

Объяснение методологических обоснований LSH

Введение

Locality-Sensitive Hashing (LSH) представляет собой вероятностный метод приближённого поиска ближайших соседей (Approximate Nearest Neighbor, ANN) в высокоразмерных пространствах. В отличие от классических хеш-функций, минимизирующих коллизии, LSH целенаправленно максимизирует вероятность коллизии для семантически близких объектов. Ключевая идея заключается в проекции данных в пространство меньшей размерности с сохранением свойств локальности: близкие объекты с высокой вероятностью попадают в один и тот же "bucket"[1].

Успех LSH в практических приложениях (поиск дубликатов, рекомендательные системы) обусловлен не только теоретическими гарантиями, но и спецификой построения хеш-структур и свойствами реальных данных. Данная работа направлена на объяснение методологических основ LSH, включая теоретическое обоснование, практические улучшения и ограничения метода.

Теоретические основы LSH

Формальное определение и базовые принципы

Формально, семейство хеш-функций \mathcal{H} называется (r, cr, p_1, p_2) -чувствительным для метрики d , если для любых двух точек p и q выполняются условия:

1. Если $d(p, q) \leq r$, то $\Pr[h(p) = h(q)] \geq p_1$,
2. Если $d(p, q) \geq cr$, то $\Pr[h(p) = h(q)] \leq p_2$,

где $c > 1$, а $p_1 > p_2$ [1].

Для усиления вероятностных свойств используются конструкции AND и OR:

- **AND-construction:** Комбинация k независимых хеш-функций усиливает строгость попадания в bucket ($p' = p_1^k$)
- **OR-construction:** Использование L независимых хеш-таблиц увеличивает вероятность попадания хотя бы в один bucket ($p' = 1 - (1 - p_1)^L$)

На практике теоретическая вероятность коллизии $P_{h_k^\ell}(r)$ часто интерпретируется как ожидаемая полнота (recall rate). Хотя $P_{h_k^\ell}(r)$ является математическим ожиданием, эмпирические исследования показывают, что для большинства реальных datasets фактическая полнота оказывается близка к теоретическому ожиданию при условии корректного выбора параметров k и ℓ и достаточного размера данных [1]. Это объясняется тем, что дисперсия полноты оказывается крайне мала благодаря агрегации множества независимых хеш-функций и действию закона больших чисел.

Вероятностные гарантии и регулирование ошибки

LSH предоставляет строгие вероятностные гарантии для качества поиска. Для найденного кандидата x можно оценить вероятность того, что он является истинным s -приближенным соседом:

$$P[\text{success}] \geq 1 - (1 - p_1^k)^L$$

Вероятность ложноположительного срабатывания (далекая точка попала в бакет) не превышает:

$$P[\text{false positive}] \leq p_2^k \cdot L$$

Параметры LSH позволяют явно регулировать точность и вероятность ошибки:

- **Ширина бакета (w):** Увеличивая w , снижаем вероятность коллизии для всех точек, но увеличиваем точность
- **Число sampled dimensions (m):** Увеличение m снижает дисперсию оценки расстояния, улучшая точность
- **Число хеш-функций (k):** Увеличение k снижает p_1 и p_2 , делая фильтрацию жестче
- **Число хеш-таблиц (L):** Увеличение L повышает вероятность нахождения истинного соседа, но увеличивает память и время

Поддержка различных метрик расстояния

LSH поддерживает различные метрики расстояния при условии выбора соответствующих LSH-семейств:

Метрики и соответствующие LSH-семейства

- **Евклидово расстояние:** Используется r -стабильное распределение (распределение Гаусса). Вероятность коллизии вычисляется как:

$$p(r) = \int_0^w \frac{1}{r} \exp\left(-\frac{t^2}{2r^2}\right) \left(1 - \frac{t}{w}\right) dt$$

- **Расстояние Хэмминга:** Используется LSH-семейство на основе битовых операций. Вероятность коллизии определяется долей совпадающих бит:

$$p(r) = 1 - \frac{r}{d}$$

где d — размерность пространства.

