

Введение в искусственный интеллект. Машинное обучение

Лекция 2. Непараметрические методы классификации и регрессии

Бабин Д.Н., Иванов И.Е., Петюшко А.А.

кафедра Математической Теории Интеллектуальных Систем

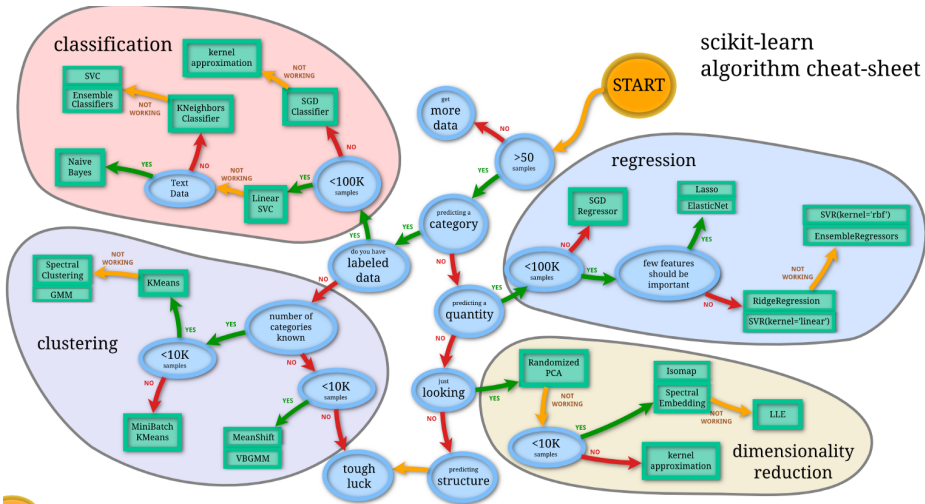
13 октября 2020 г.



- 1 Метод ближайших соседей в задаче классификации
- 2 Непараметрическая регрессия
- 3 Методы поиска ближайшего соседа

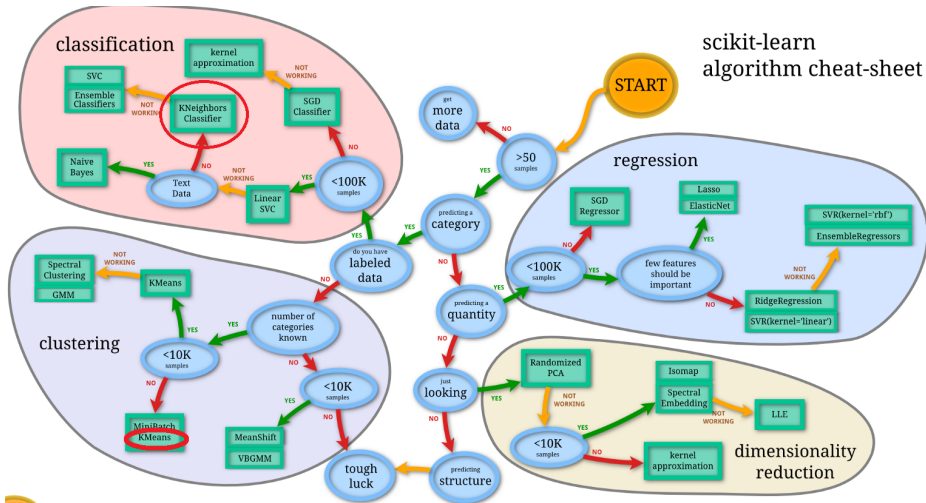


Дорожная карта Scikit-Learn¹



¹https://scikit-learn.org/stable/tutorial/machine_learning_map/

Дорожная карта Scikit-Learn¹



¹https://scikit-learn.org/stable/tutorial/machine_learning_map/

Параметрические методы

- исходят из предположения, что искомая зависимость имеет некоторый специальный вид с точностью до некоторых параметров
- параметры находятся решением оптимизационной задачи

Непараметрические методы

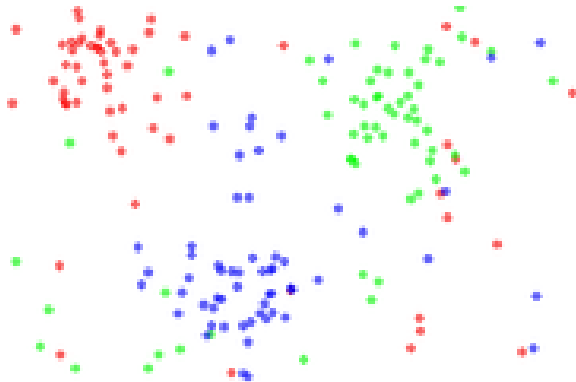
Непараметрические методы – методы не являющиеся параметрическими

