

Введение в анализ данных

Лекция 3

Метод k ближайших соседей и линейная регрессия

Евгений Соколов

esokolov@hse.ru

НИУ ВШЭ, 2021

Измерение ошибки модели

Вопросы

- Как сравнить две модели?
- Как подобрать k и метрику?

Функция потерь для классификации

- Частый выбор — бинарная функция потерь

$$L(y, a) = [a \neq y]$$

- Функционал ошибки — доля ошибок (error rate)

$$Q(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} [a(x_i) \neq y_i]$$

- Нередко измеряют долю верных ответов (accuracy):

$$Q(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} [a(x_i) = y_i]$$

Функция потерь для классификации

ВАЖНО

Accuracy — не точность!

Accuracy

$a(x)$	y
-1	-1
+1	+1
-1	-1
+1	-1
+1	+1

Accuracy

$a(x)$	y
-1	-1
+1	+1
-1	-1
+1	-1
+1	+1

Доля ошибок: 0.2

Доля верных ответов: 0.8

Accuracy

$$Q(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} [a(x_i) \neq y_i]$$

Решаем задачу выявления редкого заболевания

- 950 здоровых ($y = +1$)
- 50 больных ($y = -1$)

Модель: $a(x) = +1$

Доля ошибок: 0.05

Accuracy

$$Q(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} [a(x_i) \neq y_i]$$

- Всегда смотрите на баланс классов!
- Доля верных ответов не обязательно меняется от 0.5 до 1 для разумных моделей

Как выбрать k?

Обучающая выборка

На каком классе чаще всего ездит	Ближайшее к дому метро	Способ оплаты	Согласился повысить категорию?
Эконом	Таганская	Карта	да
Комфорт	Юго-Западная	Наличные	нет
Комфорт	Строгино	Карта	да

Применяем модель:

Эконом	Таганская	Карта	?
--------	-----------	-------	---

Как выбрать k ?

Обучающая выборка

На каком классе чаще всего ездит	Ближайшее к дому метро	Способ оплаты	Согласился повысить категорию?
Эконом	Таганская	Карта	да
Комфорт	Юго-Западная	Наличные	нет
Комфорт	Строгино	Карта	да

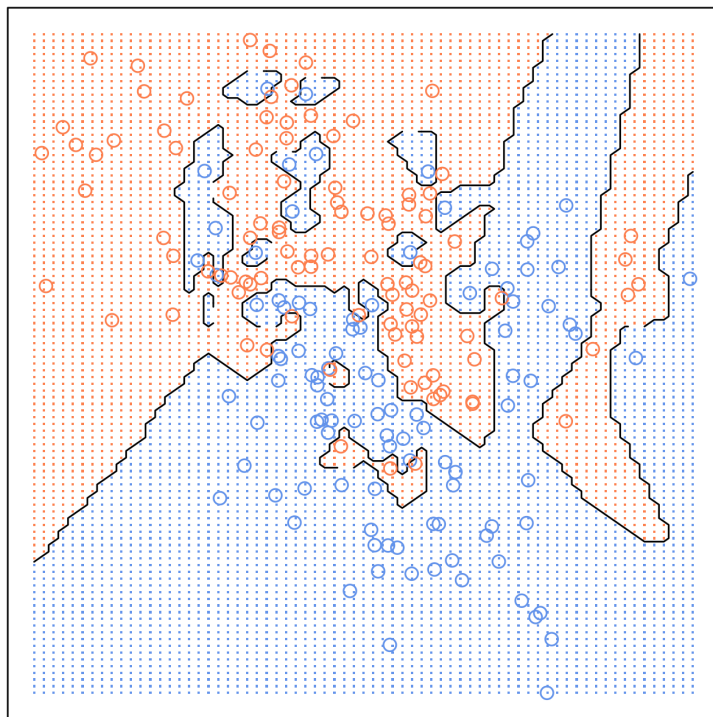
Применяем модель:

Эконом	Таганская	Карта	да
--------	-----------	-------	----

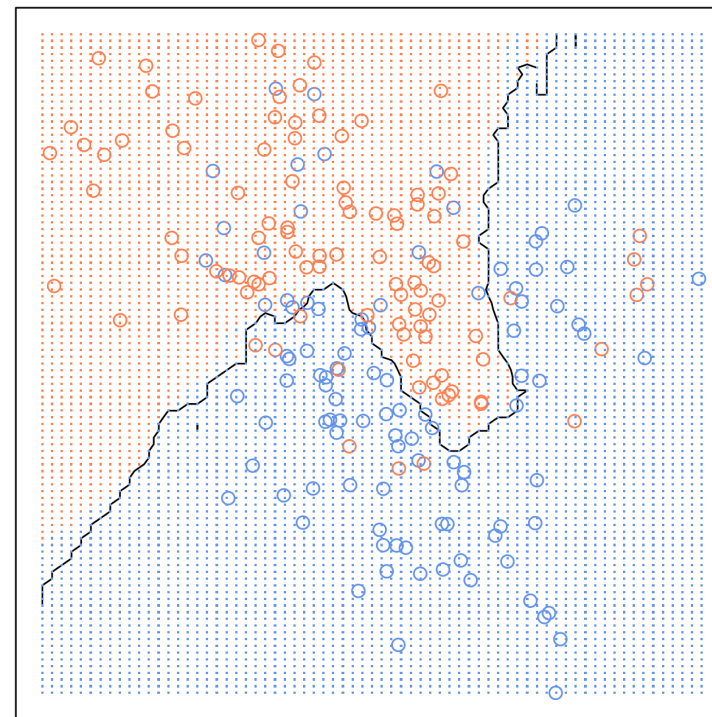
С точки зрения качества на обучающей выборке лучший выбор $k = 1$

Как выбрать k ?

1-nearest neighbours



20-nearest neighbours



<https://kevinzakka.github.io/2016/07/13/k-nearest-neighbor/>

Гиперпараметры

- Нельзя подбирать k по обучающей выборке — **гиперпараметр**
- Нужно использовать дополнительные данные

Обобщающая способность

Обобщающая способность

Как готовиться к экзамену?

Заучить все примеры с
занятий

Разобраться в предмете и
усвоить алгоритмы решения
задач

Обобщающая способность

Как готовиться к экзамену?

Заучить все примеры с
занятий

Разобраться в предмете и
усвоить алгоритмы решения
задач

Переобучение (overfitting)

Обобщение (generalization)

Обобщающая способность

Как готовиться к экзамену?

Заучить все примеры с
занятий

Разобраться в предмете и
усвоить алгоритмы решения
задач

Переобучение (overfitting)

Обобщение (generalization)

Хорошее качество на обучении
Низкое качество на новых данных

Хорошее качество на обучении
Хорошее качество на новых
данных

Отложенная выборка



Обучение



Тест

Отложенная выборка



- Слишком большое обучение — тестовая выборка нерепрезентативна
- Слишком большой тест — модель не сможет обучиться
- Обычно: 70/30, 80/20

