

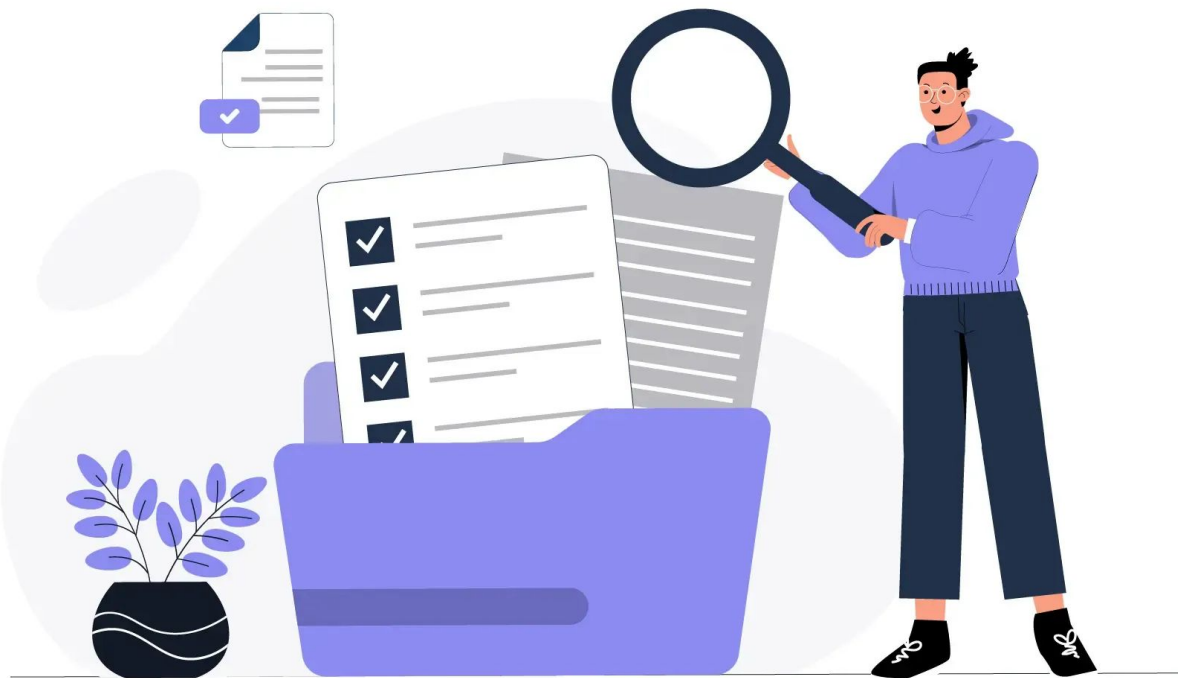
Векторный поиск

Елена Кантонистова

Далее в программе

1. Какие задачи будем решать?
2. Точный поиск по словам и по эмбедингам
3. HNSW
4. FAISS

Какие задачи будем решать?



Виды поиска

1. Точный поиск по словам
2. Точный поиск по эмбедингам
3. Приближенный поиск по эмбедингам

1. Точный поиск по словам: BM25 - способ посчитать релевантность запроса q документу D

Для каждого токена запроса q_i мы можем посчитать степень его релевантности документу D - это $BM25(q_i)$:

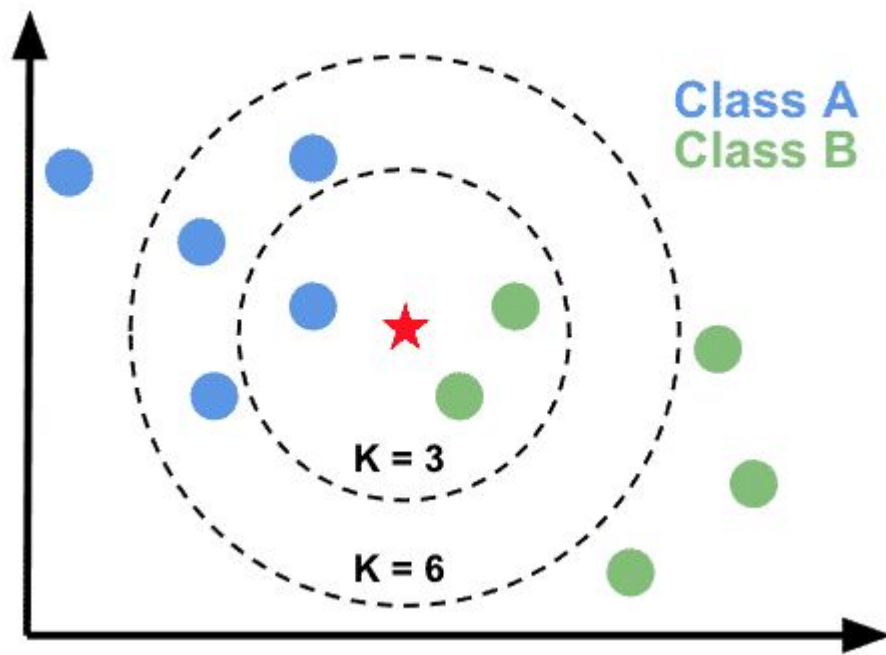
$$\sum_i^n IDF(q_i) \frac{f(q_i, D) * (k1 + 1)}{f(q_i, D) + k1 * (1 - b + b * \frac{fieldLen}{avgFieldLen})}$$

where:

