классификация
- Регрессия

- **Без учителя:**

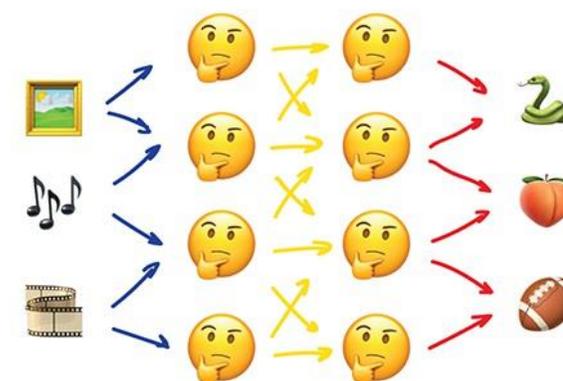
- Кластеризация
- Задача уменьшения размерности
- Задача поиска правил

Классическое обучение

- **Обучение с подкреплением**

- **Ансамбли**

- **Нейросети**



Neural Networks

«Иерархия потребностей» науки о данных



Факт:

У человека появилась идея применения МО **в новой области**

Вопрос:

Что нужно сделать человеку чтобы **с нуля** разработать свой алгоритм машинного обучения?

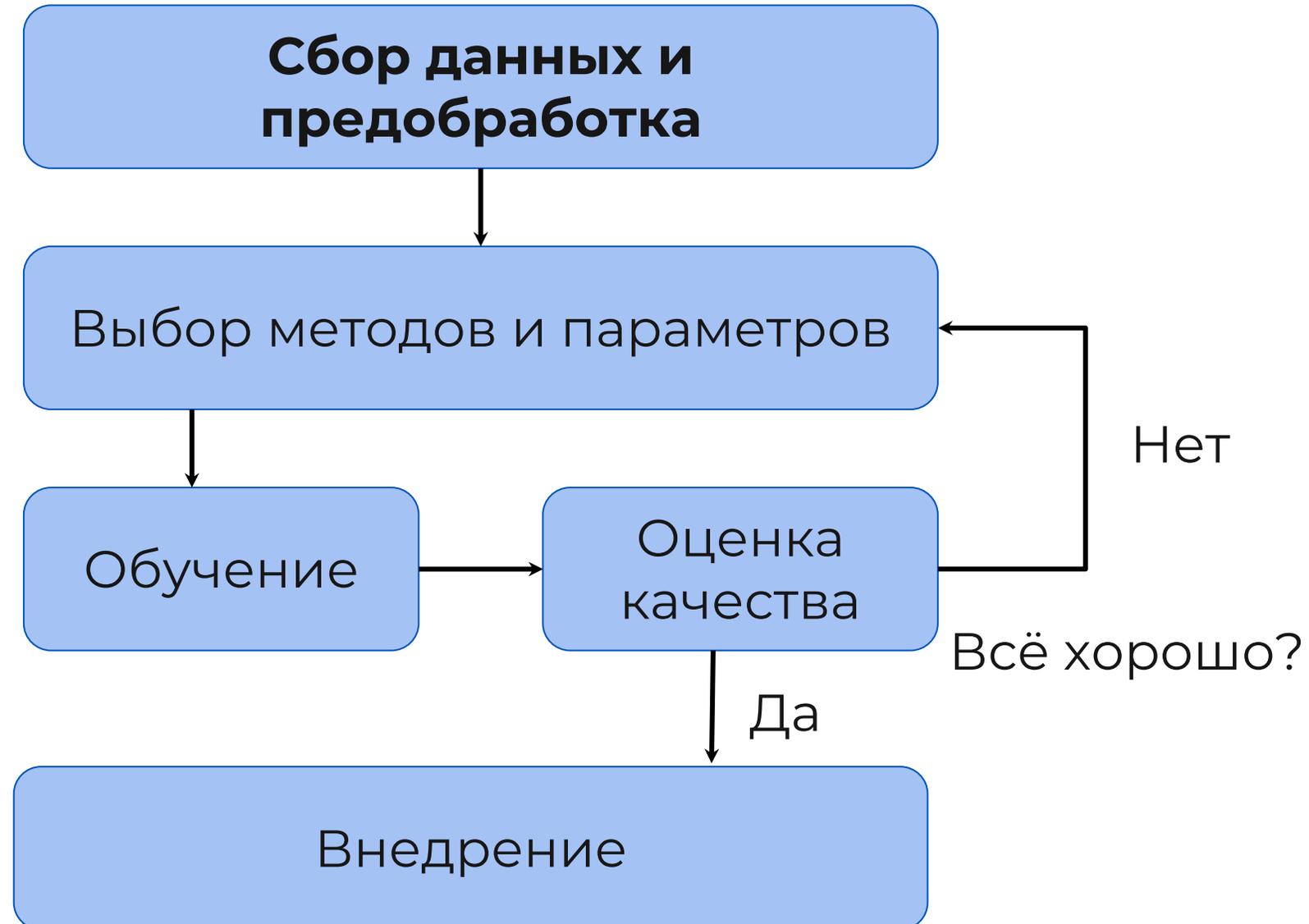
«Иерархия потребностей» науки о данных



Выбор методов и параметров

Сбор и предобработка данных

Ход работы



Сбор и предобработка данных

Сбор данных

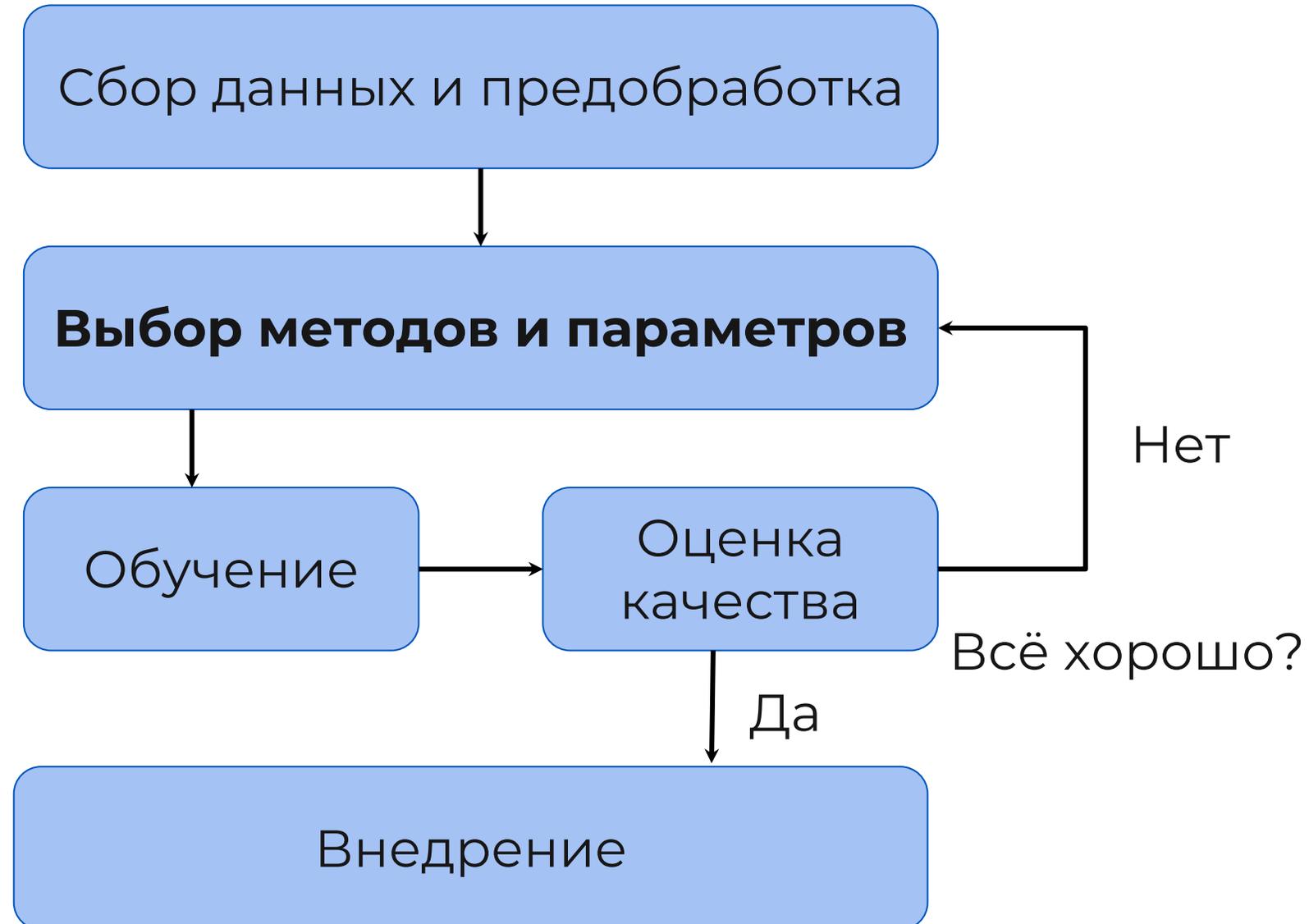
- MNIST
- Открытые базы данных
- Моделирование событий
- Данные с детекторов
- Капча

Сбор данных и
предобработка

Предобработка

- Исключение выбросов и NaN событий
- Исключение шумовых данных

Ход работы



Какой алгоритм выбрать?

Выбор методов и параметров

«Базовый» список моделей машинного обучения

- Наивный байесовский классификатор
- Дерево решений
- Метод опорных векторов
- Регрессия
- Логистическая регрессия
- Принцип главных компонент
- Другие

Какой алгоритм выбрать?

«Расширенный» список моделей машинного обучения

- Нейронная сеть
- Свёрточная нейронная сеть
- К-средних (kNN)
- Случайный лес
- Бустинг над деревьями решений
- Ансамбль моделей
- Другие