Кросс-валидация

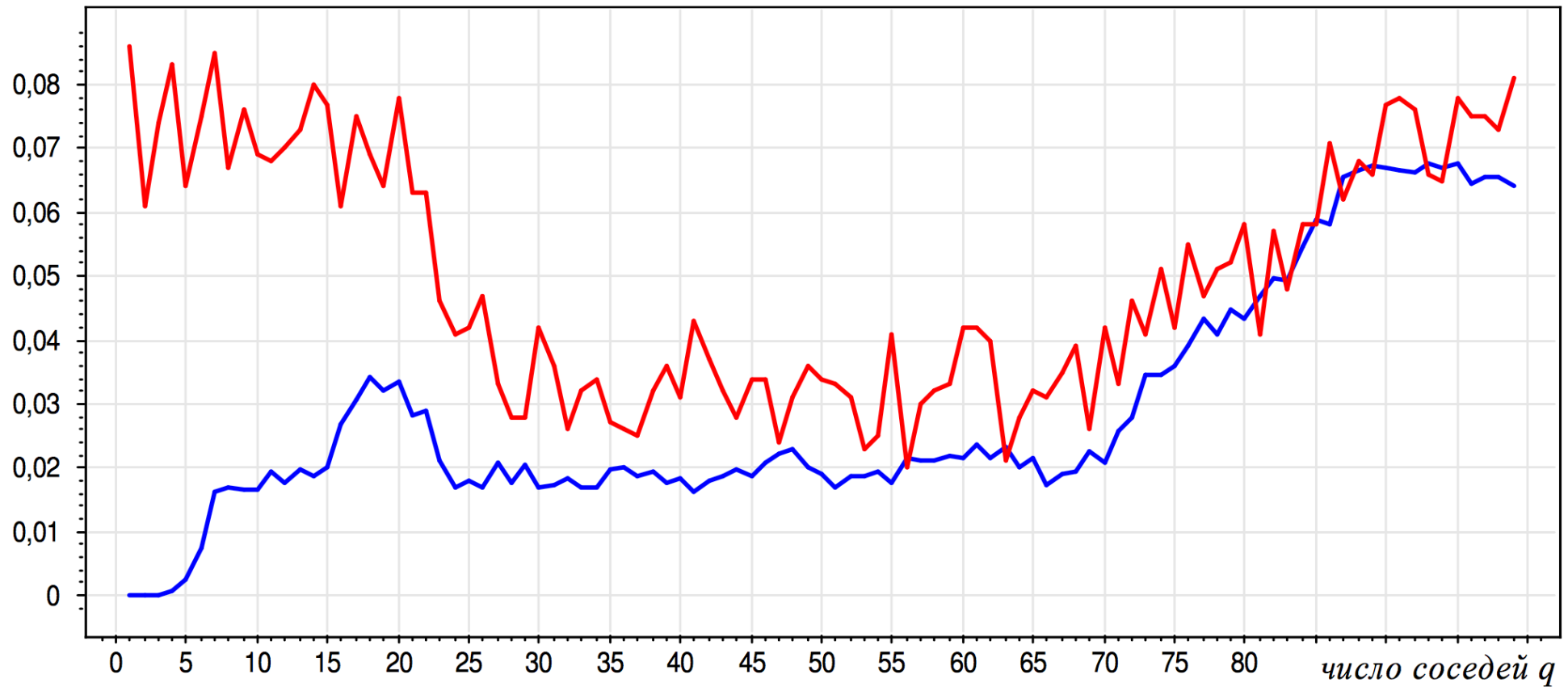


Кросс-валидация

- Надёжнее отложенной выборки, но медленнее
- Параметр — количество разбиений n (фолдов, folds)
- Хороший, но медленный вариант — $n = \ell$ (leave-one-out)
- Обычно: $n = 3$ или $n = 5$ или $n = 10$

Подбор числа соседей

частота ошибок



Чуть больше терминов

- После подбора всех гиперпараметров стоит проверить на совсем новых данных, что модель работает
- Обучающая выборка — построение модели
- Валидационная выборка — подбор гиперпараметров модели
- Тестовая выборка — финальная оценка качества модели