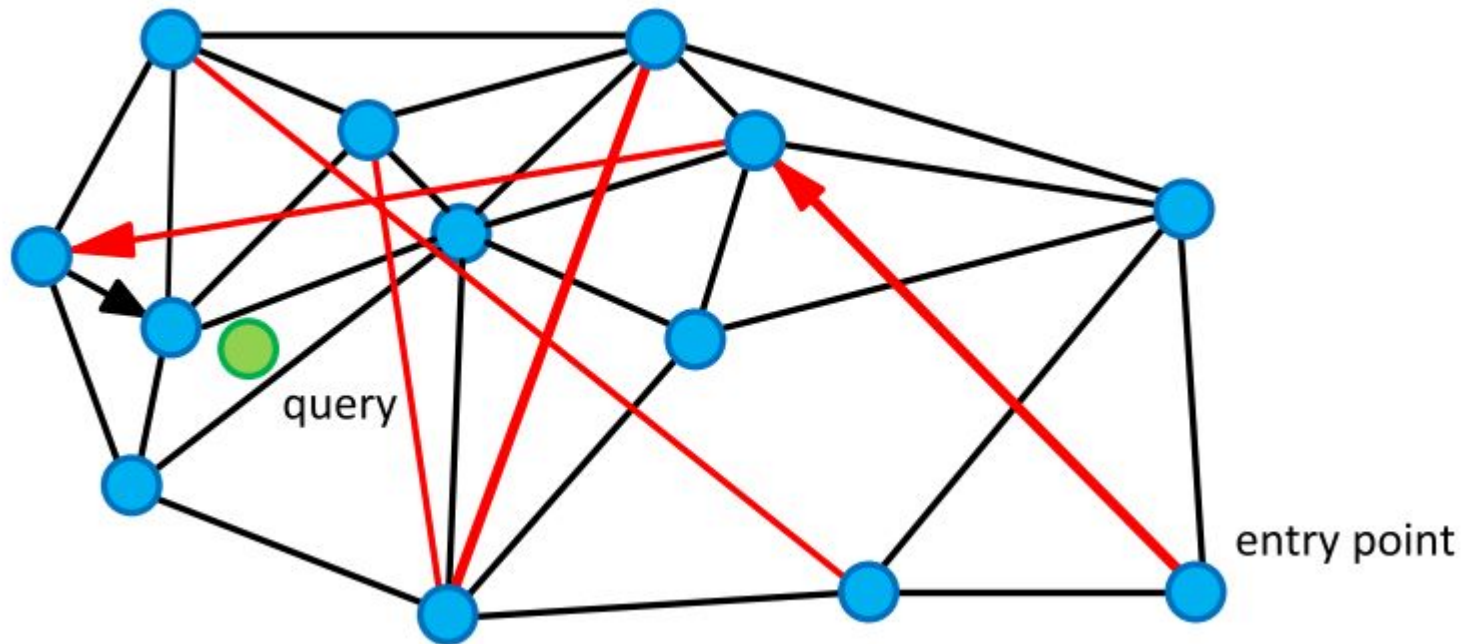
- ✓ q_i is a query term,
- ✓ $f(q_i, D)$ is q_i 's term frequency in the document D ,
- ✓ D is the length of the document,
- ✓ $avgdl$ is the average document length in the text collection,
- ✓ $k1$ and b are free parameters, usually chosen empirically (common values are $k1=2.0$ and $b=0.75$),
- ✓ $IDF(q_i)$ is the IDF for q_i .

Далее суммируем эти релевантности по всем токенам q_i запроса q .

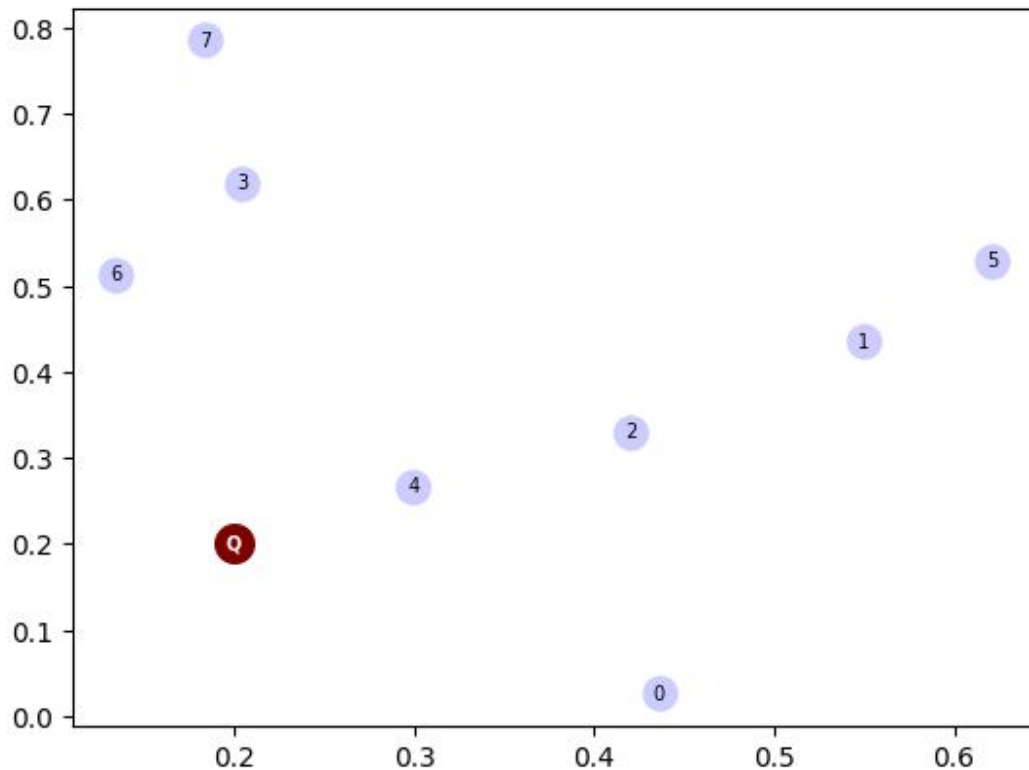
2. Точный поиск соседей (на векторах слов): KNN



