

# GRAFO DE K VECINOS MÁS CERCANOS

El artículo explora la construcción del grafo del k-vecino más cercano (KNNG) en datos de alta dimensión, destacando su uso en búsqueda de similitud, reducción de dimensiones y agrupamiento. Se centra en dos etapas principales del proceso: la construcción inicial del KNNG denominado INIT y su refinamiento mediante la propagación de vecindad denominado NBPG. El estudio presenta una comparación experimental de métodos existentes y propone alternativas para la inicialización del KNNG utilizando estructuras de indexación recientes como SW y HNSW. Además, se investiga la efectividad de la propagación de vecindad y su relación con las propiedades de los datos, concluyendo que ciertas características de los datos afectan significativamente el rendimiento de los métodos de propagación de vecindad. Los experimentos muestran que la elección del método inicial de KNNG y la estrategia de propagación de vecindad pueden tener un impacto considerable en la precisión y eficiencia del KNNG resultante. El estudio sugiere que la combinación de técnicas puede ofrecer mejoras significativas, especialmente en contextos donde la precisión y la eficiencia son críticas.

**PALABRAS CLAVE:** grafo de k vecinos más cercanos (KNNG), propagación de vecindad, datos de alta dimensión, concentración y análisis experimental

**Metodología:** El estudio se centra en dos etapas principales del proceso de construcción del KNNG: la inicialización (INIT) y la propagación del vecindario (NBPG). Se analiza la literatura existente, destacando la falta de comparaciones experimentales exhaustivas entre las soluciones propuestas y la ausencia de pruebas con estructuras de indexación modernas como SW y HNSW en la fase INIT.

**Experimentación y Análisis:**

- **Inicialización (INIT):** Se experimenta con métodos existentes y nuevas alternativas (SW y HNSW) para la creación inicial del KNNG, evaluando su impacto en la precisión del KNNG final.
- **Propagación de vecindad (NBPG):** Se investiga cómo la propagación de vecindad mejora la precisión del KNNG inicial, analizando la correlación entre las propiedades de los datos y la efectividad del NBPG.

**Hallazgos Clave:**

1. **Efectividad de Métodos INIT:** Los métodos modernos como HNSW muestran un rendimiento superior en la creación inicial del KNNG en comparación con métodos más tradicionales, ofreciendo un balance entre eficiencia y precisión.
2. **Impacto de la Propagación de vecindad:** La propagación de vecindad es crucial para refinar el KNNG, con una fuerte dependencia del tipo de datos. Los datos con alta "centralidad" o "hubness" tienden a beneficiarse menos de la propagación de vecindad.
3. **Comparaciones Experimentales:** Se proporciona una evaluación comparativa integral de los métodos, destacando que no existe una solución única que domine en todos los escenarios. La elección del método adecuado depende específicamente de las características del conjunto de datos y los requisitos de la aplicación.

**Recomendaciones:**

- Implementar HNSW para INIT cuando se busca un balance óptimo entre eficiencia y precisión.
  - Considerar la característica de "hubness" de los datos antes de elegir la estrategia de NBPG.
  - Realizar evaluaciones comparativas en contextos específicos de aplicación para determinar la mejor combinación de métodos INIT y NBPG.
- Conclusión:** El estudio ofrece insights prácticos y metodológicos para mejorar la construcción de KNNG en datos de alta dimensión, con recomendaciones específicas para maximizar la eficacia de las técnicas utilizadas. La adaptabilidad de las estrategias en función de las propiedades de los datos es crucial para optimizar el rendimiento y la aplicabilidad del KNNG en diversas aplicaciones.

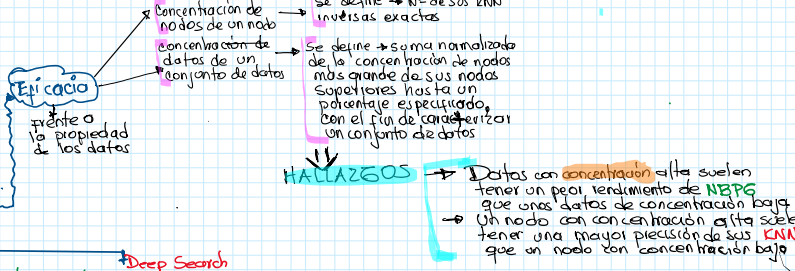
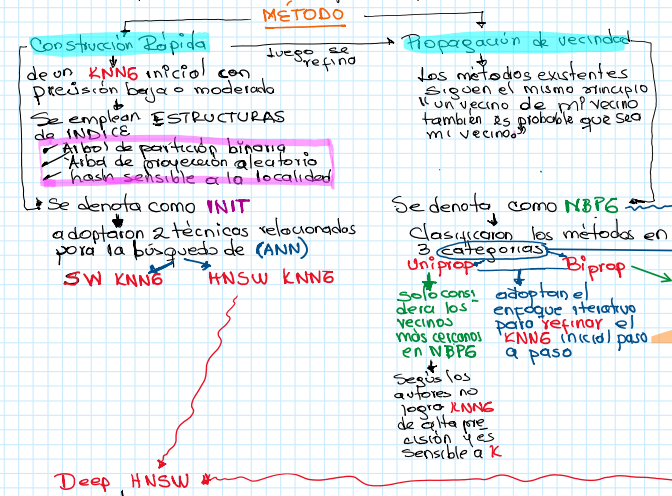
El grafo de (KNNG) se utiliza en  $\rightarrow$  visión artificial  
 $\rightarrow$  minería de datos  
 $\rightarrow$  Aprendizaje automático

