

k-Nearest Neighbors (k-NN) Implementation

Ioannis Michalainas, Iasonas Lamprinidis

October/November 2024

Abstract

This project involves implementing the k-NN algorithm to find the nearest neighbors of a set of query points Q relative to a corpus set C in a high-dimensional space. Given a set of M -point corpus data and N -point query data, both in D -dimensional space, the algorithm identifies the k-nearest neighbors of each query point. Using optimized matrix operations and the CBLAS library, our implementation efficiently handles high-dimensional distance calculations.

1 Problem Statement

The objective of this project is to implement a subroutine that computes the **k-nearest neighbors (k-NN)** of each query point in Q based on their distances to the data points in C .

To calculate the distances, we use the following formula:

$$D = \sqrt{C^2 - 2CQ^T + (Q^2)^T} \quad (1)$$

where:

- C is the set of data points (corpus).
- Q is the set of query points (query).
- D is the distance matrix containing distances between each pair of points from C and Q .

Each row of the $q \times c$ matrix D contains distances from a query point to all corpus points. We then use the quickselect algorithm to retrieve the k smallest distances and their corresponding indices in $O(n)$ time.

2 Example

In this section, we illustrate the process of generating random data points, calculating the distance matrix, and finding the k-NN using C code snippets. For clarity, let's assume we have sets C and Q in d -dimensional space.

2.1 Input

We begin by either generating random data points for both the dataset C and query set Q or by reading a .mat file. The following function creates a dataset with a specified number of points and dimensions:

```
1 void random_input(Mat* matrix, size_t points, size_t dimensions) {
2
3     srand(time(NULL) + (uintptr_t)matrix);
4
5     matrix->data = (double*)malloc(points*dimensions*sizeof(double));
6     memory_check(matrix->data);
7
8     matrix->rows = points;
9     matrix->cols = dimensions;
10
11     for (size_t i = 0; i < points * dimensions; i++) {
12         // Scale to [0, 200], then shift to [-100, 100]
13         matrix->data[i] = ((double)rand()/RAND_MAX)*200.0 - 100.0;
14     }
15 }
```

Listing 1: Generating random data points

For example, with points = 5 and dimensions = 2, a possible generated dataset could be:

$$C = \begin{bmatrix} 98.87 & 77.36 \\ 85.86 & 21.03 \\ -61.65 & -54.45 \\ -57.33 & 76.06 \\ 30.87 & 66.55 \end{bmatrix}$$

And the query dataset Q with 4 points could be:

$$Q = \begin{bmatrix} 92.90 & 21.38 \\ -76.71 & -29.80 \\ -83.48 & -40.61 \\ -46.21 & 64.69 \end{bmatrix}$$

2.2 Distance Calculation

For smaller datasets, computing the distance matrix is straightforward. Using the formula

$$D = \sqrt{C^2 - 2CQ^T + (Q^2)^T} \quad (2)$$

we compute matrix D (queries \times corpus) that represents the Euclidean distance for each query q to each corpus point c .

```
1 void calculate_distances(const Mat* C, const Mat* Q, long double* D) {
2
3     int c = (int)C->rows;
4     int d = (int)C->cols;
5
6     double* C2 = (double*)malloc(c*sizeof(double));
7     double Q2;
8     memory_check(C2);
9
10    #pragma omp parallel for
11    for(int i=0; i<c; i++) {
12        double sum = 0.0;
13        for(int j=0; j<d; j++) {
14            sum += C->data[i*d + j]*C->data[i*d + j];
15        }
16        C2[i] = sum;
17    }
18
19    double sum = 0.0;
20    for(int j=0; j<d; j++) {
21        sum += Q->data[j]*Q->data[j];
22    }
23    Q2 = sum;
24
25    double* CQ = (double*)malloc(c*sizeof(double));
26    memory_check(CQ);
27
28    cblas_dgemm(CblasRowMajor, CblasNoTrans, CblasNoTrans, c, 1, d, 1.0, C->data, d, Q->data,
29                1, 0.0, CQ, 1);
30
31    #pragma omp parallel for
32    for(int i=0; i<c; i++) {
33        D[i] = C2[i] - 2*CQ[i] + Q2;
34        if(D[i] < 0.0) {
35            D[i] = 0.0;
36        }
37    }
38
39    free(C2);
40    free(CQ);
41 }
```

Listing 2: Calculating distances using CBLAS

The computed distance matrix D between each point in C and Q is:

$$D = \begin{bmatrix} 56.30 & 7.05 & 76.74 \\ 28.89 & 107.62 & 144.42 \\ 25.85 & 119.57 & 156.71 \\ 15.91 & 77.10 & 120.14 \end{bmatrix}$$

Each entry $D[i][j]$ in this matrix represents the distance between the i -th point in Q and the j -th point in C .

2.3 Finding k-Nearest Neighbors

To find the k-nearest neighbors, we use the **quickselect** algorithm, which efficiently identifies the smallest k elements in $O(n)$ time and returns their values, as well as their indices. The ‘findKNN’ function calculates the distance from every query to every corpus point, finds the k-nearest neighbors, and stores them in matrix N .

In this example, we find the 3 nearest neighbors ($k = 3$) for each point in Q based on the distance matrix D :

$$N = \begin{bmatrix} 56.30 - id0 & 7.05 - id1 & 76.74 - id4 \\ 28.89 - id2 & 107.62 - id3 & 144.42 - id4 \\ 25.85 - id2 & 119.57 - id3 & 156.71 - id4 \\ 15.91 - id3 & 77.10 - id4 & 120.14 - id2 \end{bmatrix}$$

In this matrix: - The first row indicates that for the first query point, the 3 nearest neighbors are 56.30, 7.05, and 76.74, with indexes 0, 1, and 4 respectively. - The second row shows the 3 nearest neighbors for the second query point, and so forth.