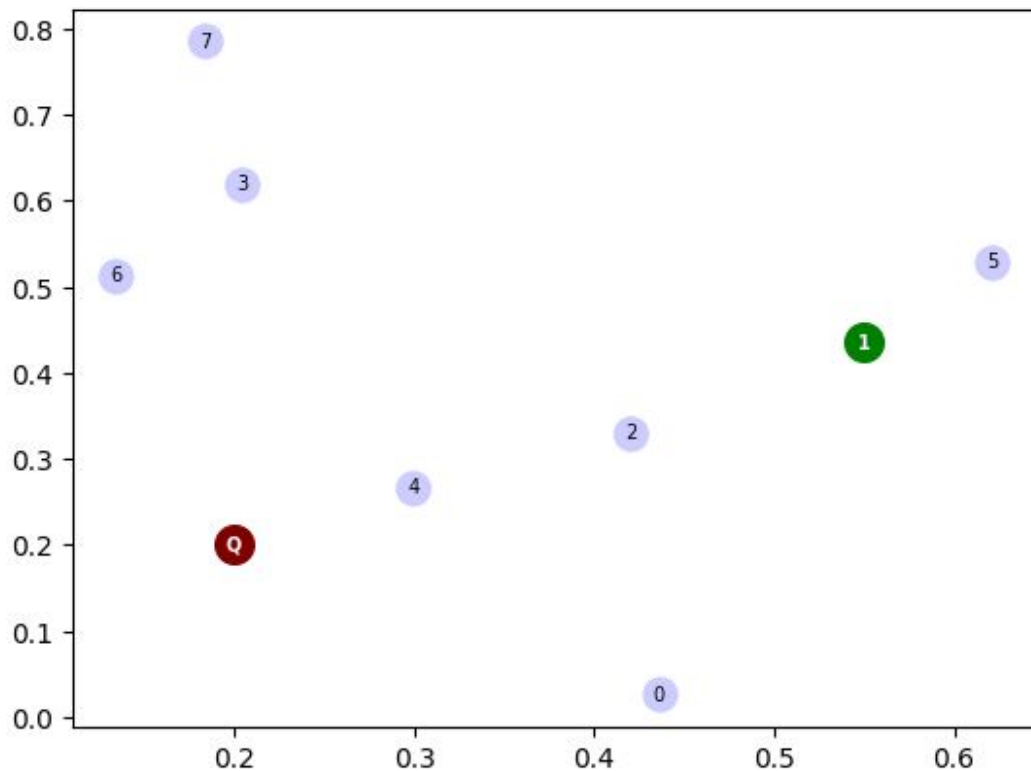
3. Approximate NN: NSW -> HNSW



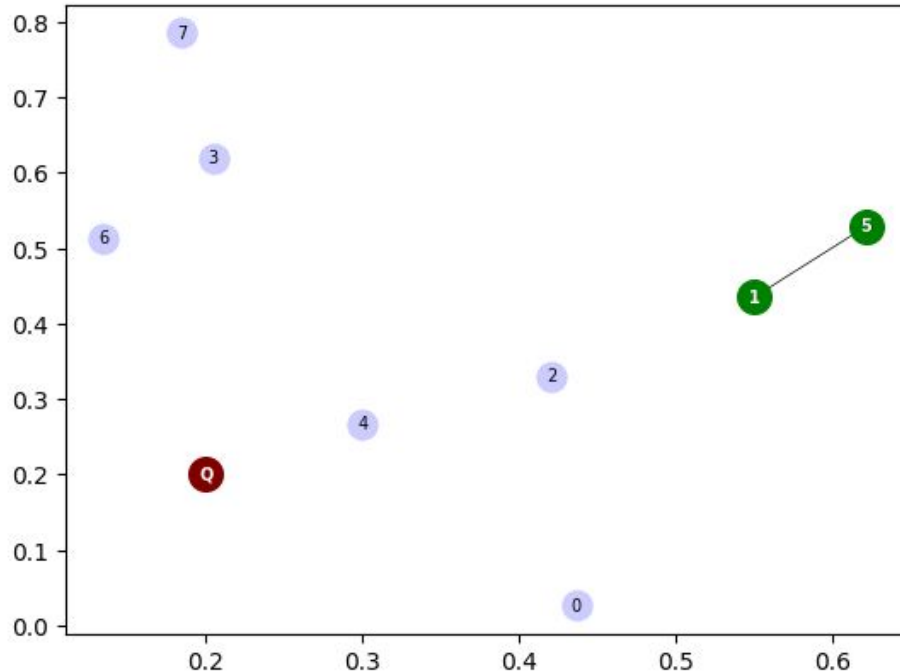
Navigable Small World (NSW): строим граф