- Метрические алгоритмы, ядерные методы



Основное предположение

- "Близкие" объекты лежат в одном классе
- Близость задаётся метрикой
- Типичный пример ²



²https://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm

Метод ближайшего соседа

- Параметр метода: метрика
- Алгоритм: по заданной метрике ищем ближайший объект в обучающей выборке и классифицируем объект так же

Преимущества

- Простота реализации (нет как таковой процедуры обучения в наивной реализации)
- Хорошая интерпретируемость

Недостатки

- Неустойчивость к выбросам
- Неоднозначность классификации при равных расстояниях до двух объектов
- Необходимость хранить всю обучающую выборку
- Алгоритм поиска вычислительно сложен (если обучающая выборка довольно большая)
- Не учитывается значение расстояния

Метод k ближайших соседей

- Параметр метода: метрика, k
- Алгоритм: по заданной метрике ищем k ближайших объектов в обучающей выборке и классифицируем объект как большинство из k объектов

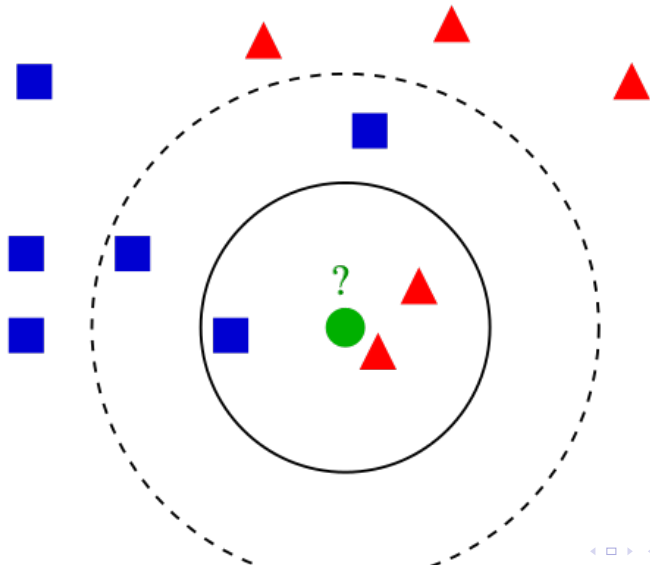
Преимущества

- Простота реализации
- Хорошая интерпретируемость
- Параметр k можно оптимизировать по скользящему контролю

Недостатки

- Неустойчивость к выбросам
- Неоднозначность классификации при равных расстояниях до двух объектов
- Необходимость хранить всю обучающую выборку
- Алгоритм поиска вычислительно сложен (если обучающая выборка довольно большая)
- Не учитывается значение расстояния

Метод k ближайших соседей



Метод k ближайших взвешенных соседей

- Параметры метода: метрика, k , веса
- Алгоритм: по заданной метрике ищем k ближайших объектов в обучающей выборке и классифицируем объект взвешенным голосованием

Преимущества

- Простота реализации
- Хорошая интерпретируемость
- Параметр k можно оптимизировать по скользящему контролю

Недостатки

- Неустойчивость к выбросам
- Неоднозначность классификации при равных расстояниях до двух объектов
- Необходимость хранить всю обучающую выборку
- Алгоритм поиска вычислительно сложен (если обучающая выборка довольно большая)
- Не учитывается значение расстояния

Метод k ближайших взвешенных соседей: выбор весов

- Веса в зависимости от порядкового номера
 - Линейно убывающие веса
 - Экспоненциально убывающие веса
 - Любая невозрастающая функция от порядкового номера
- Веса в зависимости от расстояния
 - Любая невозрастающая функция от расстояния
- Фиксированные веса объектов



Метод k ближайших взвешенных соседей среди набора эталонов

- Параметры метода: метрика, k , веса, **метод выбора эталонов**
- Алгоритм: по заданной метрике ищем k ближайших объектов среди эталонов выбранных из обучающей выборки и классифицируем объект взвешенным голосованием

Преимущества

- Простота реализации
- Хорошая интерпретируемость
- Параметр k можно оптимизировать по скользящему контролю

Недостатки

- Неустойчивость к выбросам
- Неоднозначность классификации при равных расстояниях до двух объектов
- Необходимость хранить всю обучающую выборку
- Алгоритм поиска вычислительно сложен
- Не учитывается значение расстояния

Задача

Получить примерно такое же качество работы алгоритма при меньшем количестве хранимых данных.