3 Fine Lines

3.1 Parallelism

If matrices C and Q are very large, there is a chance that matrix D (queries \times corpus) does not fit in memory. To combat this, we calculate matrix D in segments (slices), use each slice to find the k -NN of the corresponding query, and then discard the slice to make room in memory. This way, we prevent running out of memory. To boost execution speed, we parallelize the workload of computing the D slice and running quickselect on it.

```
1 void findKNN(Mat* C, Mat* Q, Neighbor* N, int k) {
2
3     int const c = C->rows;
4     int const q = Q->rows;
5     int const d = C->cols;
6
7     #pragma omp parallel for
8     for(int i=0; i<q; i++) {
9
10        long double* D = (long double*)malloc(c*sizeof(long double));
11        memory_check(D);
12
13        calculate_distances(C, &(Mat){.rows = 1, .cols = d, .data = Q->data + i*d}, D);
14
15        int* indices = (int*)malloc(c*sizeof(int));
16        memory_check(indices);
17
18        for(int j=0; j<c; j++) {
19            indices[j] = j;
20        }
21
22        quickSelect(D, indices, 0, c-1, k, N + i*k);
23        free(D);
24        free(indices);
25    }
26 }
```

Listing 3: Routine for computing the k smallest distances

3.2 Problem Minimization

When the datasets become arbitrarily large, we encounter some problems: A problem that arises as we scale up is computational complexity. To remedy the extreme computational cost and the curse of high dimensions, we use random projection, based on the Johnson–Lindenstrauss lemma, which proves we can reduce dimensions by projecting our matrices to lower dimensions.

```
1 void random_projection(Mat* M, int t, Mat* RP) {
2
3     srand(time(NULL) + (uintptr_t)M);
4
5     int n = (int)M->rows;
6     int d = (int)M->cols;
7
8     RP->rows = n;
9     RP->cols = t;
10    RP->data = (double*)malloc(n*t*sizeof(double));
11
12    double* R = (double*)malloc(d*t*sizeof(double));
13    for (int i=0; i<d*t; i++) {
14        // Rademacher distribution (-1 or +1)
15        R[i] = (rand()%2 == 0 ? -1 : 1) / sqrt((double)t);
16    }
17
18    cblas_dgemm(CblasRowMajor, CblasNoTrans, CblasNoTrans, n, t, d, 1.0, M->data, d, R, t,
19               0.0, RP->data, t);
20
21    free(R);
22 }
```

Listing 4: Routine for applying Random Projection

It is of great importance to choose the target dimensionality t carefully when using random projection. Based on the lemma, $t = \frac{4\log(c)}{\epsilon^2}$, where ϵ is the allowed error and c is the number of corpus points, is a good approximation to retain accuracy while significantly reducing computational cost.

```
1 if((c*d>1000000 && d>100) || d>500) {
2
3     Mat C_RP, Q_RP;
4     double const e = 0.1;
5     int const t = 4*log((double)c) / (e*e);
6     printf("Target dimension (t) for random projection: %d\n", t);
7     random_projection(&C, t, &C_RP);
8     random_projection(&Q, t, &Q_RP);
9 }
```

Listing 5: Approximation Check

Applying this minimization technique is valuable only when the dataset is large. A custom metric to calculate when to use approximation would be $c \times d > 10^6$ when $d > 100$ or $d > 500$.

3.3 Details

We avoid computing the square root of the distances we calculate in ‘calculate_ddistances’. Instead, we compare squared distances 0 and $b > 0$: $a > b \iff a^2 > b^2$) and then return their square root. This way, instead of $c \times q$ sqrt calculations, we make $k \times q$ (substantially less).

```
1 void quickSelect(long double* arr, int* indices, int left, int right, int k, Neighbor*
   result) {
2
3     if(left <= right) {
4         int pivotIndex = partition(arr, indices, left, right);
5
6         if(pivotIndex == k-1) {
7
8             for(int i=0; i<k; i++) {
9                 result[i].distance = sqrt(arr[i]);
10                result[i].index = indices[i];
11            }
12            return;
13        } else if(k-1 < pivotIndex) {
14            quickSelect(arr, indices, left, pivotIndex-1, k, result);
15        } else {
16            quickSelect(arr, indices, pivotIndex+1, right, k, result);
17        }
18    }
19 }
```

Listing 6: Routine for applying Random Projection

4 Summary

The following steps summarize the process:

- a. Calculate squared terms for data points C and query points Q .
- b. Compute the dot product CQ^T to facilitate the distance calculation.
- c. Combine results to obtain the distance matrix D .
- d. Identify the k-nearest neighbors for each query point using quickselect.

5 Conclusion

This project demonstrates an efficient k-NN algorithm implementation that uses advanced matrix operations and sorting techniques. This approach enables accurate neighbor searches in high-dimensional spaces, achieving a balance between clarity and computational efficiency.

6 Tools

In this project, we used the following tools:

1. The C programming language, its compiler, and standard libraries.
2. The linear algebra library OpenBLAS.
3. The file reading libraries HDF5 and Matio.
4. The parallel programming libraries OpenMP, OpenCilk, and Pthreads.
5. The version control software Git.
6. The GNU/Linux OS.
7. The Neovim text editor.

7 Sources

- https://en.wikipedia.org/wiki/Random_projection
- https://en.wikipedia.org/wiki/Johnson%E2%80%93Lindenstrauss_lemma
- <https://cs.stanford.edu/people/mmahoney/cs369m/Lectures/lecture1.pdf>
- <https://github.com/erikbern/ann-benchmarks?tab=readme-ov-file>