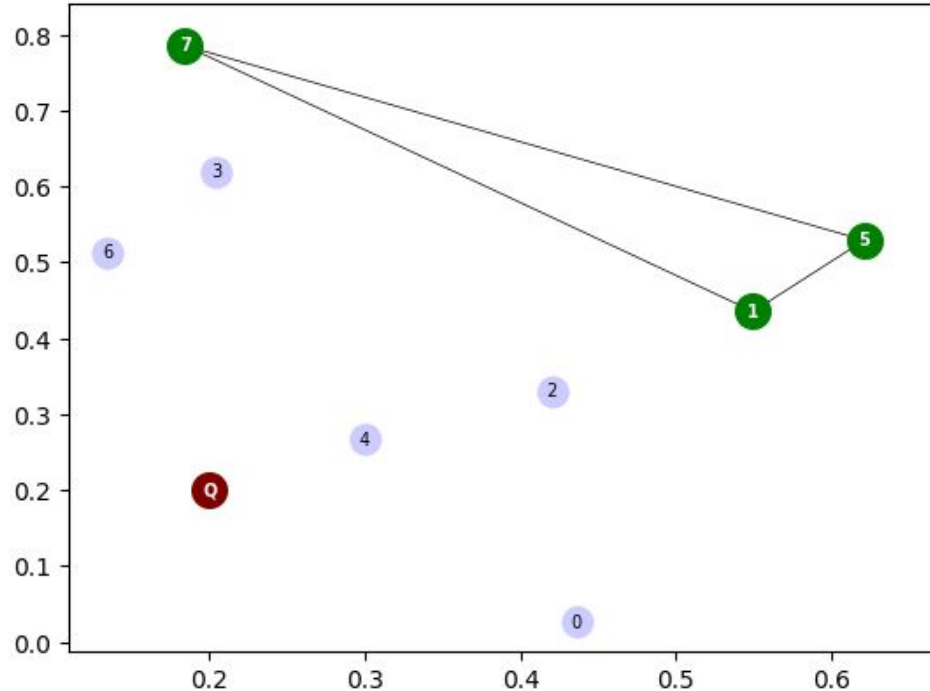
Берем случайную вершину (эмбеддинг) и добавляем в граф



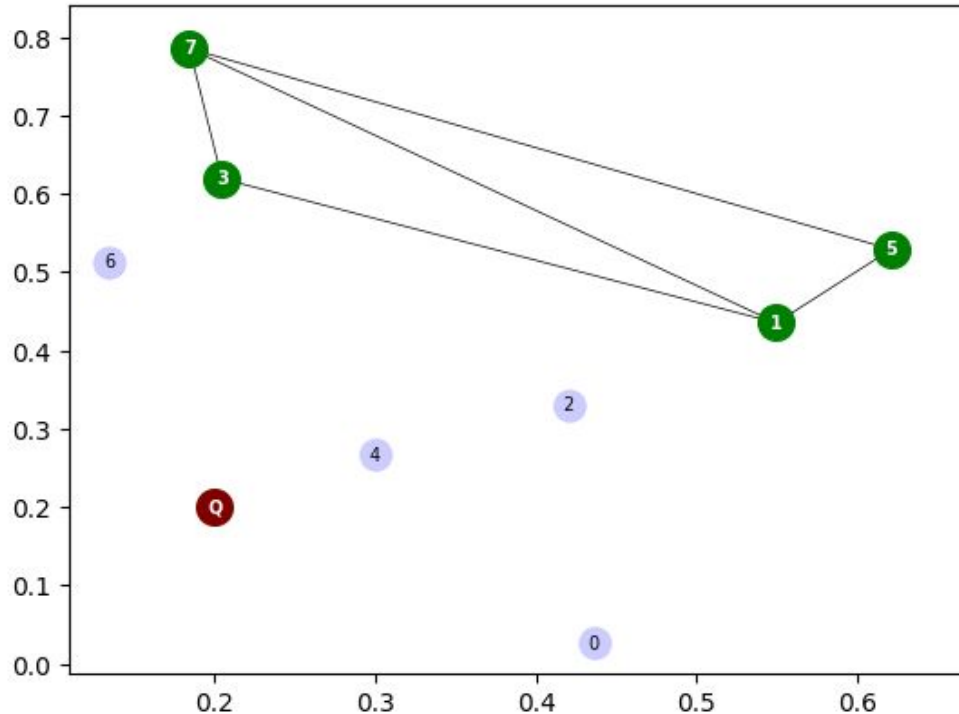
Каждую следующую случайную вершину соединяем с двумя ближайшими, которые уже есть в графе