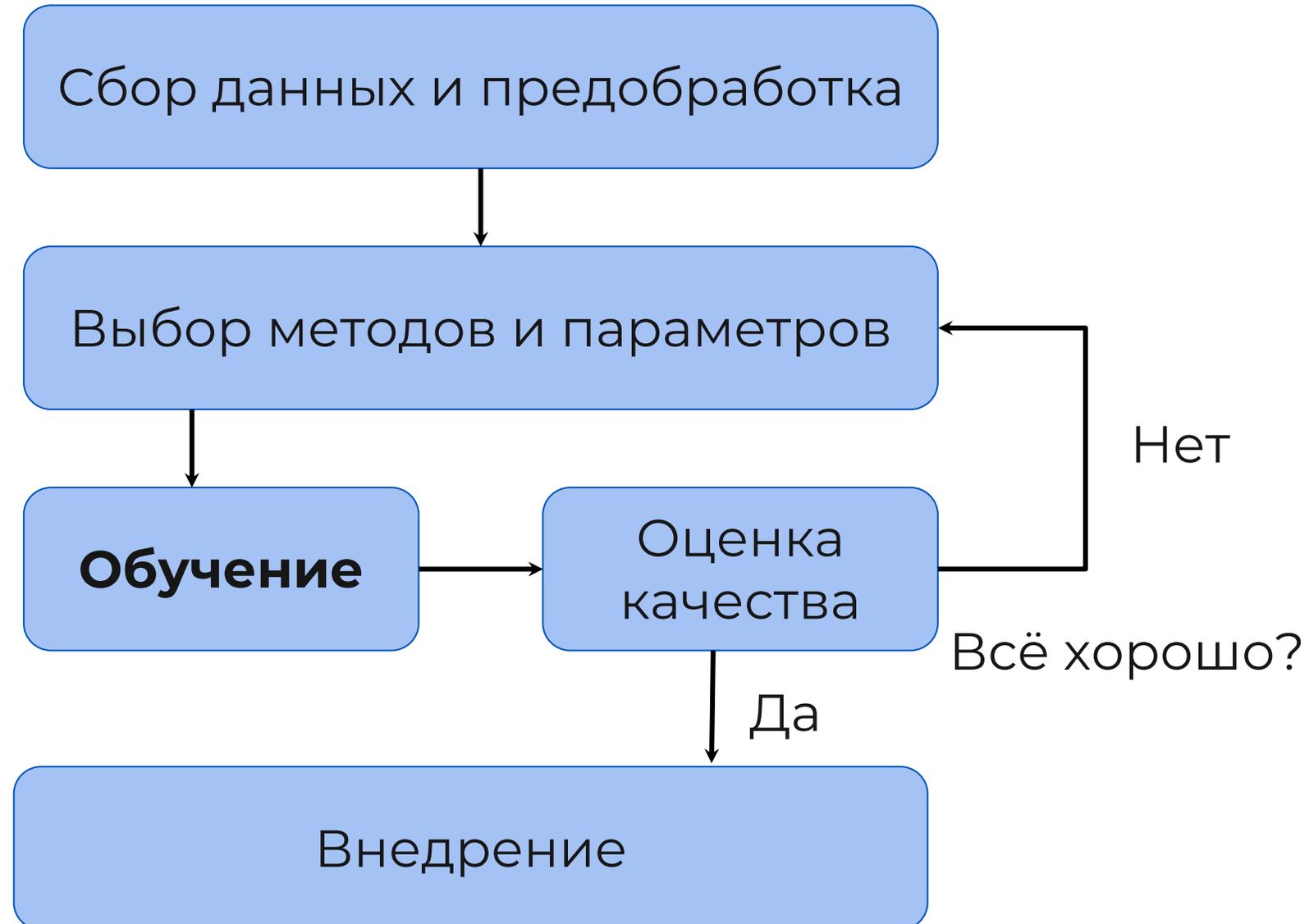
Выбор методов и параметров

Какой алгоритм выбрать?



Выбор методов и параметров

Ход работы



Тренировка модели

Тренировка модели с учителем

- Данные размечены
- Данные разделены на тренировочный и тестовый наборы
- Необходимо либо предсказать класс, либо предсказать значение

Обучение



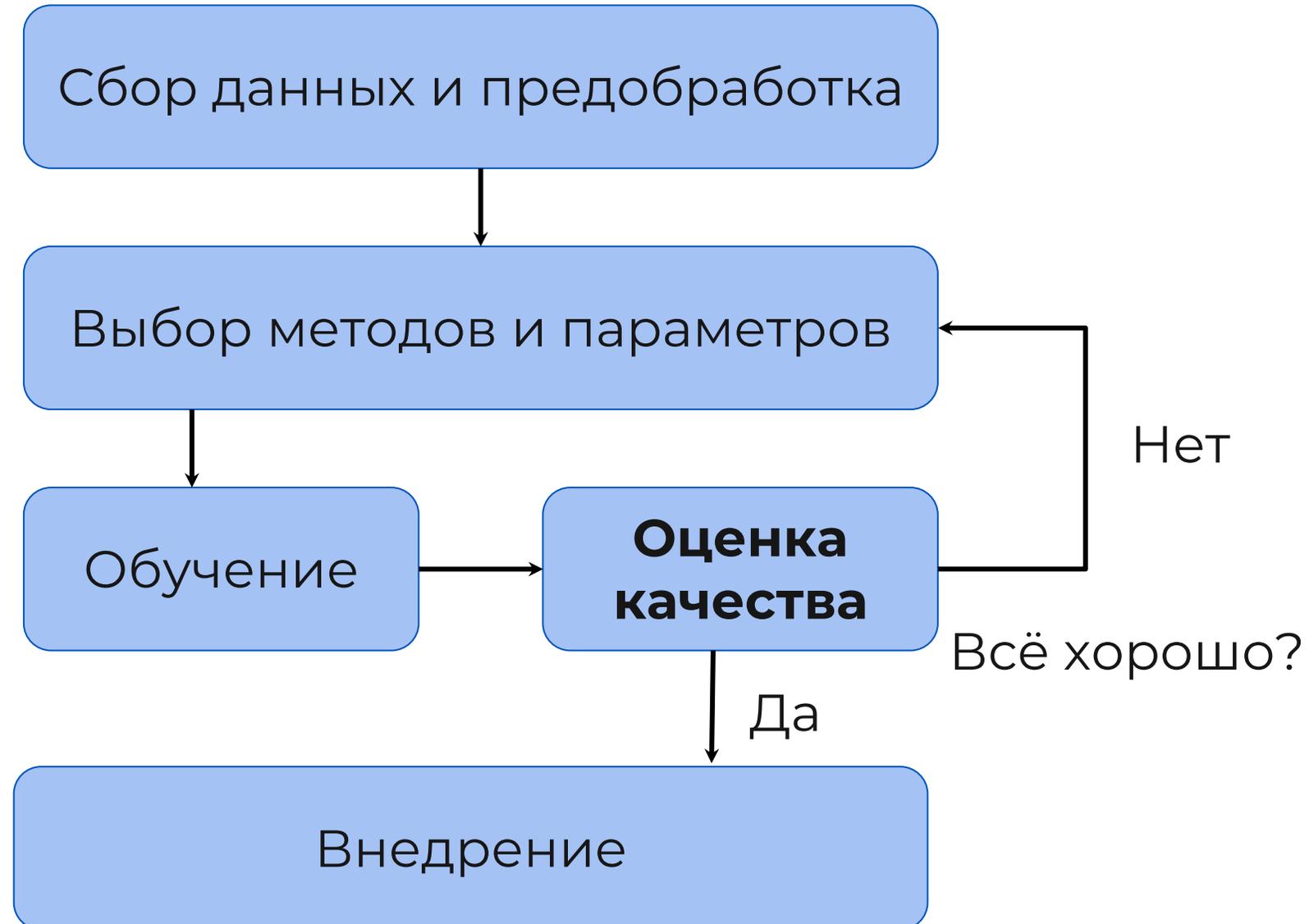
Тренировка модели без учителем

- Данные не размечены
- Единый набор данных для (само)обучения модели
- Необходимо найти классы объектов, зависимости в данных, выявить последовательности

Обучение

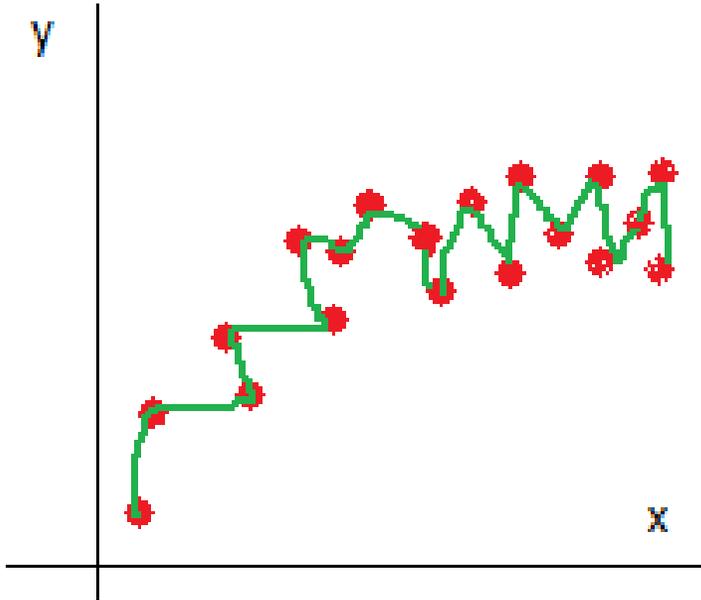


Ход работы

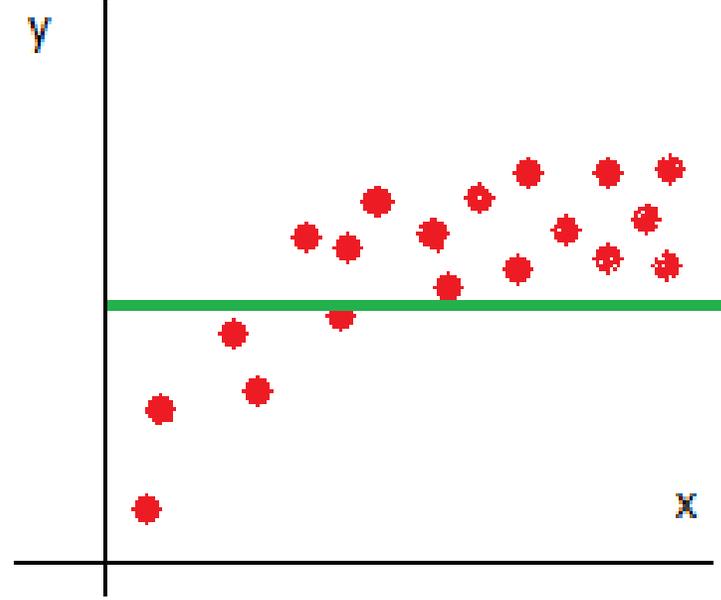


Как оценивать эффективность модели?

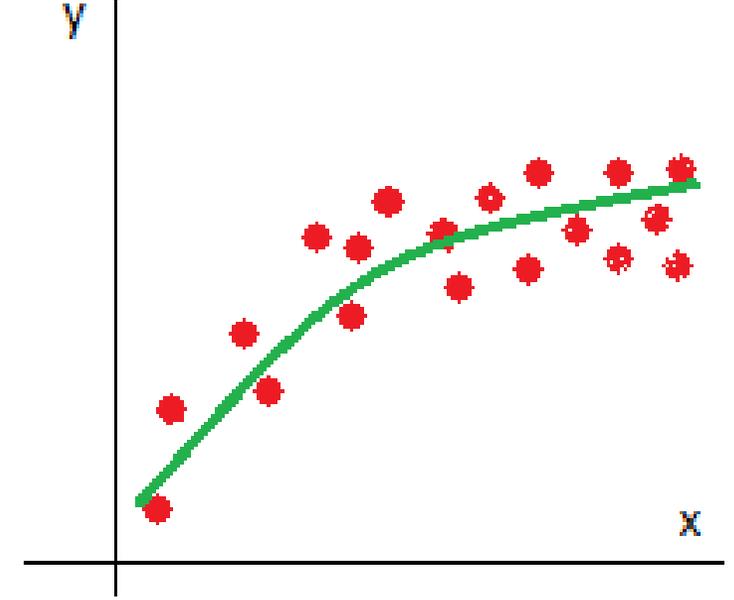
Оценка
качества



Overfitting (High Variance)



Underfitting (High Bias)



Just Right

Метрики

Оценка
качества



	$Y = 1$	$Y = 0$
$X = 1$	True Positive (TP)	False Positive (FP) (Type I error)
$X = 0$	False Negative (FN) (Type II error)	True Negative (TN)