- **Косинусное расстояние:** Применяются гиперплоскостные LSH-функции ($\text{sign}(\langle r, x \rangle)$). Вероятность коллизии зависит от угла между векторами:

$$p(r) = 1 - \frac{\arccos(\text{cosine_sim})}{\pi}$$

Сравнительная характеристика метрик

Метрика	LSH-семейство	Вероятность коллизии	Оптимальные сценарии
Евклидово	р-стабильное распределение	$p(r) = \int_0^w \frac{1}{r} \exp\left(-\frac{t^2}{2r^2}\right) \left(1 - \frac{t}{w}\right) dt$	Общие данные, изображения, векторы
Расстояние Хэмминга	Битовые хэши	$p(r) = 1 - \frac{r}{d}$	Тексты, категориальные данные, хэши
Косинусное	Гиперплоскостное	$p(r) = 1 - \frac{\arccos(\text{cosine_sim})}{\pi}$	Тексты, рекомендации, NLP

Таблица 1: Характеристики LSH для разных метрик расстояния

Методологические улучшения

Ускорение вычислений: FastLSH

Классические реализации LSH для евклидова пространства (E2LSH) требуют вычисления скалярных произведений полномерных векторов, что приводит к сложности $O(n)$ на одну хеш-функцию. Метод FastLSH [2] сочетает случайное проецирование и случайную выборку признаков (random sampling). Вместо полного n -мерного вектора для каждого хеша случайно выбирается подмножество из m координат ($m \ll n$), к которому применяется стандартная LSH-функция. Это сокращает вычислительную сложность до $O(m)$.

Теоретически показано, что при $m \rightarrow \infty$ распределение расстояний в sampled-пространстве сходится к исходному. При конечных m расхождения моментов распределения незначительны, что сохраняет LSH-свойства. Эксперименты подтверждают, что даже при $m = 30$ (при исходной размерности в сотни/тысячи) качество поиска практически не ухудшается, а скорость возрастает в десятки раз.

Оптимизация поиска: roLSH

Другое направление оптимизации — сокращение дорогостоящих обращений к данным (I/O operations). Метод roLSH [3] предлагает две стратегии:

- **roLSH-samp**: На этапе индексации выполняется сэмплирование тестовых запросов для определения статистики оптимальных радиусов поиска. Это позволяет начинать поиск не с минимального радиуса, а с предварительно оцененного значения.
- **roLSH-NN**: Используется neural network для предсказания оптимального радиуса поиска для конкретного запроса на основе его хешей. Это позволяет для каждого запроса стартовать с радиуса, близкого к оптимальному, минимизируя последующие расширения.

Оба метода значительно сокращают количество random disk seeks — основного bottleneck в external memory implementations.

Ограничения и дальнейшие направления

Несмотря на эффективность, LSH имеет ряд ограничений:

- **Зависимость от параметров**: Производительность сильно зависит от выбора параметров (k , L , w и др.), которые часто подбираются эмпирически.
- **Чувствительность к данным**: Эффективность может снижаться на данных со специфическими распределениями или аномалиями.
- **Вычислительная сложность**: Хотя LSH сокращает общее число сравнений, стоимость построения индекса может быть высокой для очень больших datasets.

Перспективные направления исследований включают:

- Разработку data-dependent методов, адаптирующихся к распределению данных
- Автоматический подбор параметров (k , L , w) под целевую точность/скорость
- Гибридные подходы, сочетающие LSH с другими методами индексации
- Реализацию поддержки различных метрик расстояния через наследование от базового класса LSH

Заключение

Locality-Sensitive Hashing является мощным инструментом для приближённого поиска ближайших соседей в высокоразмерных пространствах. Его методологическая основа сочетает

вероятностные гарантии с практической эффективностью. Теоретически обоснованная вероятность коллизии находит подтверждение в реальных приложениях благодаря эффекту больших данных и тщательному подбору параметров.

Современные улучшения, такие как FastLSH и roLSH, address ключевые limitations классических схем, делая LSH ещё более scalable и практичным. Поддержка различных метрик расстояния и возможность явного регулирования точности делают LSH универсальным инструментом для решения задач поиска похожих объектов в больших массивах данных.

Список литературы

- [1] Lu, K., Wang, H., Xiao, Y., Song, H. Why locality sensitive hashing works: A practical perspective // Information Processing Letters. – 2018. – Т. 136. – С. 49-58.
- [2] Tan, Z., Wang, H., Xu, B., Luo, M., Du, M. Fast Locality Sensitive Hashing with Theoretical Guarantee // arXiv preprint arXiv:2309.15479v1. – 2020.
- [3] Jafari, O., Nagarkar, P., Montano, J. Improving Locality Sensitive Hashing by Efficiently Finding Projected Nearest Neighbors // arXiv preprint arXiv:2006.11284v1. – 2020.