Возможно получить улучшение качества, так как в процессе выбора эталонов будут удалены выбросы.

Идеи

- Кластеризация объектов
- Жадный алгоритм



Выбор эталонов кластеризацией k средних (k-means)

Задача

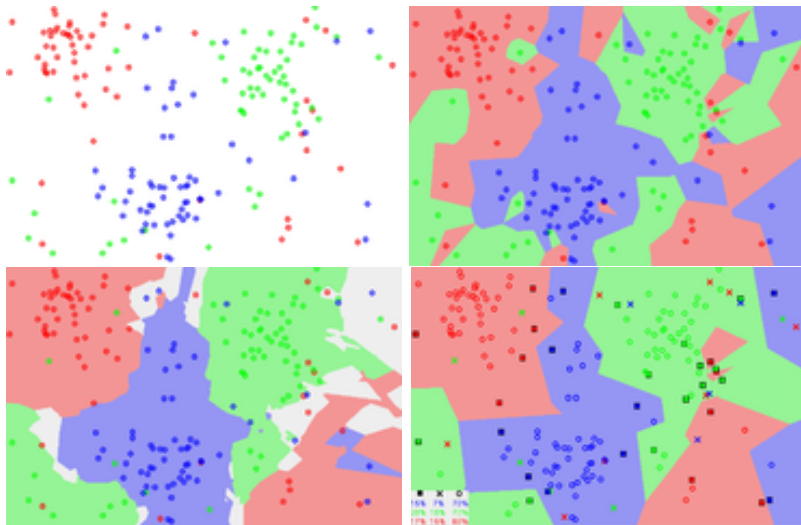
$$V = \sum_{i=1}^k \sum_{x \in S_i} (x - \mu_i) \rightarrow \min_{S_i},$$

где k — число кластеров, S_i — полученные кластеры, μ_i — центр масс S_i кластера.

Алгоритм

- 1 Случайно выбираются k элементов из выборки и объявляются центроидами
- 2 Для фиксированных центроидов каждый элемент выборки относится к одному из кластеров
- 3 Для фиксированных кластеров вычисляются центроиды
- 4 Пункты 2,3 повторяются до сходимости

Примеры 1-нп, 5-нп, 1-нп с выбором эталонов



Дополнительные модификации: RadiusNN

Идея

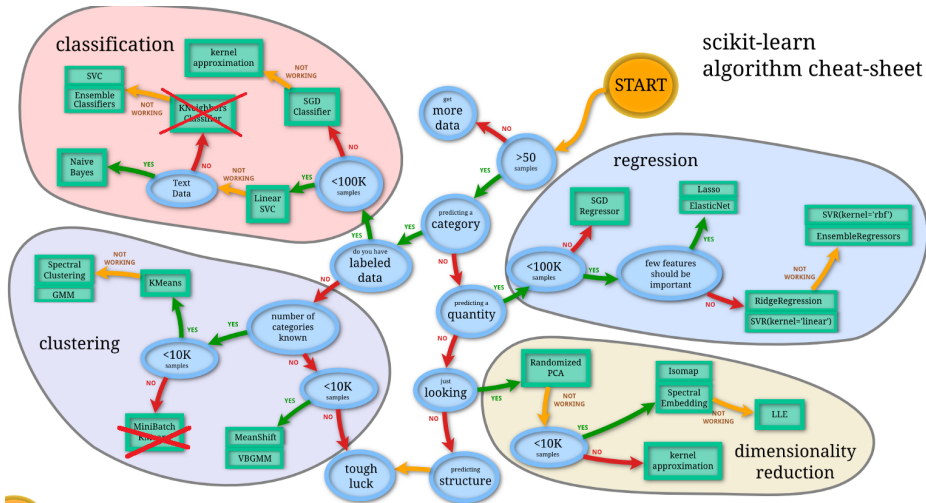
Часть имеет смысл искать соседей на расстоянии не больше чем некоторый радиус R

Параметр R

Вместо входного параметра количества соседей используется радиус



Дорожная карта Scikit-Learn³



³https://scikit-learn.org/stable/tutorial/machine_learning_map/

- Метод ближайших соседей – простой и хорошо интерпретируемый метод классификации
- Метод имеет большое число вариаций для настройки
 - Подбор метрики (metric learning)
 - Число ближайших соседей
 - Веса во взвешенном варианте метода
 - Алгоритм подбора эталонов