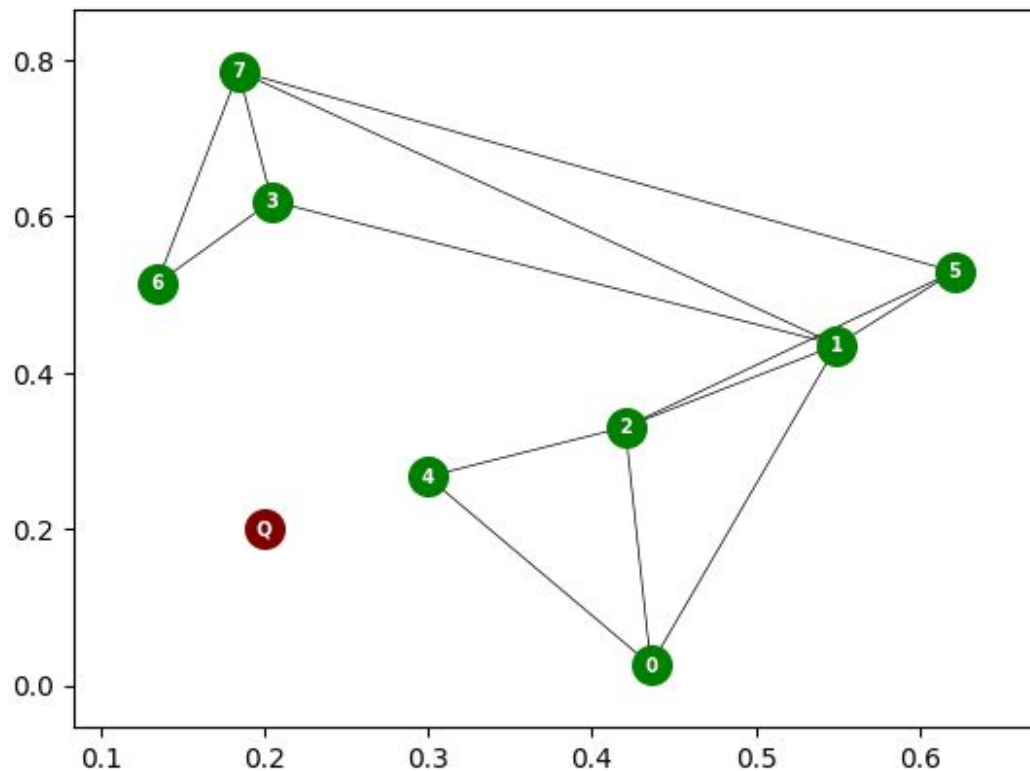
Каждую следующую случайную вершину соединяем с двумя ближайшими, которые уже есть в графе



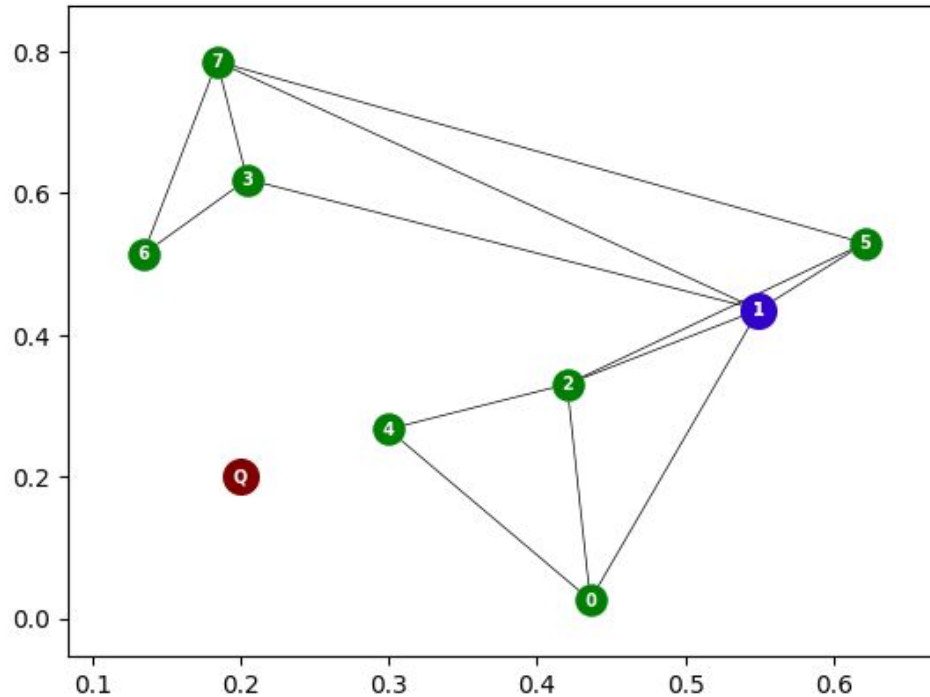
Каждую следующую случайную вершину соединяем с двумя ближайшими, которые уже есть в графе