X – результат работы алгоритма,
 Y – истинная метка

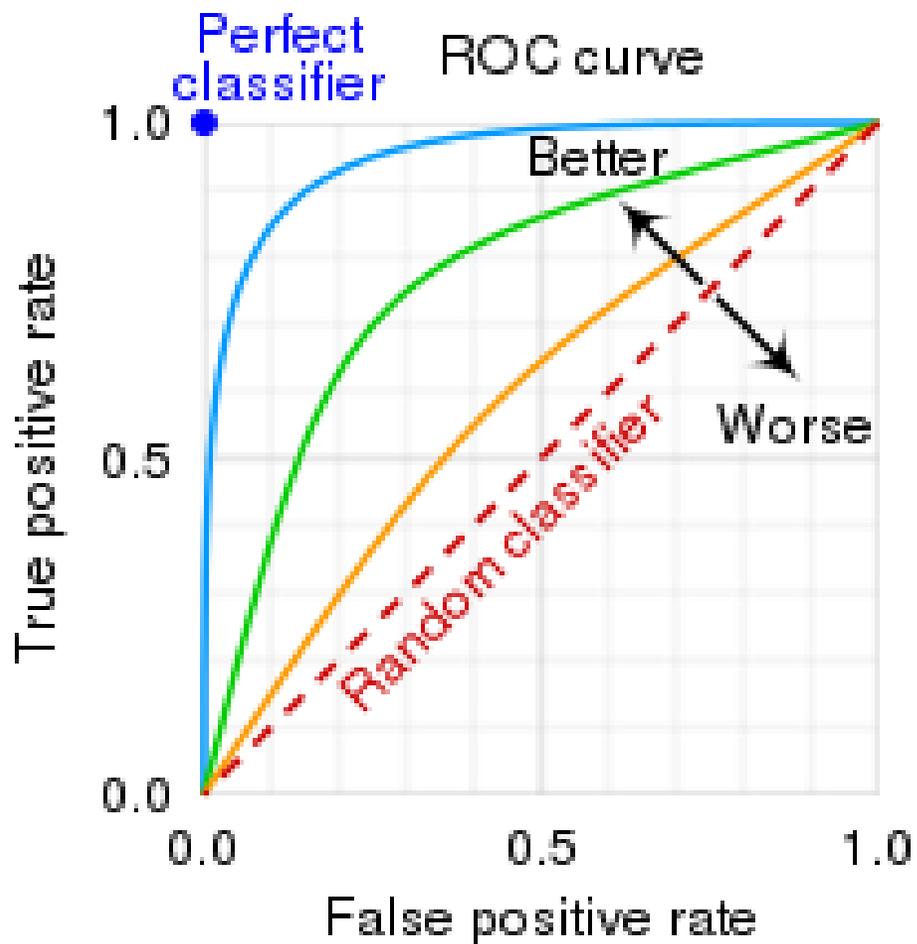
$$\text{Accuracy} = \frac{1}{n} \sum_{i=1}^n \{x_i = y_i\} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Precision} = \frac{\sum_{i=1}^n \{x_i = y_i = 1\}}{\sum_{i=1}^n \{x_i = 1\}} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{\sum_{i=1}^n \{x_i = y_i = 1\}}{\sum_{i=1}^n \{y_i = 1\}} = \frac{TP}{TP + FN}$$

$$F - \text{score} = 2 * \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

ROC – кривая и AUC



X – результат работы алгоритма, Y – истинная метка

	Y = 1	Y = 0
X = 1	True Positive (TP)	False Positive (FP) (Type I error)
X = 0	False Negative (FN) (Type II error)	True Negative (TN)

$$TPR = \frac{TP}{TP + FN} = Recall$$

$$FPR = \frac{FP}{FP + TN}$$

Метрики задачи классификации

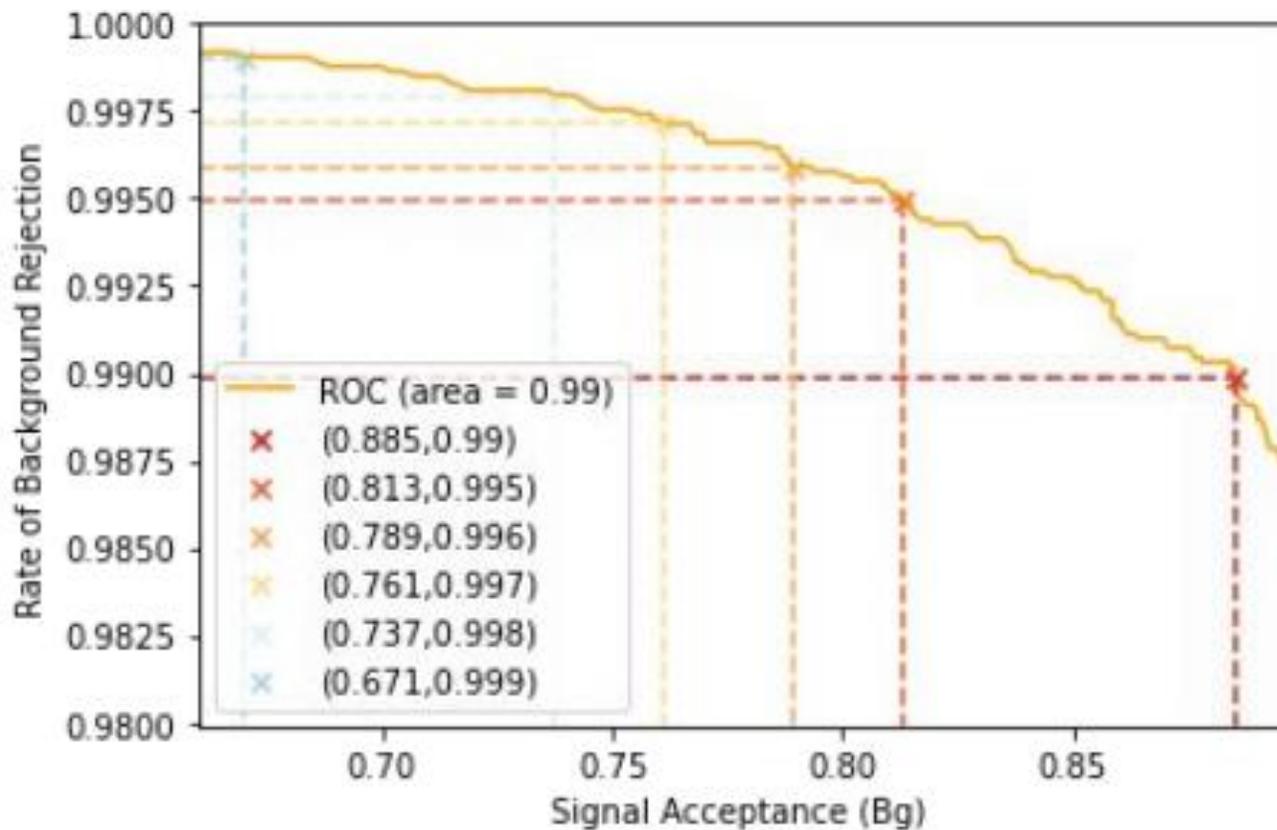
Оценка
качества



ROC – кривая и AUC

$$BRR = 1 - FPR = 1 - \frac{FP}{FP + TN} = \frac{TN}{FP + TN}$$

$$SA = TPR = \frac{TP}{TP + FN} = \text{Recall}$$



	Y = 1	Y = 0
X = 1	True Positive (TP)	False Positive (FP) (Type I error)
X = 0	False Negative (FN) (Type II error)	True Negative (TN)

X – результат работы алгоритма, Y – истинная метка

Метрики задачи регрессии

Оценка
качества



$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - x_i|$$

y – истинное значение целевой переменной
 x – предсказанное значение целевой переменной

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - x_i)^2$$

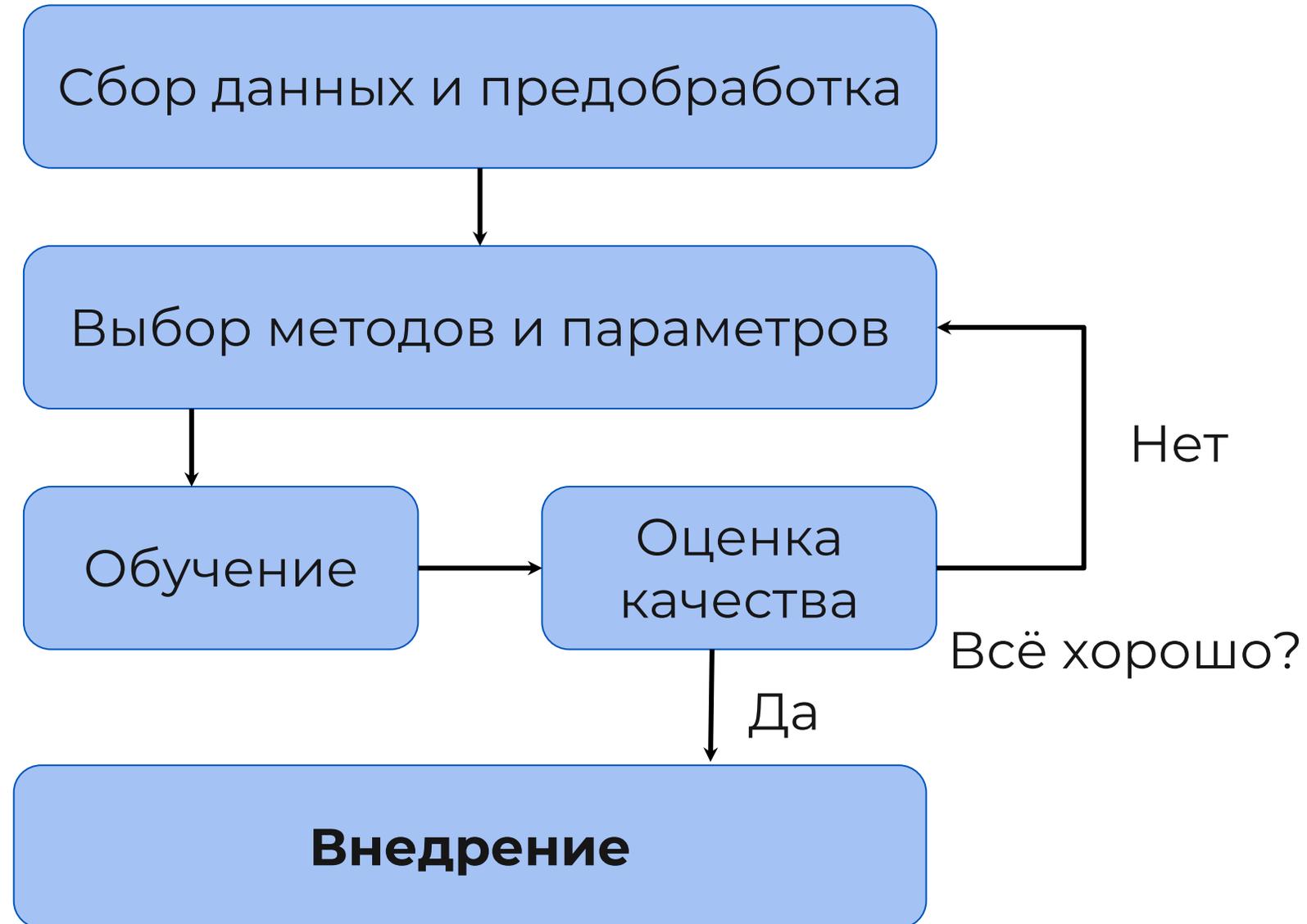
$$R^2 \equiv 1 - \frac{\sum_i (y_i - x_i)^2}{\sum_i (y_i - \bar{y})^2}$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - x_i)^2}$$