- Главный минус параметрических моделей, что для описания зависимости необходимо иметь параметрическую модель
- В случае невозможности подбора адекватной модели имеет смысл пользоваться непараметрическими регрессионными методами

Предположение

Близким объектам соответствуют близкие ответы



Простейшая модель

Приближаем искомую зависимость константой в некоторой окрестности

Формула Надарая-Ватсона

Если в окрестности точки несколько объектов из обучающей выборки, то разумно использовать взвешенное среднее в качестве предсказания алгоритма

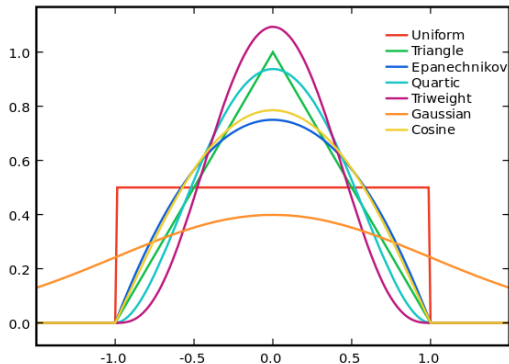
$$a(x) = \frac{\sum_i y_i \omega_i(x)}{\sum_i \omega_i(x)},$$

где $\omega_i(x) = K_h(x, x_i)$, а функция K_h называется ядром с шириной окна сглаживания h .



Примеры ядер

- $K_h(x, x_i) = K\left(\frac{\|x - x_i\|}{h}\right)$
- Типичные примеры ⁴



⁴[https://ru.wikipedia.org/wiki/Ядро_\(статистика\)](https://ru.wikipedia.org/wiki/Ядро_(статистика))

Напоминание: Вывод выражения среднеквадратичной ошибки

Определения

Пусть $y = y(x) = f(x) + \varepsilon$ — целевая зависимость, где $f(x)$ — детерминированная функция, $\varepsilon \sim N(0, \sigma^2)$ и $a(x)$ — алгоритм машинного обучения.

Полагаем, что ε и a — независимые ($Ea\varepsilon = EaE\varepsilon$). $Ey = Ef$, $Dy = D\varepsilon = \sigma^2$.

Разложение квадрата ошибки

$$\begin{aligned} E(y - a)^2 &= E(y^2 + a^2 - 2ya) = Ey^2 + Ea^2 - 2Eya = \\ &= Ey^2 + Ea^2 - 2E(f + \varepsilon)a = Ey^2 + Ea^2 - 2Efa - 2E\varepsilon a = \\ &= Ey^2 - (Ey)^2 + (Ey)^2 + Ea^2 - (Ea)^2 + (Ea)^2 - 2fEa = \\ &= Dy + Da + (Ey)^2 + (Ea)^2 - 2fEa = Dy + Da + (Ef)^2 - 2fEa + (Ea)^2 = \\ &= Dy + Da + (E(f - a))^2 = \sigma^2 + \text{variance}(a) + \text{bias}^2(f, a) \end{aligned}$$

Разброс и смещение для kNN

Разброс

$$\begin{aligned} \text{Variance}(a) &= D \left(\frac{1}{k} \sum_{i=1}^k y(x_{(i)}) \right) = \frac{1}{k^2} D \left(\sum_{i=1}^k y(x_{(i)}) \right) = \\ &= \frac{1}{k^2} D \left(\sum_{i=1}^k (f(x_{(i)}) + \varepsilon_i) \right) = \frac{1}{k^2} D \left(\sum_{i=1}^k f(x_{(i)}) \right) + \frac{1}{k^2} D \left(\sum_{i=1}^k \varepsilon_i \right) = \\ &= 0 + \frac{1}{k^2} k \sigma^2 = \frac{\sigma^2}{k} \end{aligned}$$

Смещение

$$\text{bias}^2(f, a) = (E(f(x_0) - a(x_0)))^2 = \left(f(x_0) - \frac{1}{k} \sum_{i=1}^k f(x_{(i)}) \right)^2$$

Bias-variance разложение для kNN

$$Error(x_0) = E(a(x_0) - f(x_0))^2 = \left(f(x_0) - \frac{1}{k} \sum_{i=1}^k f(x_{(i)}) \right)^2 + \frac{\sigma^2}{k} + \sigma^2$$

- С ростом k разброс уменьшается
- С ростом n сдвиг уменьшается



- Главное преимущество непараметрической регрессии — это отсутствие предположений о виде модели зависимости
- Метод имеет большое число вариаций для настройки
 - Подбор метрики (metric learning)
 - Число ближайших соседей
 - Веса во взвешенном варианте метода
 - Ширину окна сглаживания





Где могут быть полезны методы поиска ближайших соседей?