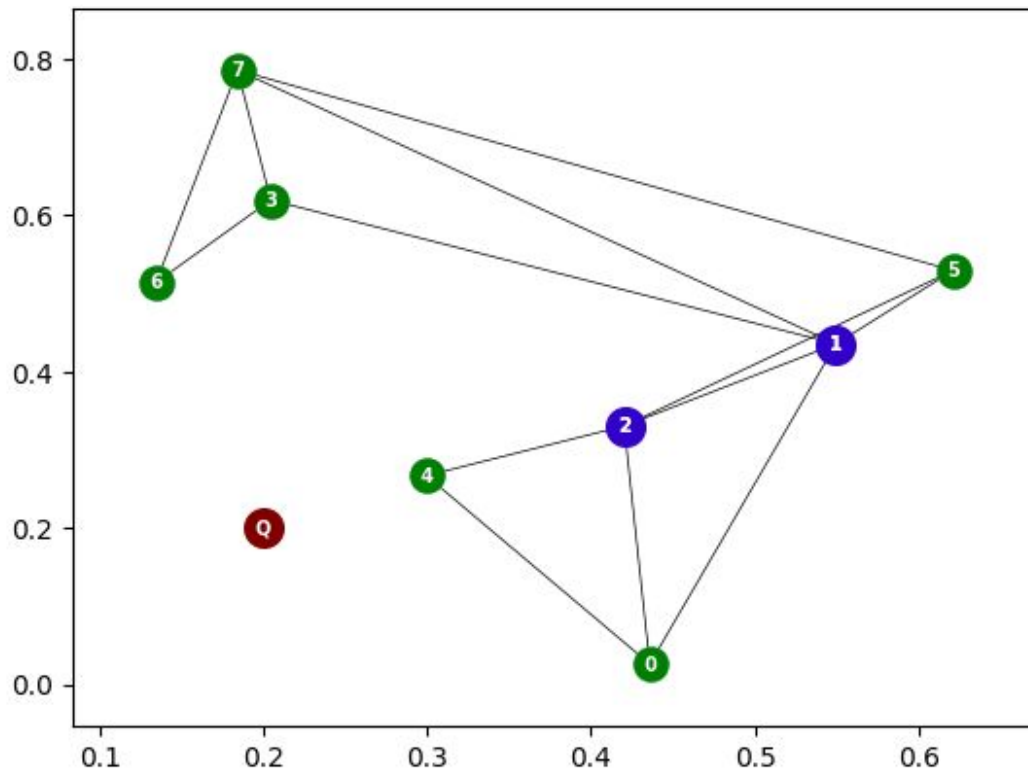
Navigable Small World (NSW): получаем граф



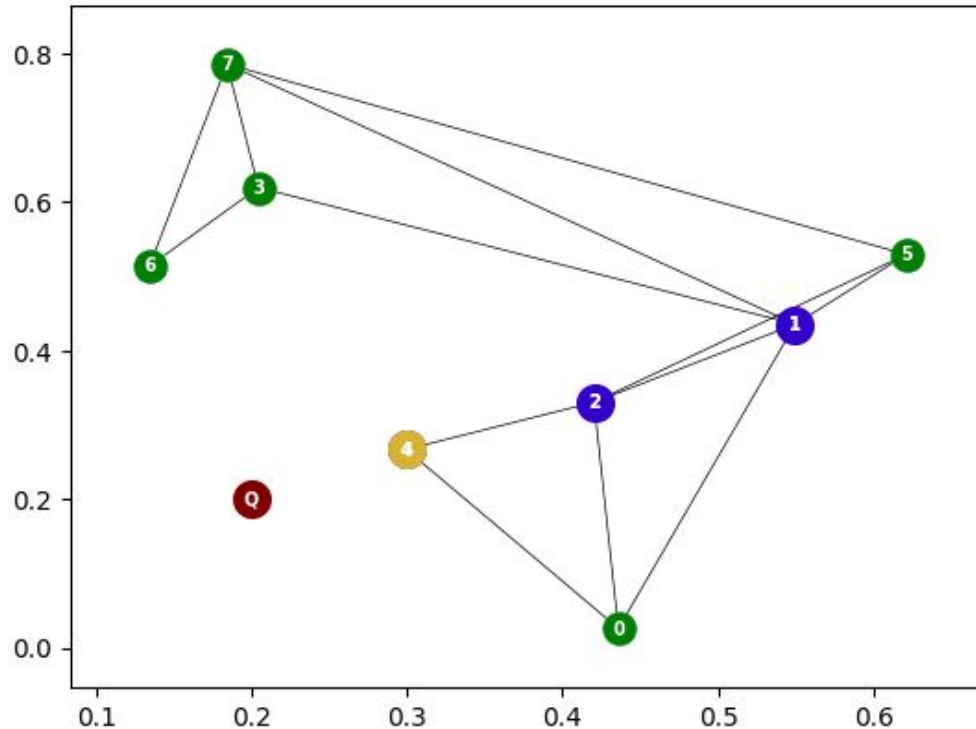
NSW: поиск по графу - берем случайную вершину и считаем расстояние от q до нее и ее соседей