$$MAPE = \frac{100\%}{n} \sum_{i=1}^n \left| \frac{y_i - x_i}{y_i} \right|$$

$$\text{My metric} = \frac{1}{n} \sum \dots$$

Ход работы



Развертывание модели

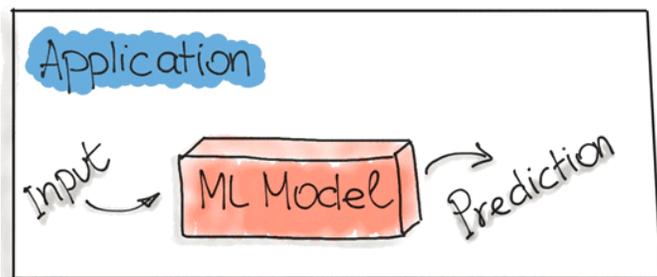
Внедрение



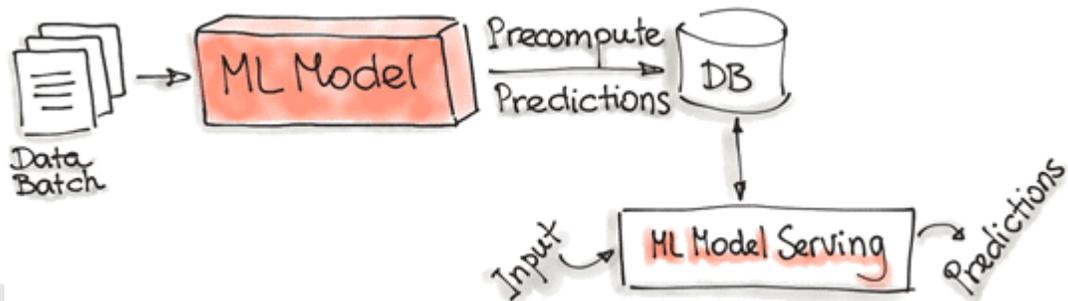
Модель как услуга



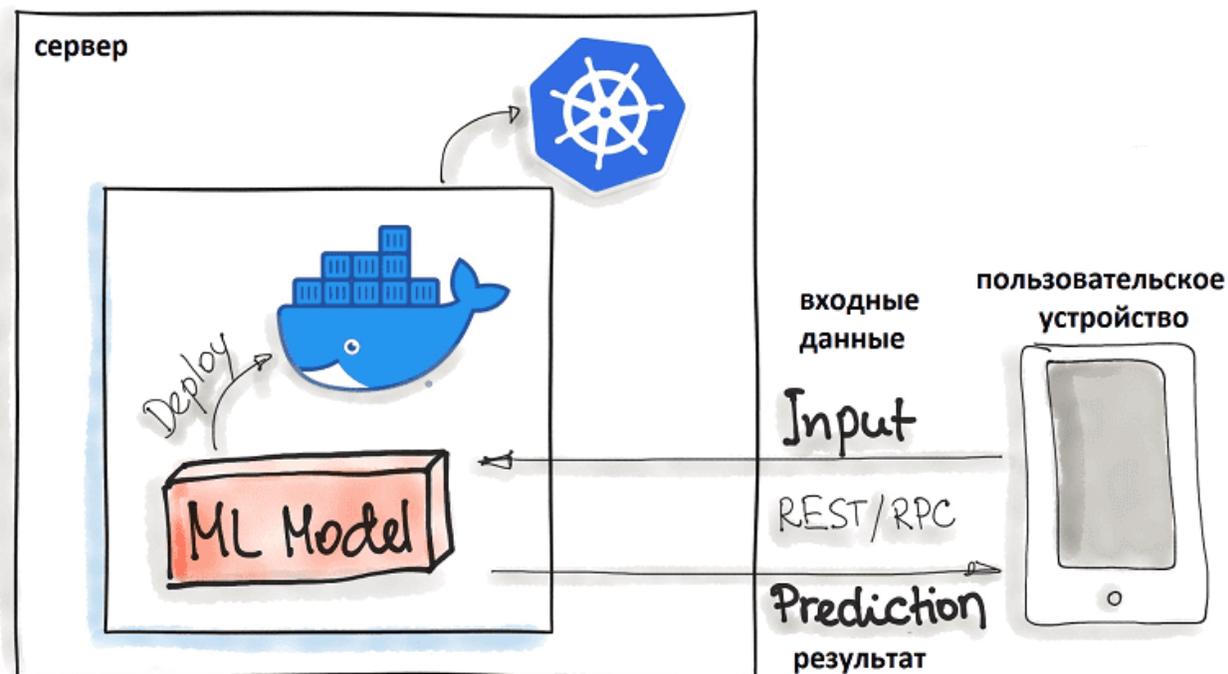
Модель как зависимость



Предварительный расчет



Развертывание с помощью Docker-контейнеров



**Конец 1 части.
Вопросы?**



Методы статистического анализа

Часть 2: Модели



27.10.2022

План второй части

Будут рассмотрены следующие модели

- Наивный байесовский классификатор
- Дерево решений
- Метод опорных векторов
- Регрессия
- Логистическая регрессия
- Принцип главных компонент

План второй части

Будут рассмотрены следующие модели

- **Наивный байесовский классификатор**
- Дерево решений
- Метод опорных векторов
- Регрессия
- Логистическая регрессия
- Принцип главных компонент

Наивная байесовская классификация

привет...	1829
валера ...	1710
нет ...	1191
куда ...	1012
небо ...	985
огурцы ...	873
говорить...	747
третий ...	739

нормальные
письма

виагра ...	1552
казино ...	1492
100% ...	1320
кредит...	1184
скидка ...	985
нажми ...	873
free ...	747
доход ...	739

спам-письма

672 раза

«КОТИК»

13 раз

Простейший спам-фильтр

(использовались года до 2010)

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

формула Байеса



не спам

Наивный Байес

Наивная байесовская классификация



Алгоритм работы

Пусть имеется набор данных в формате «Погодные условия» - «Игра». Задача: на основе погодных условий спрогнозировать возможность проведения матча. Для этого необходимо проделать следующие шаги:

- **Преобразовать** набор данных в частотную таблицу.
- Создать **таблицу** правдоподобия, рассчитав соответствующие вероятности.

Например: $P(\text{облачно}) = 0,37$, $P(\text{игра состоится}) = 0.55$

- Рассчитать **апостериорную** вероятность для каждого класса по теореме Байеса:

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Задача: состоится ли матч при солнечной погоде?

$$P(\text{состоится}|\text{солнечно}) = \frac{P(\text{солнечно}|\text{состоится}) P(\text{состоится})}{P(\text{солнечно})}$$

Проблема: отравление Байеса

План второй части

Будут рассмотрены следующие модели

- Наивный байесовский классификатор
- **Дерево решений**
- Метод опорных векторов
- Регрессия
- Логистическая регрессия
- Принцип главных компонент

Дерево принятия решений

Давать ли кредит?

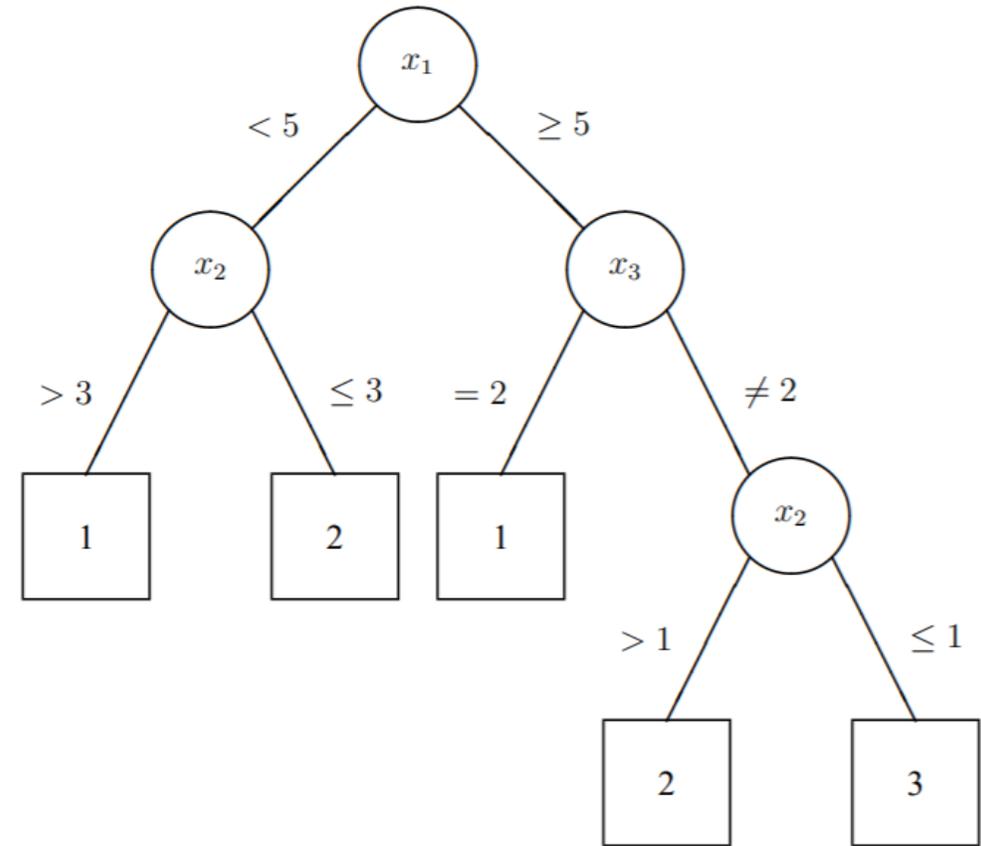


Дерево принятия решений

Алгоритм построения

Построение дерева складывается из следующих шагов:

- **Цикл** по всем признакам:
 - Цикл по упорядоченным **значениям** каждого признака:
 - **Поиск** наилучшего разделения среди всех значений этого признака на основании некоторого критерия.
 - **Разбиение** на левую и правую ветвь.
 - **Повторение** процедуры внутри каждой из ветвей вплоть до срабатывания критерия остановки.