KNNG  $\rightarrow$  Dado un conjunto de datos de dimensión  $d: D \subset \mathbb{R}^d$  y un entero positivo  $K$  trata cada punto de  $U \in D$  como un nodo del grafo y crea aristas desde  $U$  hasta sus  $K$  vecinos más cercanos (KNN) en  $D \setminus \{U\}$

En este artículo utilizamos punto, vector y nodo de manera intercambiable.

(KNNG) se utiliza para varias operaciones:

- \* Búsqueda aproximada de  $K$  vecinos más cercanos (ANN)
- \* Los métodos de búsqueda de (ANN) de última generación como  $\rightarrow$  DPS, NSG
- Para ello primero construyen (KNNG) preciso en  $D$  y luego ajustan el vecindario para cada nodo en la parte superior de (KNNG)



**Deep HNSW**

es una nueva combinación que presenta el mejor equilibrio entre EFICIENCIA, PRECISIÓN y MEMORIA en la mayoría de los casos

**PROBLEMA DE INVESTIGACION**

Sean  $D \subset \mathbb{R}^d$  un conjunto de datos que consta de  $n$   $d$ -dimensionales vectores reales. El KNNG en  $D$  se define como

**Def 1.** KNNG Sea  $D \subset \mathbb{R}^d$  y un entero positivo  $K$ , sea  $G = (V, E)$  el KNNG de  $D$  donde  $V$  es el conjunto de nodos y  $E$  es el conjunto de aristas. Cada nodo  $v \in V$  representa de manera única un vector en  $D$ . Existe una arista dirigida  $(u, v) \in E$  si y solo si  $v$  es uno de los KNN de  $u$  en  $D$  ¿?



Fig. 2. An illustrative KNNG with  $n = 6$  and  $k = 2$ .

En la figura se muestra el KNNG de un conjunto de



Fig. 2. An illustrative KNNG with  $n = 6$  and  $k = 2$ .

En la figura se muestra el **KNNG** de un conjunto de datos con 6 puntos de datos. Para cada nodo, tiene 2 aristas salientes que apuntan a sus 2NN.

El artículo se centra en soluciones en memoria sobre vectores densos de alta dimensión con la distancia euclidiana como medida.

→ Marco de Construcción

```

Algoritmo 1: Marco de construcción de KNNG
Entrada: Conjunto de datos D y un entero positivo k.
Salida: KNNG G
1 INIT: generar un gráfico inicial G rápidamente;
2 NBPQ: refinar G por propagación de vecindad; 3 devolver G;
  
```

También analizamos sus requerimientos de memoria y complejidad temporal. Para los requerimientos de memoria, solo contamos las estructuras de datos auxiliares de cada método; ignorando las estructuras comunes, es decir, los datos D y el KNNG G.

→ KNNG vs Grupo de Proximidad

```

Algoritmo 2: Búsqueda en el gráfico (H, q, k, efSearch, ep)
Entrada: H, q, k, efSearch y ep
Salida: KNN de q, deje
que el grupo sea el conjunto candidato e inserte ep en el grupo;
2 L = max(k, efSearch) e i = 0; 3 mientras i <
efSearch haga 4 u = pool[i] y
marque u como expandido;
5 /* Procedimiento: expand(q, u, H) para
cada vecino v de u en H hacer
6 /* Procedimiento: update(pool, v)
pool.add(v, dist(q, v)); ordenar
7 pool en orden ascendente de dist(q, v); si pool.size()
> L, entonces pool.resize(L);
8
9 i = índice del primer nodo no expandido u en el grupo;
10
11
12 devuelve los primeros k nodos del grupo;
  
```

→ Búsqueda de (ANN)

Sea H un grupo de proximidad y q una consulta. La búsqueda comienza en uno o más nodos ep específicos o seleccionados aleatoriamente, que primero se introduce en un conjunto de candidatos (una lista de nodos expandidos).