Метод k ближайших соседей с
весами

kNN: применение

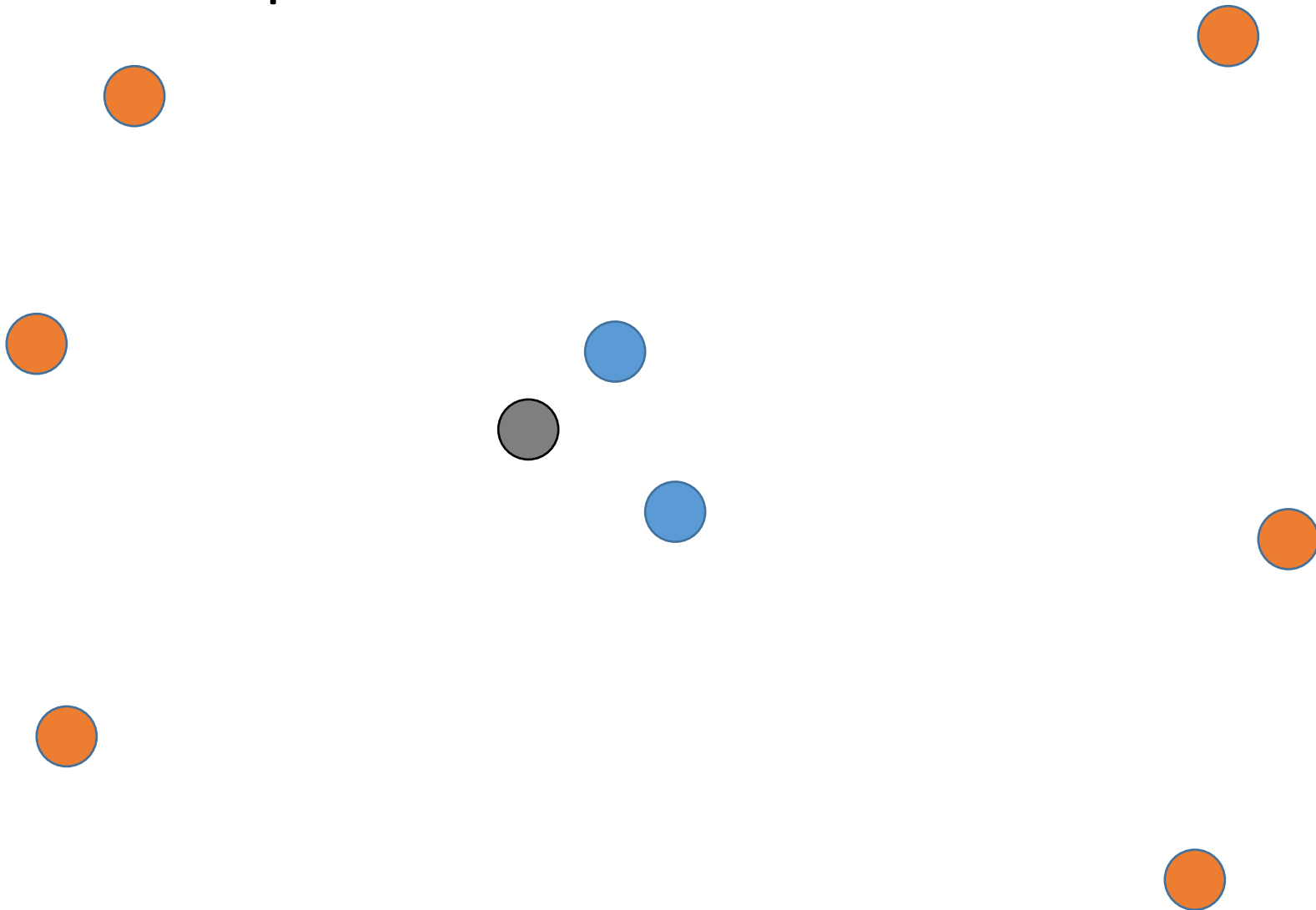
Дано: новый объект x

Применение модели:

- Сортируем объекты обучающей выборки по расстоянию до нового объекта:
 $\rho(x, x_{(1)}) \leq \rho(x, x_{(2)}) \leq \dots \leq \rho(x, x_{(\ell)})$
- Выбираем k ближайших объектов: $x_{(1)}, \dots, x_{(k)}$
- Выдаём наиболее популярный среди них класс:

$$a(x) = \arg \max_{y \in \mathbb{Y}} \sum_{i=1}^k [y_{(i)} = y]$$

Проблема с расстояниями



Взвешенный knn

$$a(x) = \arg \max_{y \in \mathbb{Y}} \sum_{i=1}^k w_i [y_{(i)} = y]$$

Варианты:

- $w_i = \frac{k+1-i}{k}$
- $w_i = q^i$
- Не учитывают сами расстояния

Взвешенный knn

$$a(x) = \arg \max_{y \in \mathbb{Y}} \sum_{i=1}^k w_i [y_{(i)} = y]$$

Парзеновское окно:

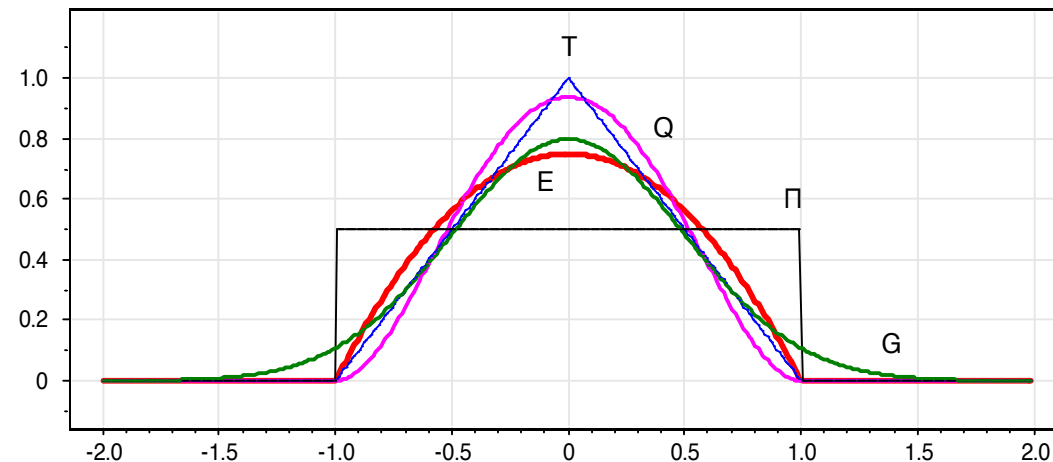
- $w_i = K \left(\frac{\rho(x, x_{(i)})}{h} \right)$
- K — ядро
- h — ширина окна

Ядра для весов

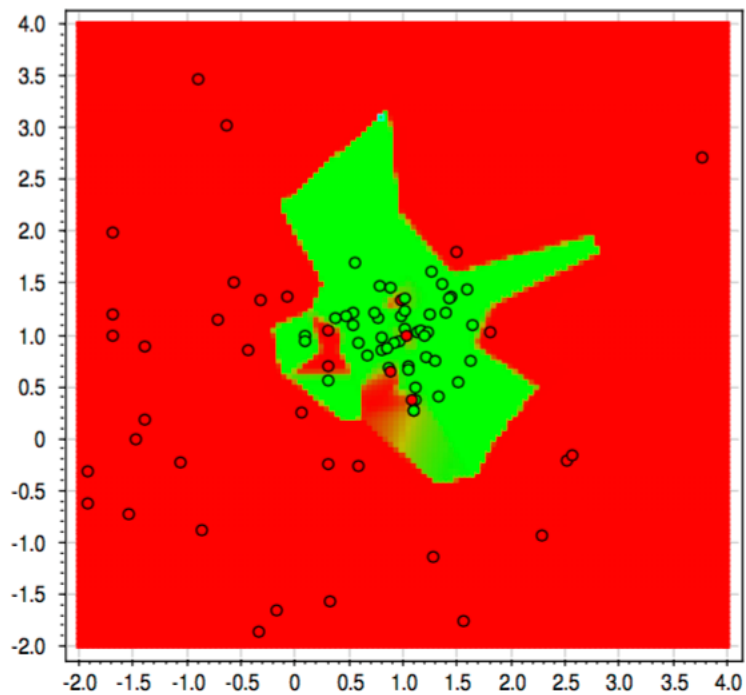
- Гауссовское ядро:

$$K(z) = (2\pi)^{-0.5} \exp\left(-\frac{1}{2}z^2\right)$$

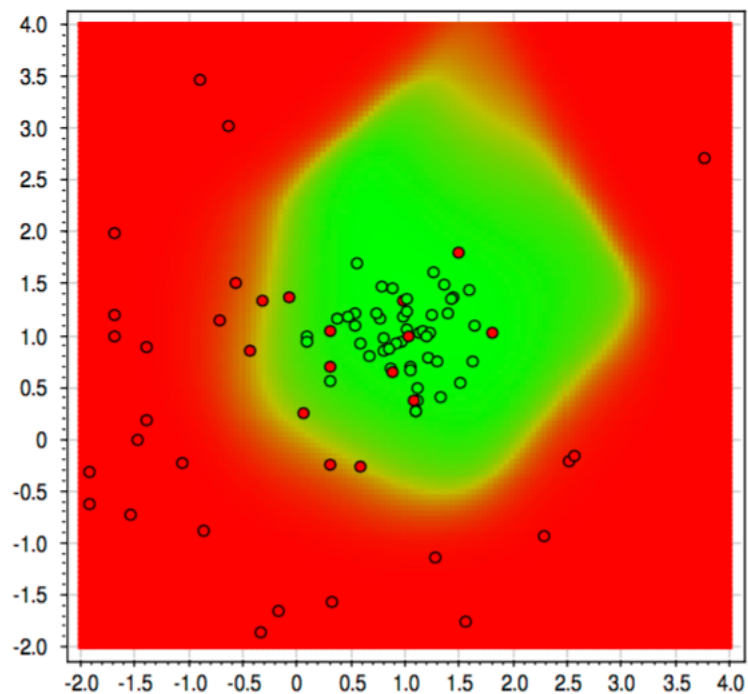
- И много других:



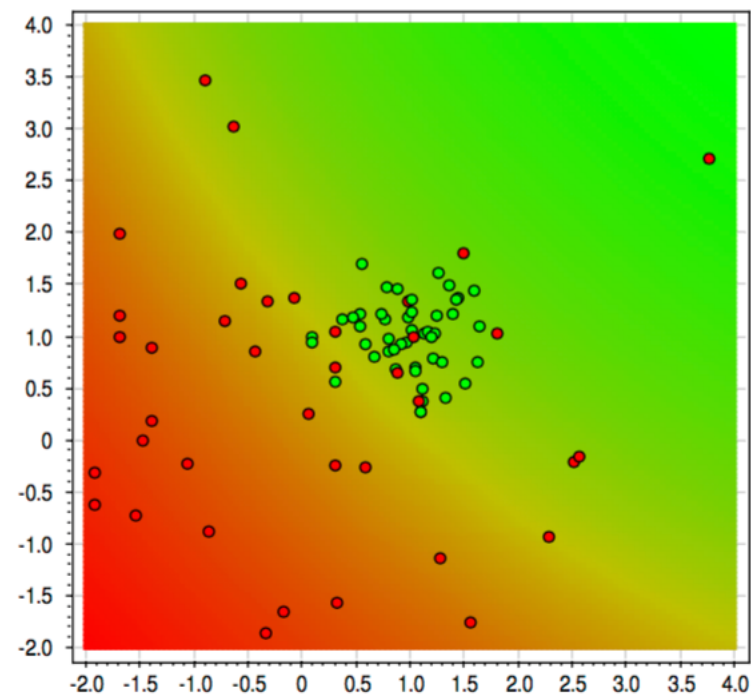
Ядра для весов



$h = 0.05$



$h = 0.5$



$h = 5$