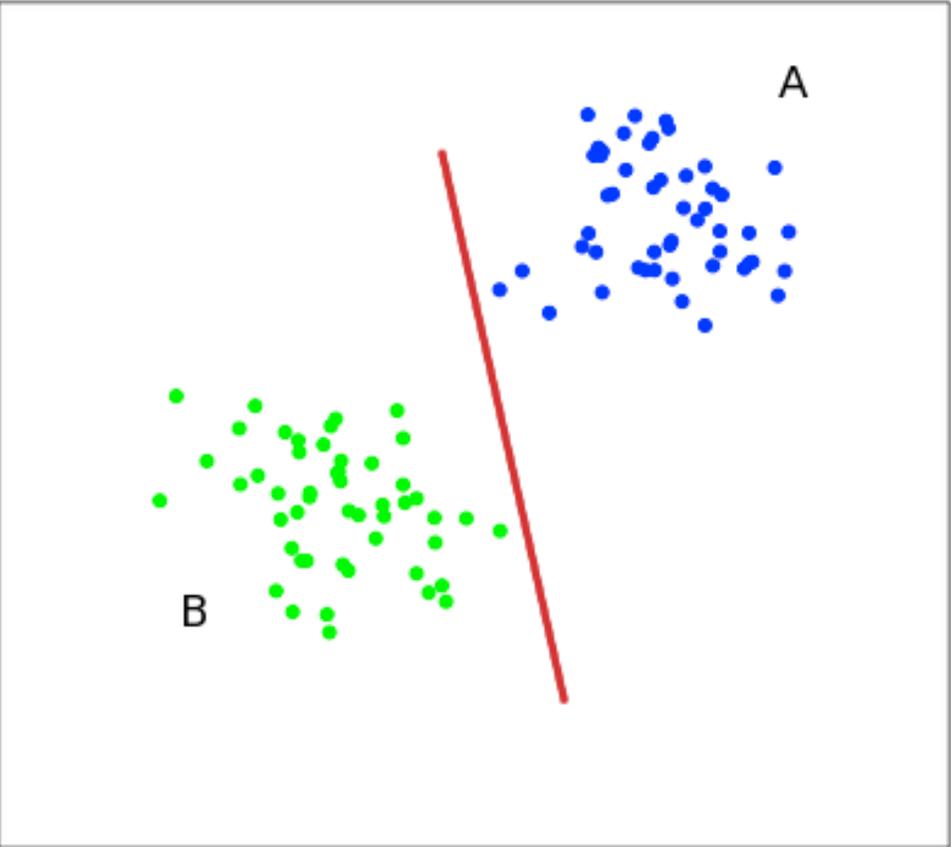


План второй части

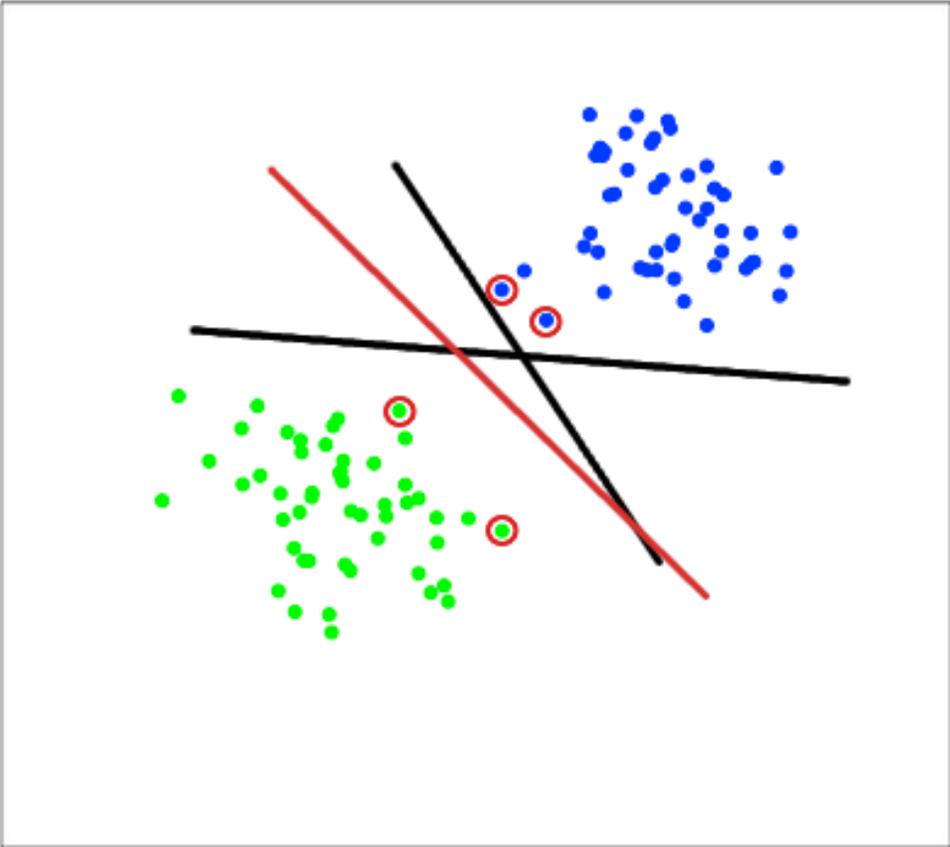
Будут рассмотрены следующие модели

- Наивный байесовский классификатор
- Дерево решений
- **Метод опорных векторов**
- Регрессия
- Логистическая регрессия
- Принцип главных компонент

Метод опорных векторов



?



Метод опорных векторов

Принцип работы

Пусть есть обучающая выборка $(x_1, y_1), \dots, (x_m, y_m), x_i \in R^n, y_i \in \{-1, 1\}$

$F(\vec{x}) = \text{sign}((\vec{w}, \vec{x}) + b)$, где

w – нормальный вектор к разделяющей (гипер)плоскости,

b – вспомогательный вектор

$F(\vec{x}) = 1$ или -1 в зависимости от класса события

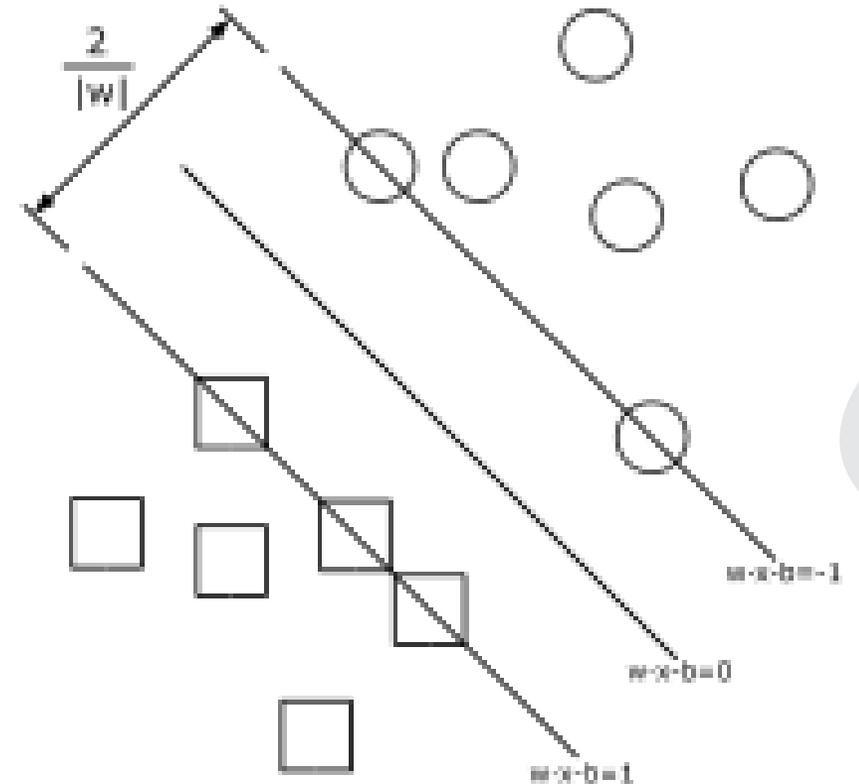
Метод опорных векторов

Принцип работы

- Максимальное расстояние до опорных векторов
- Точки не должны попасть в область полосы

$$\begin{cases} \mathit{arg} \min_{\vec{w}, b} \|\mathbf{w}\|^2, \\ y_i((\vec{w}, \vec{x}) + b) \geq 1, i = \overline{1, m} \end{cases} \Rightarrow$$

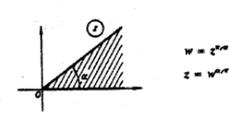
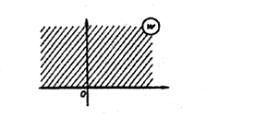
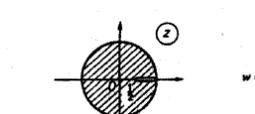
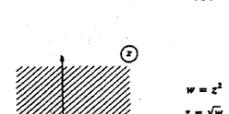
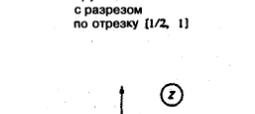
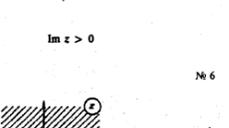
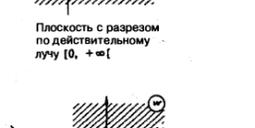
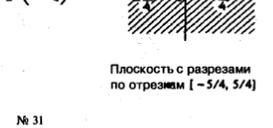
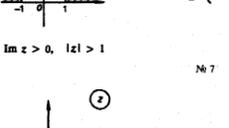
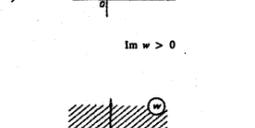
метод множителей Лагранжа

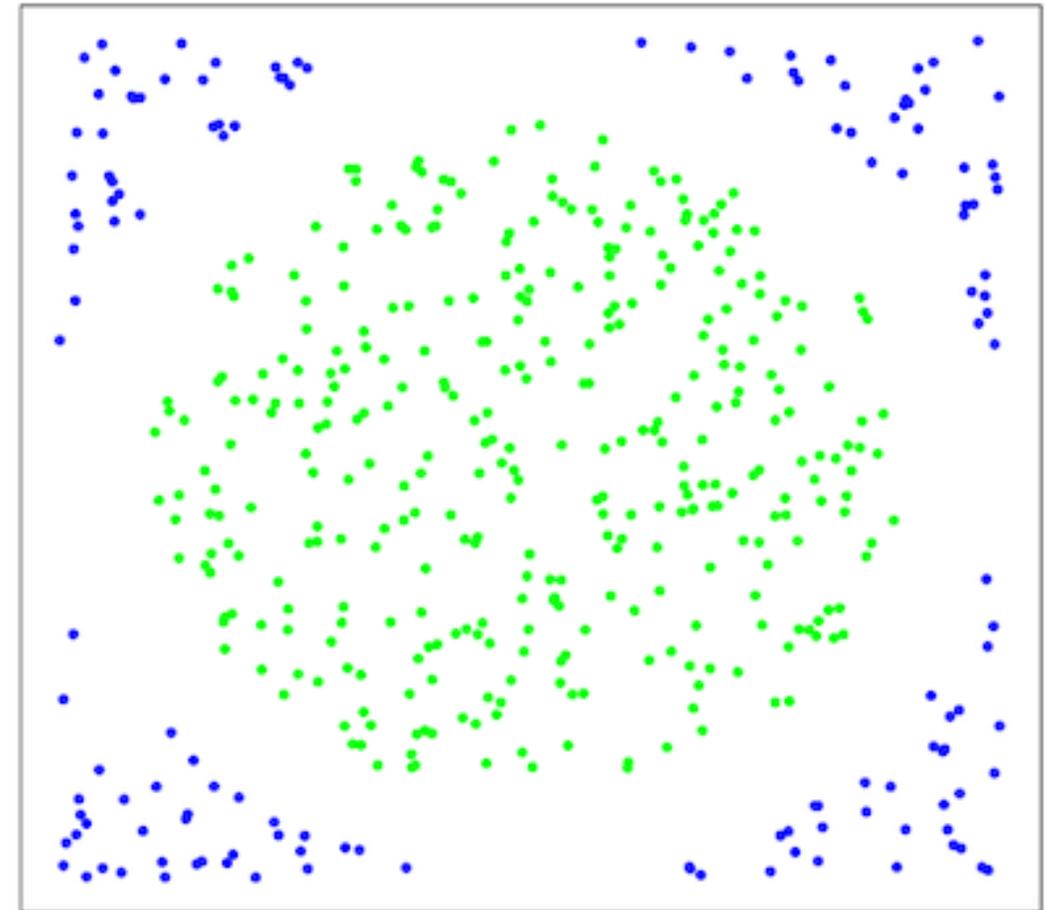


Метод опорных векторов

Вопрос: А как быть если данные не совсем хорошо линейно разделимы?