Sea  $L = \max(k, efSearch)$  el tamaño de pool. En cada iteración de la línea 4-11, encontramos el primer nodo no expandido  $u$  en pool y luego expandimos  $u$  en las líneas 6-10 lo que se denota como  $expand(q, u, H)$ . Esta expansión trata a cada vecino  $v$  de  $u$  en H como un candidato para definir el pool en la línea 8-10, lo que se denota como  $update(pool, v)$ . Una vez que se han examinado los primeros nodos  $efSearch$  en pool, el proceso de búsqueda finaliza y devuelve los primeros  $k$  nodos en el pool. Obviamente  $efSearch$  es clave para el rendimiento de la búsqueda. Un  $efSearch$  grande aumenta tanto el costo como la precisión.

③ → Construcción de Grafos Iniciales

Métodos para INIT

Entoque basado en Particiones

particionar los puntos de datos en grupos suficientemente pequeños pero más similares.

Juego cada punto encuentra sus **KNN** en el grupo o grupos más próximos y se genera un **KNNG** inicial.

Incluye → La división múltiple **LSH KNNG**

Métodos q adaptan diferentes técnicas de partición

Entoque basado en mundos pequeños

Los puntos de datos se dividen en grupos lo suficientemente pequeños pero similares entre sí.

Consiste en crear un gráfico de mundo pequeño navegable entre los puntos de datos, cuya estructura es útil para identificar los **KNN** de cualquier punto de datos.

técnicas → **SW KNNG** **TRANS KNNG**

→ División Aleatoria múltiple

→ wang et al → propusieron este enfoque para construir grafos de vecindad aprox iniciales de base y luego unirlos para obtener un **KNNG** inicial.

Particione los puntos de datos en D en 2 grupos disjuntos.

Para asignar puntos cercanos al mismo grupo, la partición se realiza en la línea de la división principal.

El conjunto de datos D se particiona en 2 grupos disjuntos  $D_1$  y  $D_2$  que se dividen nuevamente de manera recursiva hasta que el grupo a dividir contiene como máximo  $1/d_{div}$  puntos. Luego se construye un subgrafo de vecindad para cada grupo  $O(d * T_{div}^2)$ .

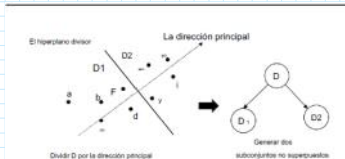


Fig. 3. Ilustre la división de D por la dirección principal (representada por la flecha continua). Como se muestra con la línea discontinua, el hiperplano de división que es perpendicular a la dirección principal divide a D en dos grupos disjuntos,  $D_1$  y  $D_2$ , que se dividen de forma recursiva. La figura debe recrearse.

Una única división produce un grafo de vecindad aproximado base que contiene una serie de subgrupos aislados y no puede conectar un punto con sus vecinos más cercanos que se encuentran en subgrupos diferentes. Por lo tanto, la división múltiple utiliza divisiones aleatorias  $L_{div}$  que permite calcular la dirección principal a partir del conjunto grupo.

Para cada punto  $u \in D$  tiene un conjunto de **KNN** de cada división y por lo tanto los conjuntos  $L_{div}$  totales de **KNN**.

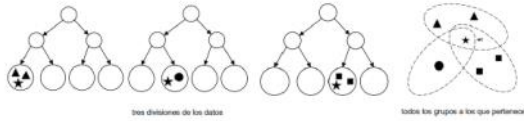


Fig. 4. Esta figura ilustra tres divisiones aleatorias de un conjunto de datos D. Para el punto de datos y denotado como una estrella, sus KNN se calculan en el grupo de cada una. división y luego se combinaron para devolver los mejores KNN.

### Cost Analysis

Dado que esas divisiones  $L_{div}$  se pueden generar de forma independiente, el requisito de memoria es  $O(n)$  donde  $n$  es el tamaño de la división correcta. Para cada división, la función recursiva `partition cost`  $O(n * d * \log n)$  y el procedimiento de `brute-force` en todos los grupos finalmente genera los datos se requiere  $O(n * d * L_{div})$ . Con  $L_{div}$  divisiones, el costo es múltiple para  $cost$  of multiple division. Podemos ver que  $L_{div}$  y  $T_{div}$  son clave.

### 3.2 LSH KINGS → Zhang et al propusieron el método

Utiliza la técnica de hash para dividir  $D$  en grupos de igual tamaño y luego construye un grupo de vecindad en cada grupo usando el método `brute-force`. Se crean divisiones `hash` para construir la base `hash` de grupos de vecindad aproximada que luego se combinan para generar un `KINGS` de mayor precisión.

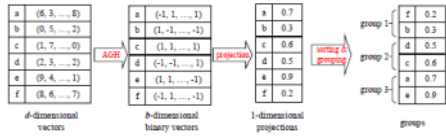


Fig. 5. The process of creating a division by LSH in LSH KINGS.