Точные

- Полный перебор
- К-мерное дерево (KD-tree)
- Метрическое дерево (ball-tree)

Приближенные

- Locality sensitive hashing (LSH)
- Navigable Small World (NSW)
- HNSW ⁵

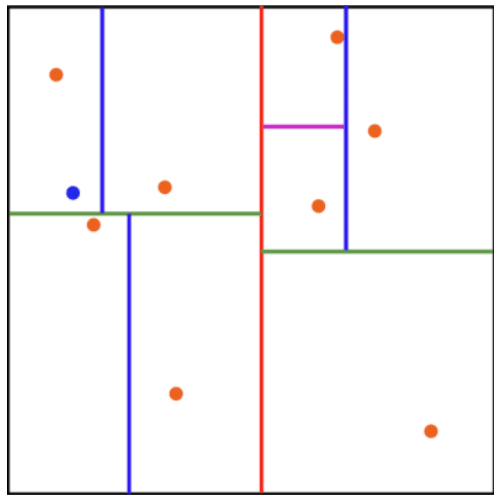
⁵<https://arxiv.org/abs/1603.09320>

Алгоритм построения

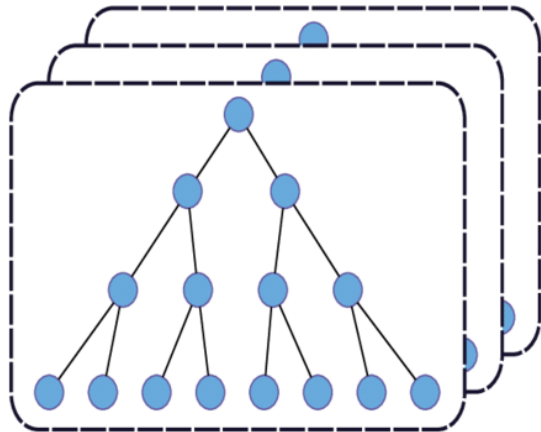
- 1 Если количество элементов меньше некоторого порогового значения, то отбивается лист и в него помещаются все элементы, в противном случае переходим к следующему пункту
- 2 Случайно выбирается признак, по которому будет разделение. По этому признаку ищется медиана
- 3 Все объекты с выбранным признаком левее медианы идут в левое поддерево, остальные в правое
- 4 Для левого и правого поддерева применяется та же процедура построения

⁶Bentley, J. L. (1975). Multidimensional binary search trees used for associative searching. Communications of the ACM. 18 (9): 509–517. doi:10.1145/361002.361007

К-мерное дерево



Kd-Tree in 2D



Multiple Randomized Kd-Trees

Поиск ближайшего соседа в K -мерном дереве

Алгоритм поиска I

Для нашего запроса идём по дереву и в соответствующем листе ищем нужное количество ближайших соседей

Идея

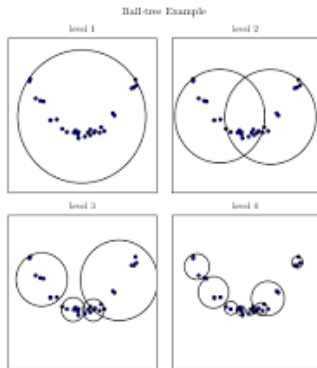
Если расстояние до дальнего ближайшего соседа меньше, чем расстояние до разделяющей гиперплоскости, то это означает, что во втором поддереве ближайших соседей нет.

Алгоритм поиска II

- 1 Выполняем шаги алгоритма I, считая в каждой вершине расстояние до разделяющей гиперплоскости
- 2 Делаем обратный ход алгоритма, если расстояние до разделяющей гиперплоскости меньше, чем расстояние до дальнего ближайшего соседа