Ответ: обратиться к ТФКП

 <p>№ 5</p>	$w = z^{2n}$ $z = w^{1/2n}$  <p>№ 5</p>	 <p>№ 30</p>	$w = \frac{1}{2} \left(z + \frac{1}{z} \right)$  <p>№ 30</p>
 <p>№ 6</p>	$w = z^2$ $z = \sqrt{w}$  <p>№ 6</p>	 <p>№ 31</p>	$w = \frac{1}{2} \left(z + \frac{1}{z} \right)$  <p>№ 31</p>
 <p>№ 7</p>	$w = \frac{1}{2} \left(z + \frac{1}{z} \right)$  <p>№ 7</p>	 <p>№ 32</p>	$w = \frac{z}{i}$  <p>№ 32</p>
 <p>№ 8</p>	$w = \frac{1}{2} \left(z + \frac{1}{z} \right)$  <p>№ 8</p>	<p>Внешность круга с разрезами</p>	

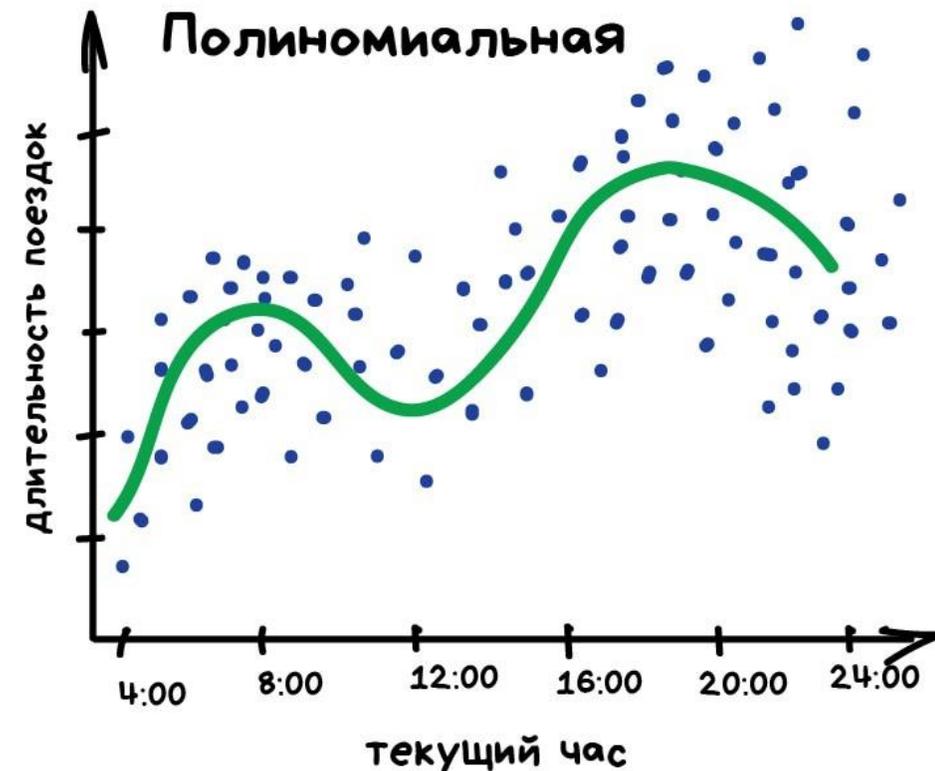
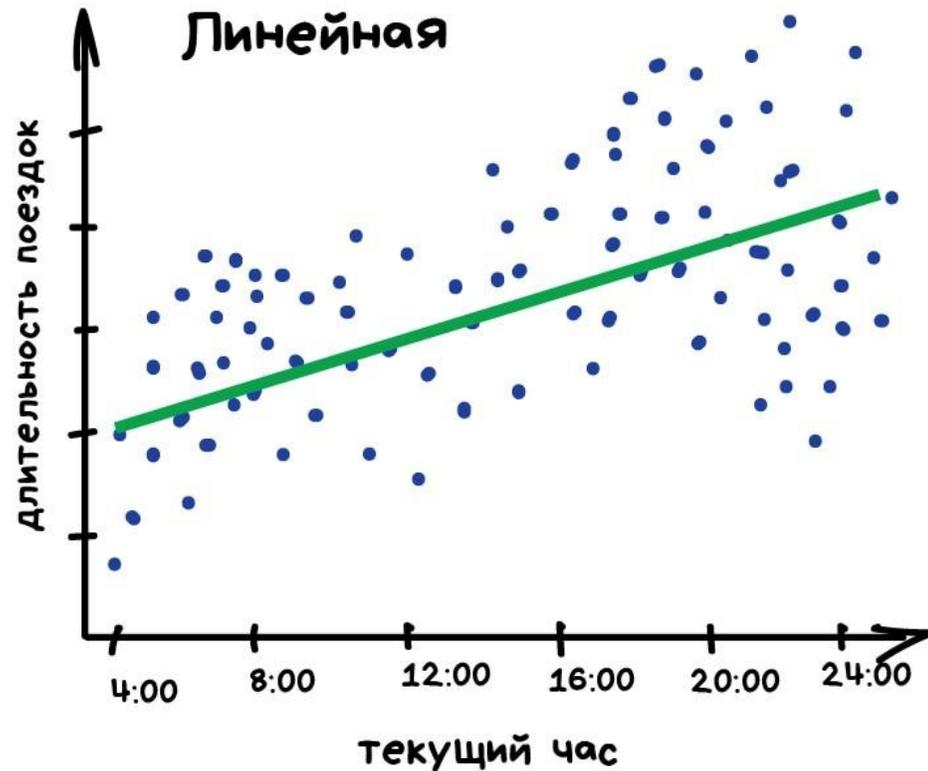


План второй части

Будут рассмотрены следующие модели

- Наивный байесовский классификатор
- Дерево решений
- Метод опорных векторов
- **Регрессия**
- Логистическая регрессия
- Принцип главных компонент

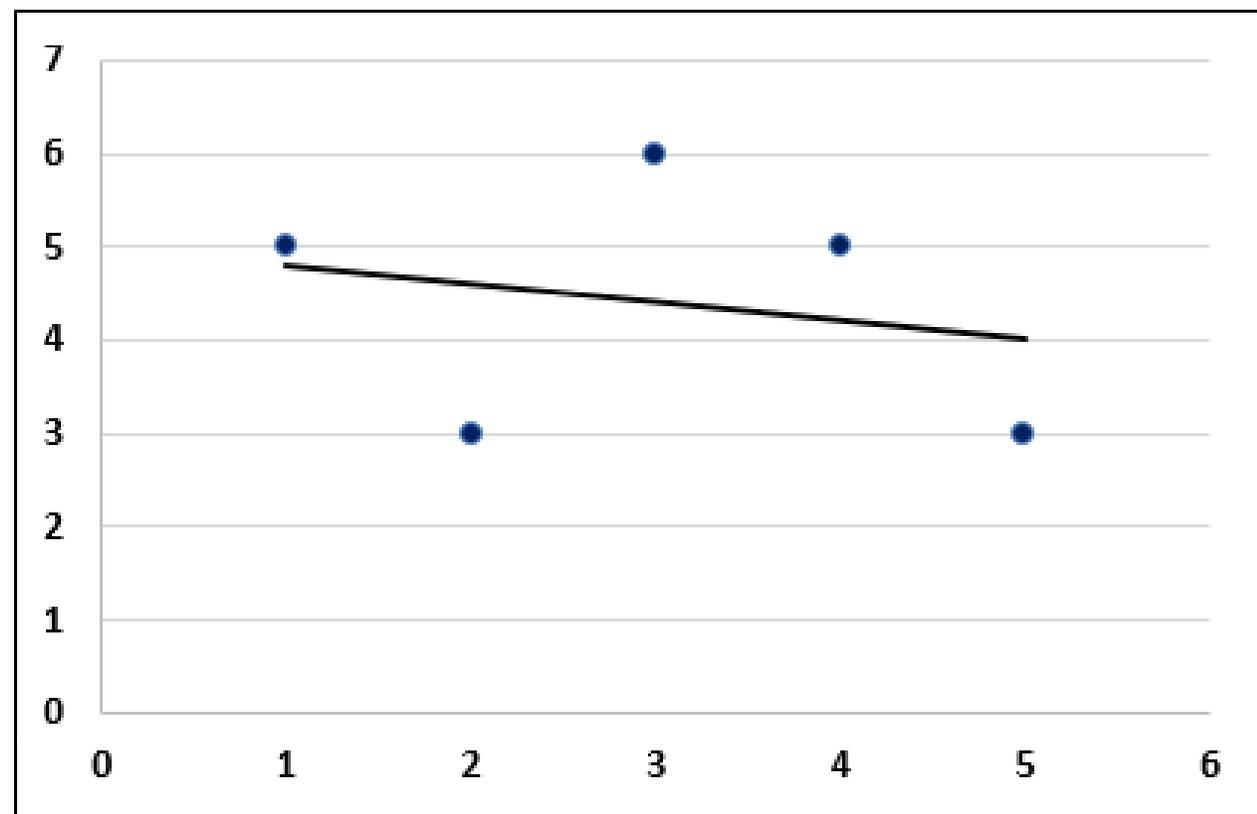
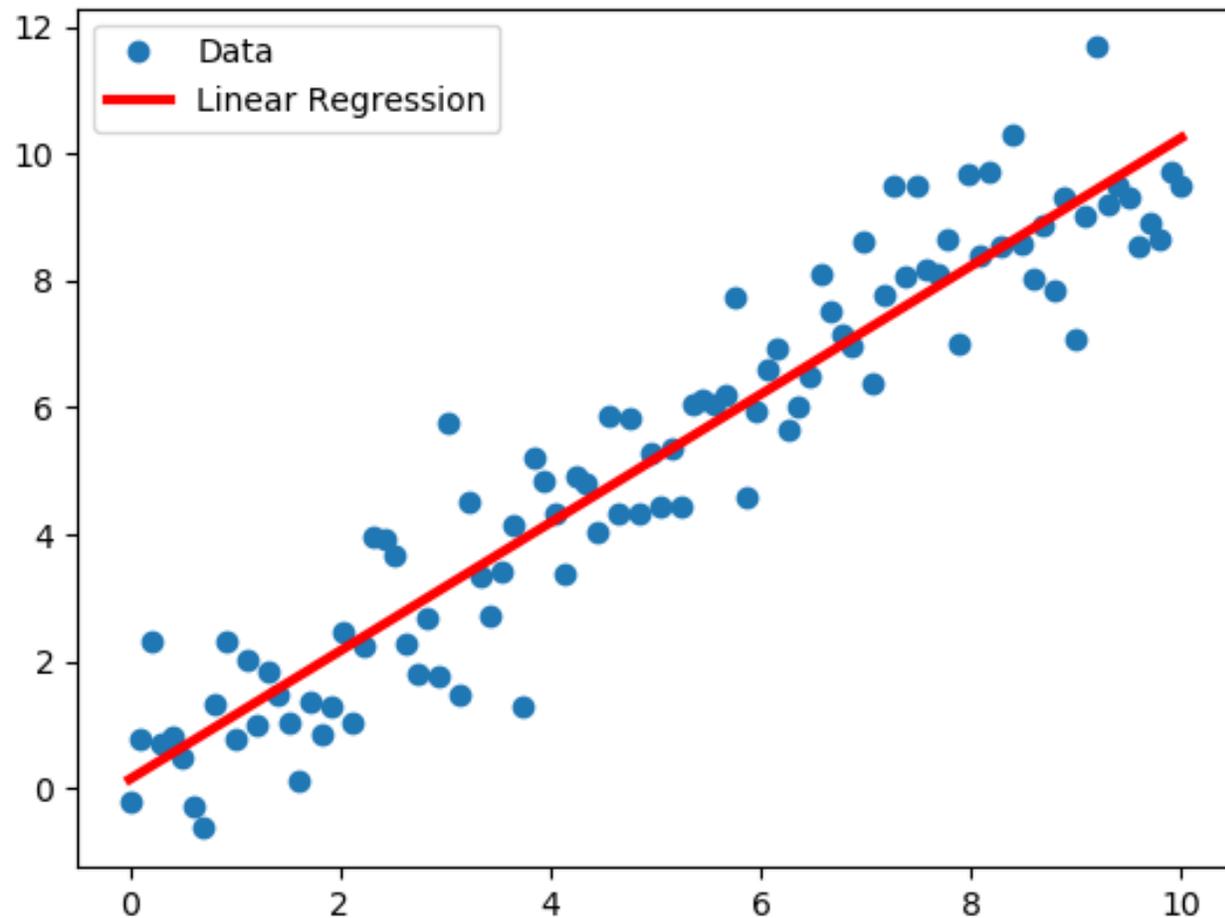
Предсказываем пробки



Вопрос: Что Вы знаете о линейной регрессии?