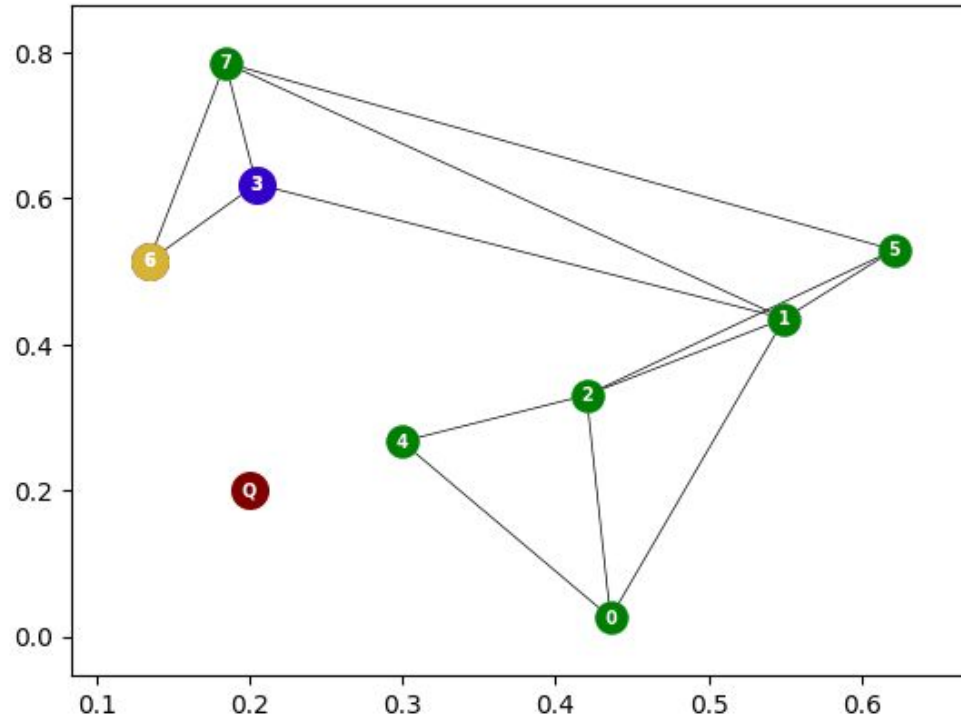
NSW: поиск по графу - идем в ближайшую к q из соседей



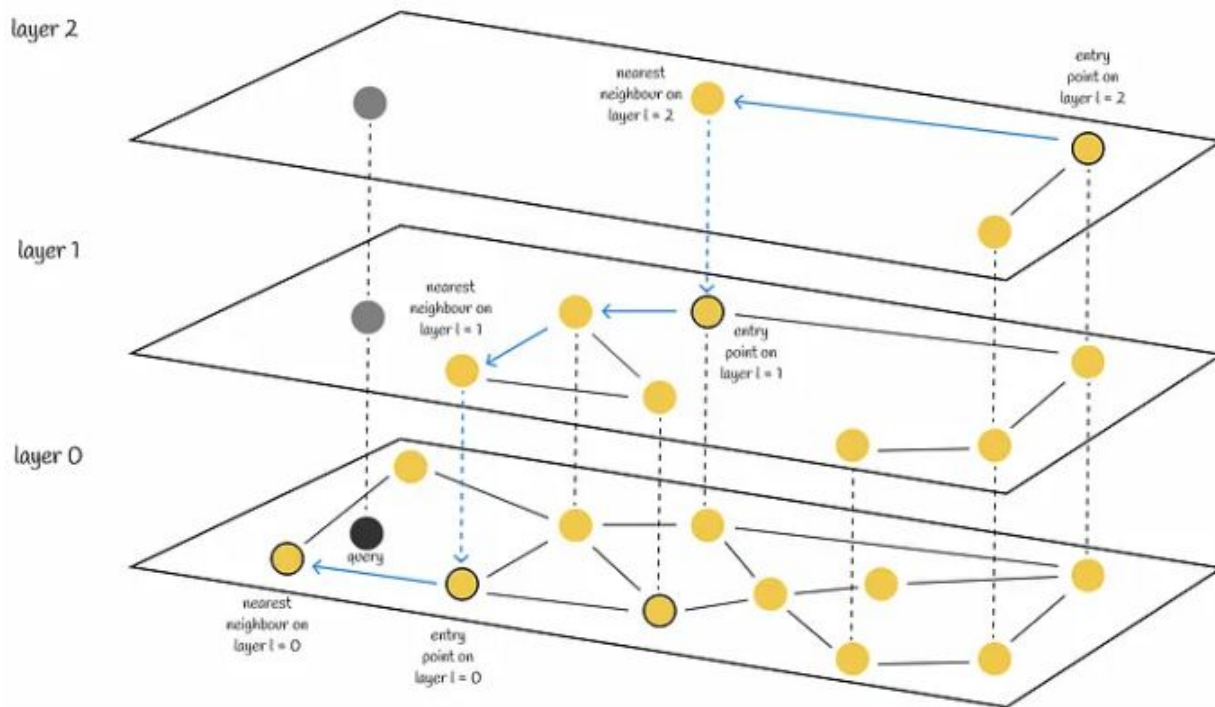
NSW: поиск по графу - идем в ближайшую к q из соседей и из нее повторяем процесс



NSW: итог - получаем приближенного ближайшего соседа. Почему? Зависит от стартовой точки.



HNSW - Hierarchical NSW - когда запросов много!



HNSW - Hierarchical NSW - когда запросов много!

Строим уровневые графы:

- Каждый вектор при добавлении получает случайный максимальный уровень L - то есть точка добавляется на уровни $1, 2, \dots, L$
- На каждом уровне точка соединяется с M ближайшими соседями (по выбранной метрике)

Добавление точки (объекта) в граф:

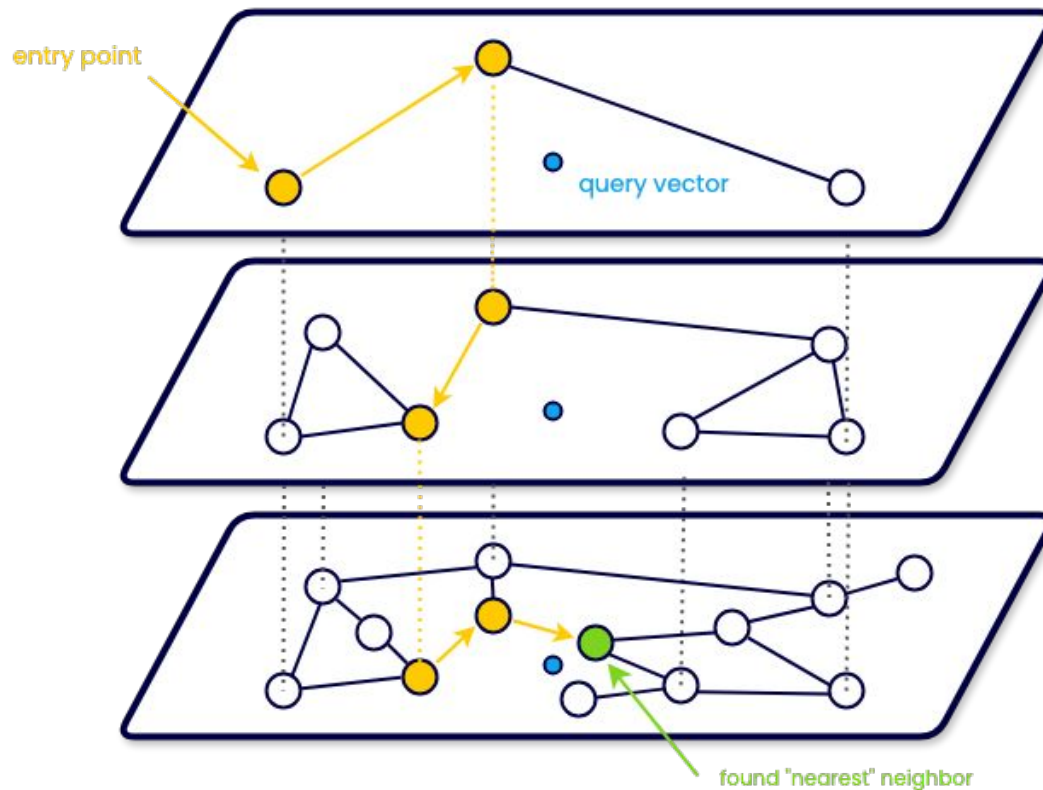
- Задаем точке случайный уровень L
- Добавляем точку на верхний уровень и соединяем с M ближайшими соседями
- Далее добавляем точку на следующий уровень и ищем там M ближайших соседей и так далее

HNSW: поиск по графу

Когда приходит эмбединг запроса:

- Начало поиска — стартуем с произвольной (или заранее выбранной) точки на самом верхнем уровне
- Грубая навигация — двигаемся по рёбрам, пока не найдём ближайшую точку на этом уровне
- Спуск вниз — переходим на уровень ниже, используя найденную точку как старт
- Точный поиск — на нижнем (нулевом) уровне ищем ближайших соседей до нужного k .

HNSW: поиск по графу



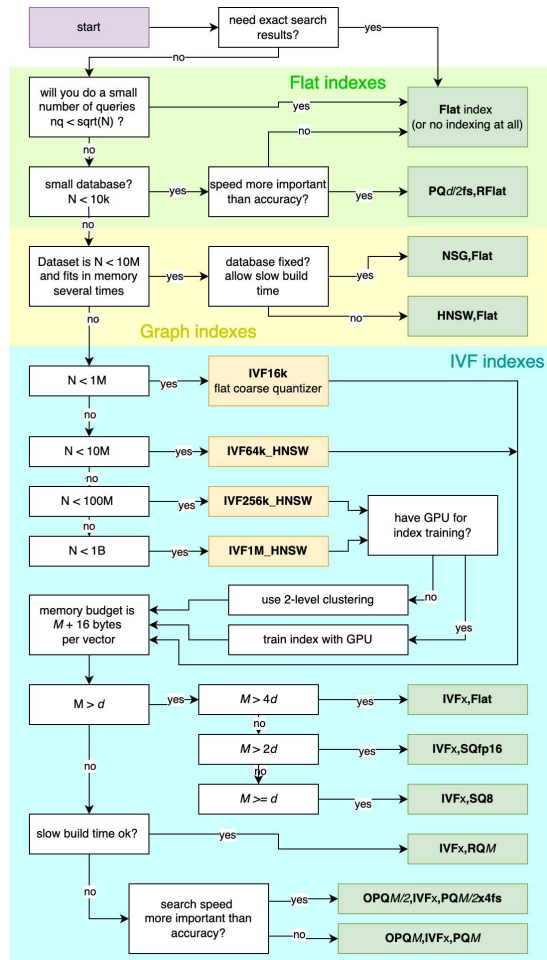
Библиотеки: FAISS

FAISS (Facebook AI Similarity Search) — это библиотека от Facebook для эффективного поиска ближайших соседей (ANN — Approximate Nearest Neighbors) в больших коллекциях векторов.

Индекс в FAISS

- ✓ Индекс - это структура, в которой физически лежат эмбединги
- ✓ Вид структуры зависит от выбранного алгоритма
 - ✓ Flat = массив
 - ✓ HNSW = многоуровневый граф
 - ✓ IVF = кластеры
 - ✓ PQ = сжатые векторы в кодовых блоках

FAISS



Практика

https://colab.research.google.com/drive/1o-zBXgLiPQ_b4rO2UgBu2F82EkR-8dWp?usp=sharing