Идея

Использовать вместо полугиперплоскостей шары



⁷Omohundro, Stephen M. (1989), Five Balltree Construction Algorithms

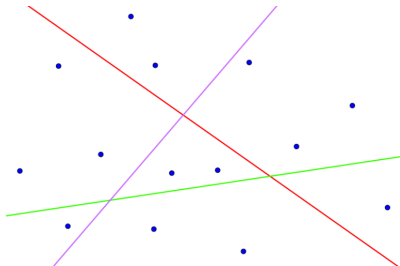
Locality Sensitive Hashing (LSH) ⁸

Идея

Разделить пространство, используя хэш-функции

Пример

В качестве семейства функций можно рассмотреть гиперплоскости



⁸<https://codeforces.com/blog/entry/54080?locale=ru>

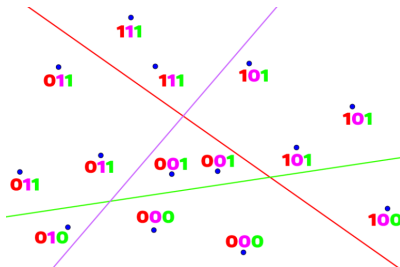
Locality Sensitive Hashing (LSH) ⁸

Идея

Разделить пространство, используя хэш-функции

Пример

В качестве семейства функций можно рассмотреть гиперплоскости



⁸<https://codeforces.com/blog/entry/54080?locale=ru>

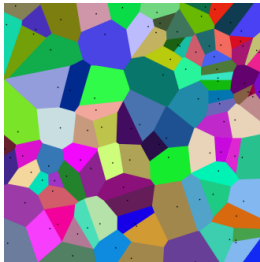
Жадные алгоритмы на графе

Жадный алгоритм

Итеративно ищем ближайшего соседа в графе

Теорема

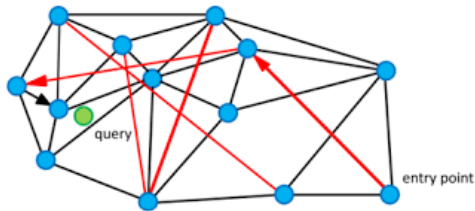
Для графа Делоне жадный алгоритм решает задачу поиска ближайшего соседа точно



Navigable Small World (NSW) ⁹

Идея

Поиск по графу типа Small World



Гарантии

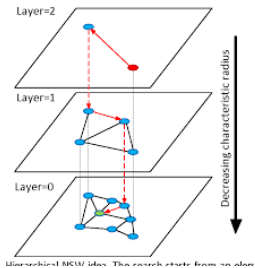
Теоретических гарантий нет, поэтому алгоритм поиска запускают несколько раз в зависимости от требуемой точности

⁹Y. Malkov, A. Ponomarenko, A. Logvinov, and V. Krylov, Approximate nearest neighbor algorithm based on navigable small world graphs, Information Systems, vol. 45, pp. 61-68, 2014.

Hierarchical Navigable Small World (HNSW) ¹⁰

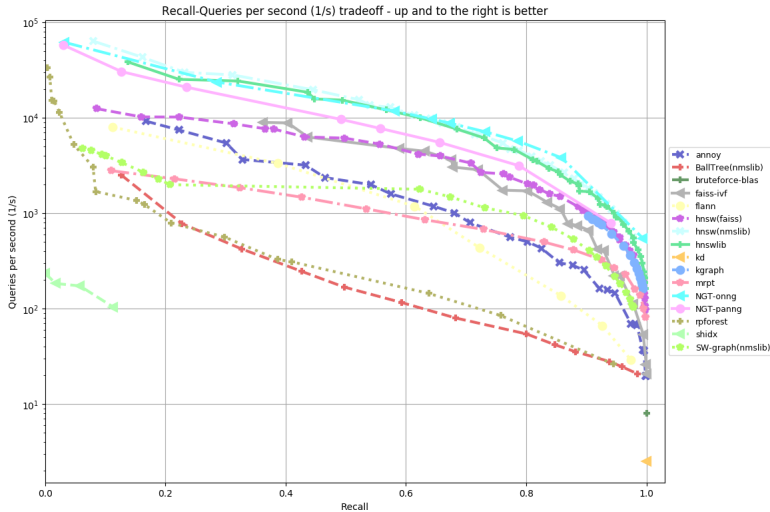
Идея

В графе типа Small World можно выделить подграфы меньшего размера и сделать поиск итеративным



¹⁰<https://arxiv.org/abs/1603.09320>

Сравнение методов поиска ближайших соседей ¹¹



¹¹<https://github.com/erikbern/ann-benchmarks>

- Метод поиска ближайших соседей — важная задача теории алгоритмов
- Нужно помнить, что есть методы в среднем быстрее, чем полный перебор
- Для современных индустриальных систем характерно использование не точных, но очень быстрых алгоритмов поиска





Спасибо за внимание!