Регрессия

Метод наименьших квадратов



Метод наименьших квадратов

Алгоритм решения

Задача: минимизация функции $F(\mathbf{a}, \mathbf{b}) = \sum_{i=1}^n (y_i - (ax_i + b))^2$

Решение:
$$\begin{cases} \frac{\partial F(\mathbf{a}, \mathbf{b})}{\partial a} = 0 \\ \frac{\partial F(\mathbf{a}, \mathbf{b})}{\partial b} = 0 \end{cases} \Leftrightarrow \begin{cases} -2 \sum_{i=1}^n (y_i - (ax_i + b))x_i = 0 \\ -2 \sum_{i=1}^n (y_i - (ax_i + b)) = 0 \end{cases} \Leftrightarrow \begin{cases} a \sum_{i=1}^n x_i^2 + b \sum_{i=1}^n x_i = \sum_{i=1}^n x_i y_i \\ a \sum_{i=1}^n x_i + \sum_{i=1}^n b = \sum_{i=1}^n y_i \end{cases}$$

Ответ:
$$\begin{cases} a = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \\ b = \frac{\sum_{i=1}^n y_i - a \sum_{i=1}^n x_i}{n} \end{cases}$$

Полиномиальная регрессия

Задача: найти коэффициенты β_i в разложении:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_m x_i^m + \epsilon_i$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^m \\ 1 & x_2 & x_2^2 & \dots & x_2^m \\ 1 & x_3 & x_3^2 & \dots & x_3^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^m \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_m \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_n \end{bmatrix} \Leftrightarrow \vec{y} = X\vec{\beta} + \vec{\epsilon}$$

Тогда примерно коэффициенты полиномиальной регрессии

можно оценить как: $\hat{\vec{\beta}} = (X^T X)^{-1} X^T \vec{y}$

План второй части

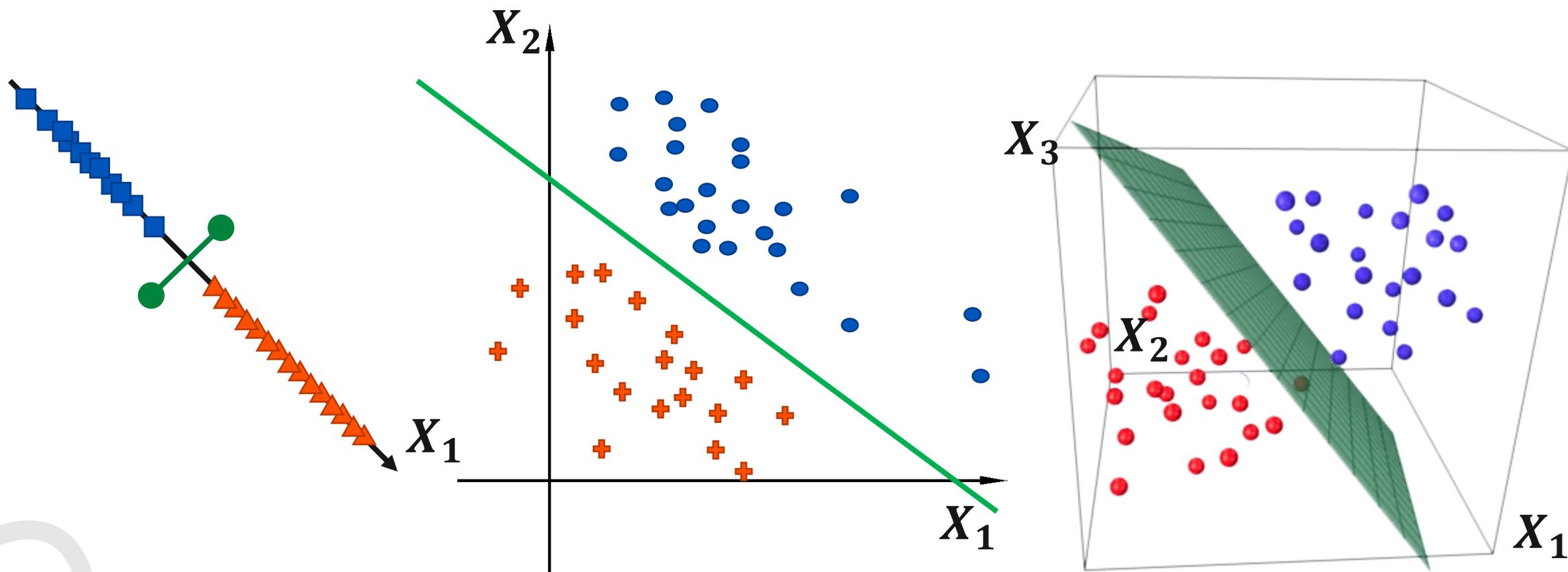
Будут рассмотрены следующие модели

- Наивный байесовский классификатор
- Дерево решений
- Метод опорных векторов
- Регрессия
- **Логистическая регрессия**
- Принцип главных компонент

Логистическая регрессия классификация

Пусть есть 2 класса событий: **0** и **1**. P_0 и P_1 – вероятность принадлежности к определенному классу.

$$P_1 = 1 - P_0; P_1, P_2 \in [0, 1]$$



Логистическая регрессия

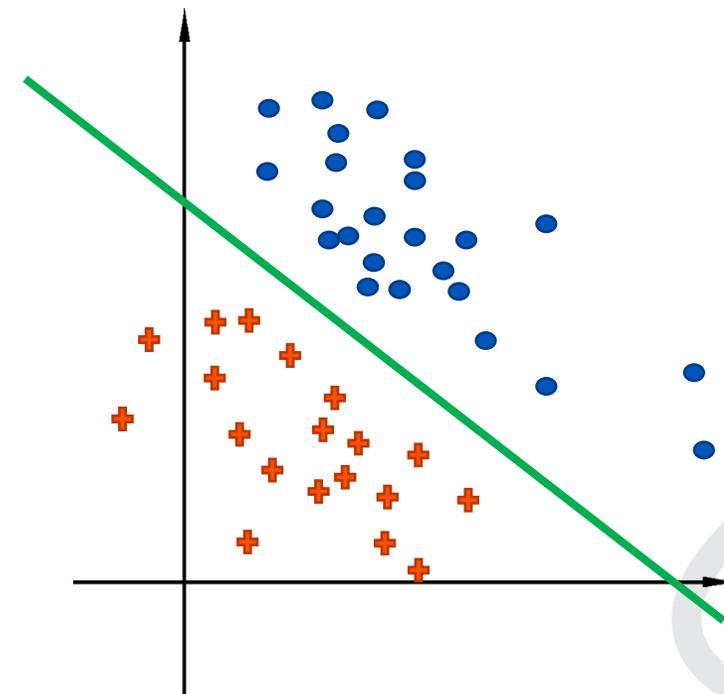
Рассмотрим задачу в двумерном случае

Пусть найдено уравнение прямой: $\beta_0 + \beta_1 x_1 + \beta_2 x_2 = 0$

Рассмотрим точку (a, b) . Подставим ее в уравнение прямой:

$t = \beta_0 + \beta_1 a + \beta_2 b$, причем t может принимать 3 значения:

$$\left[\begin{array}{l} t < 0, \text{ тогда } (a, b) \text{ принадлежит классу } 0, P_1 \in [0, 0.5) \\ t > 0, \text{ тогда } (a, b) \text{ принадлежит классу } 1, P_1 \in (0.5, 1] \\ t = 0, \text{ тогда } (a, b) \text{ лежит на прямой, } P_1 = 0.5 \end{array} \right.$$



Логистическая регрессия

Имеем: функция с выходным значением $(-\infty, +\infty)$.

Вопрос: как перевести выходное значение в вероятность?

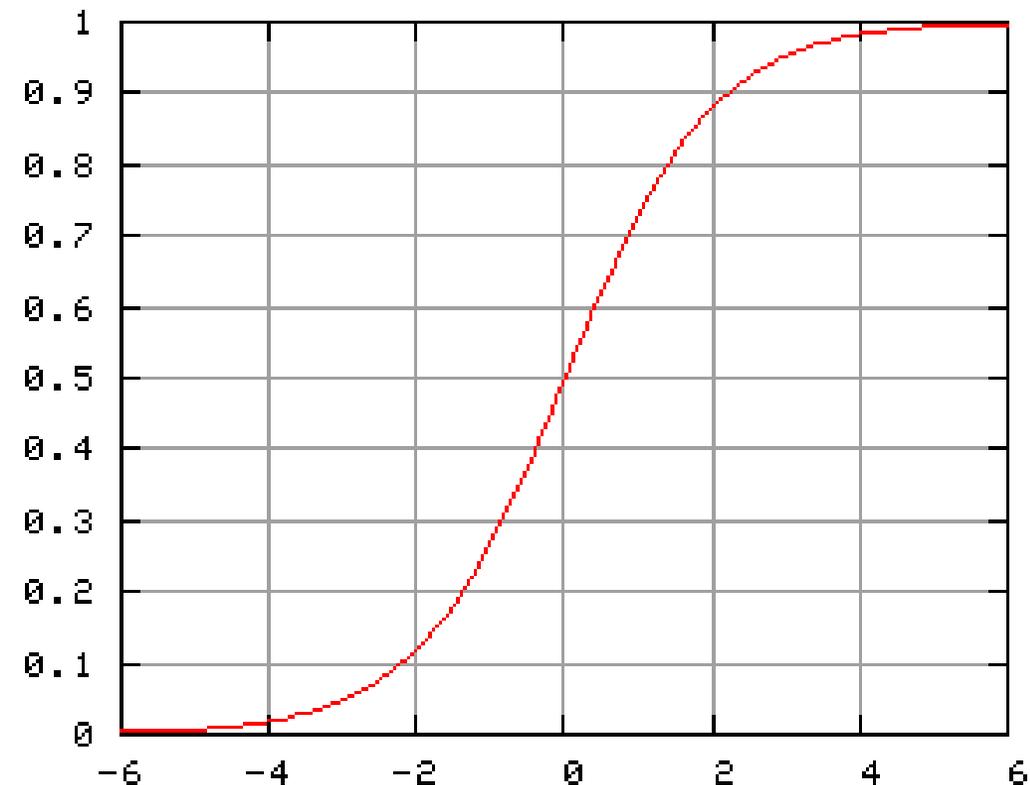
$$OR(X) = \frac{P(X)}{1 - P(X)} \in [0, +\infty)$$

$$\log(OR(X)) = \log\left(\frac{P(X)}{1 - P(X)}\right) \in (-\infty, +\infty)$$

Приняв $\log(OR(X)) = t, t \in (-\infty, +\infty)$

Получаем: $e^t = OR(X) = \frac{P(X)}{1 - P(X)}$, тогда

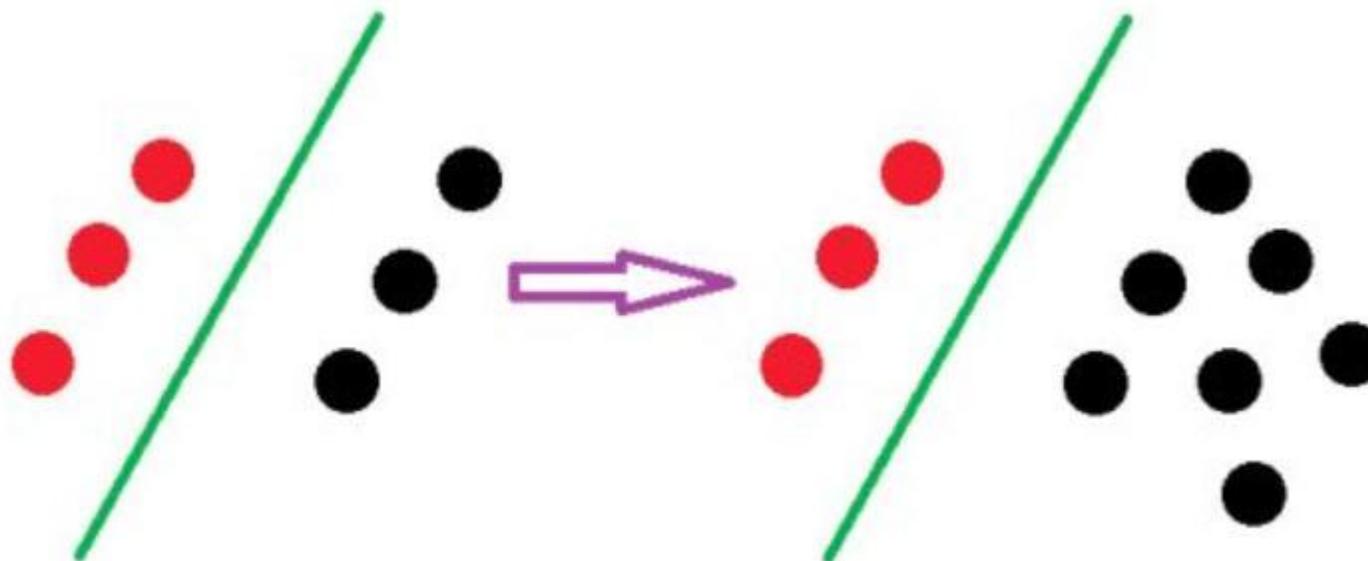
$$P(X) = \frac{e^t}{1 + e^t} = \frac{1}{1 + e^{-t}} - \text{логистическая функция}$$



LR vs. SVM

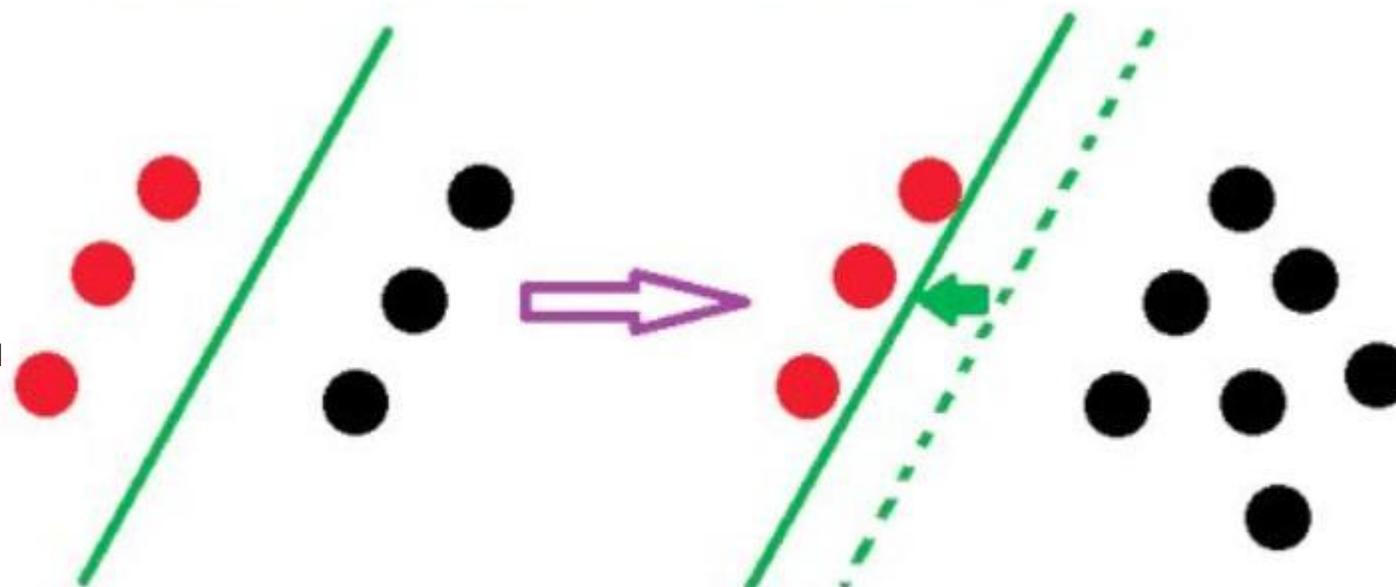
SVM

Геометрические свойства



LR

Статистические подходы

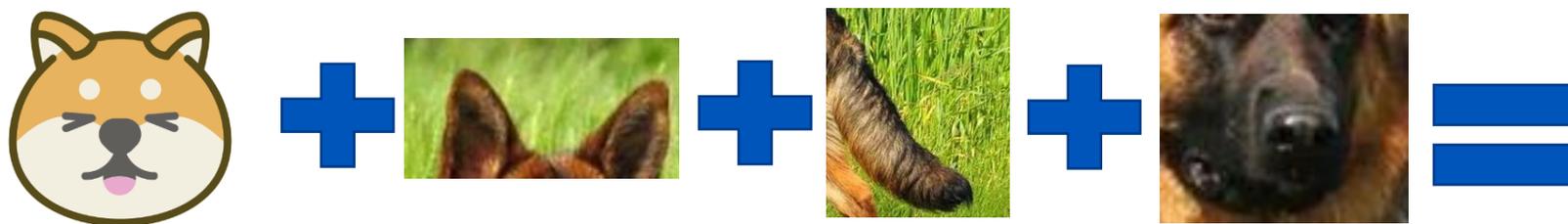


План второй части

Будут рассмотрены следующие модели

- Наивный байесовский классификатор
- Дерево решений
- Метод опорных векторов
- Регрессия
- Логистическая регрессия
- **Принцип главных компонент**

Принцип уменьшения размерности



Овчарка

Метод главных компонент

Принцип работы алгоритма

Пусть имеется матрица переменных X ($I \times J$), где I – число образцов, J – число признаков, $J \gg I$.

Обозначим $t_\alpha = p_{\alpha 1} x_1 + \dots + p_{\alpha J} x_J$ – линейная комбинация исходных переменных x_j .

$e_{(G)}$	1	2	3	4	5	6	7	8	$v_{(G)}$
$I =$	1	1	1	0	0	0	0	0	1
	1	0	0	1	1	0	0	0	2
	0	0	1	0	0	1	1	0	3
	0	0	0	1	0	0	0	0	4
	0	1	0	0	1	1	0	1	5
	0	0	0	0	0	0	1	1	6

Метод главных компонент

$$t_{\alpha} = p_{\alpha 1} x_1 + \dots + p_{\alpha J} x_J$$

Тогда исходную матрицу X можно разложить следующим образом:

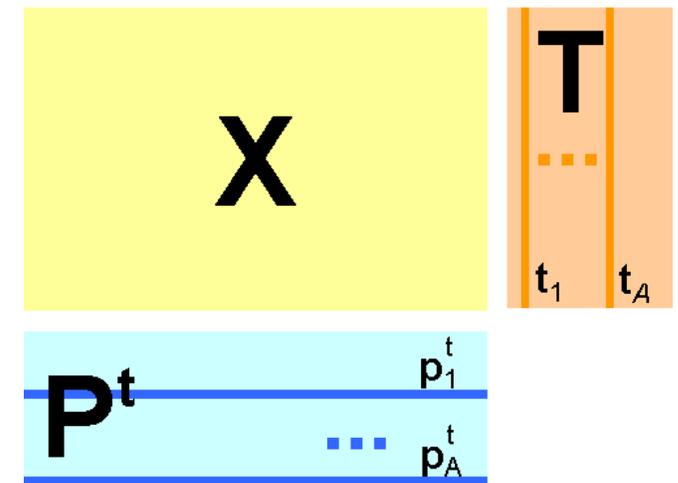
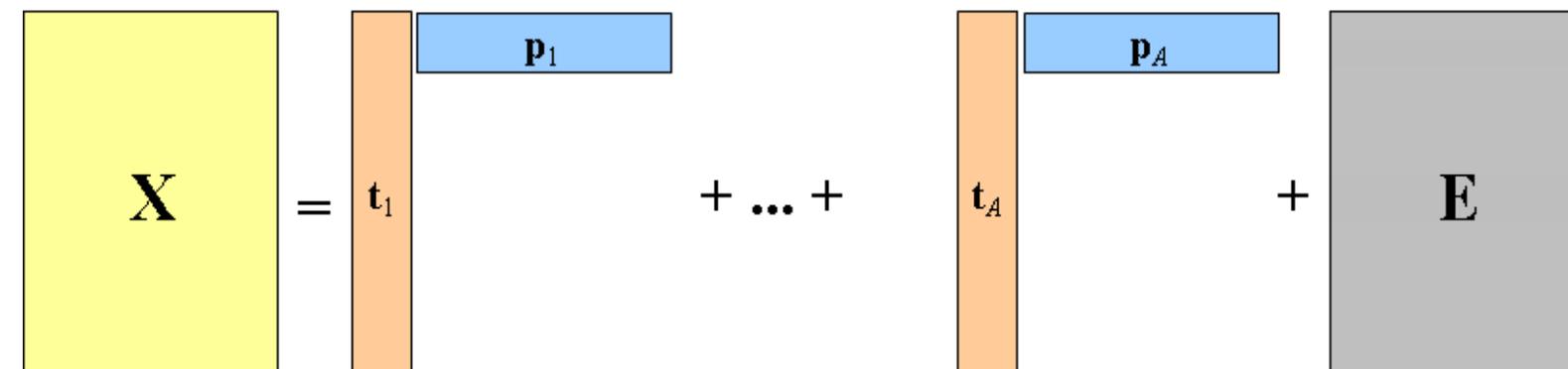
$$X = TP^T + E = \sum_{\alpha=1}^A t_{\alpha} p_{\alpha}^T + E$$

Матрица \mathbf{T} называется матрицей *счетов* (scores). Ее размерность $(I \times A)$.

Матрица \mathbf{P} называется матрицей *нагрузок* (loadings). Ее размерность $(J \times A)$.

\mathbf{E} – это матрица *остатков*, размерностью $(I \times J)$.

Важное свойство – **ортогональность**.



**Конец 2 части.
Вопросы?**





ВСЁ!



Источники



Красивые картиночки:

https://vas3k.ru/blog/machine_learning/

Внедрение:

<https://www.bigdataschool.ru/blog/mlops-deployment-patterns-and-strategies.html>

Регрессионные модели:

<https://habr.com/ru/company/ods/blog/323890/>

Различие SVM vs. LR:

<https://progler.ru/blog/razlichiya-mezhdu-mashinoy-opornyh-vektorov-i-logisticheskoy-regressiey>