

k-Nearest Neighbors (k-NN) Implementation

Ioannis Michalainas

October/November 2024

Abstract

This C project implements the k-NN algorithm to identify the nearest neighbors of a set of query points Q relative to a corpus set C in a high-dimensional space. Given a corpus of c points and q query points, both in d -dimensional space, the algorithm efficiently determines the k-nearest neighbors for each query point. The implementation uses optimized matrix operations, parallel processing, and approximation techniques to handle large datasets and high-dimensional distance calculations efficiently.

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 . *This implementation also works when $Q = C$.*

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 (or .hdf5). 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 . In larger datasets, we typically compute matrix D in chunks of size 300, as this size was found to be optimal for balancing parallelism on our machine while avoiding unnecessary overhead. Further details on the calculation of matrix D for larger datasets will be provided in the following chapter. For smaller input datasets, matrix D is computed in its entirety.

```
1 void calculate_distances(const Mat* C, const Mat* Q, int start_idx, int end_idx, long
   double* D) {
2
3   int c = C->rows;
4   int d = C->cols;
5   int batch_size = end_idx - start_idx;
6
7   double* C2 = (double*)malloc(c*sizeof(double));
8   memory_check(C2);
9   #pragma omp parallel for
10  for(int i=0; i<c; i++) {
11    double sum = 0.0;
12    for(int j=0; j<d; j++) {
13      sum += C->data[i*d + j]*C->data[i*d + j];
14    }
15    C2[i] = sum;
16  }
17
18  double* Q2 = (double*)malloc(batch_size*sizeof(double));
19  memory_check(Q2);
20  #pragma omp parallel for
21  for(int i=0; i<batch_size; i++) {
22    double sum = 0.0;
23    for(int j=0; j<d; j++) {
24      sum += Q->data[(start_idx+i)*d + j] * Q->data[(start_idx+i)*d + j];
25    }
26    Q2[i] = sum;
27  }
28
29  double* CQ = (double*)malloc(c*batch_size*sizeof(double));
30  memory_check(CQ);
31
32  cblas_dgemm(CblasRowMajor, CblasNoTrans, CblasTrans, c, batch_size, d,
33             1.0, C->data, d, Q->data + start_idx * d, d, 0.0, CQ, batch_size);
34
35  #pragma omp parallel for collapse(2)
36  for(int i=0; i<c; i++) {
37    for(int j=0; j<batch_size; j++) {
38      D[j*c + i] = C2[i] - 2.0*CQ[i*batch_size + j] + Q2[j];
39
40      if(D[j*c + i] < 0.0) {
41        D[j*c + i] = 0.0;
42      } else if(D[j*c + i] < 1) {
43        D[j*c + i] = sqrt(D[j*c + i]);
44      }
45    }
46  }
47
48  free(C2);
49  free(Q2);
50  free(CQ);
51 }
```

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*.

```
1 void quickSelect(long double* arr, int* indices, int left, int right, int k, Neighbor*
2     result) {
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                 if(arr[i] < 1) {
10                    result[i].distance = arr[i];
11                } else {
12                    result[i].distance = sqrt(arr[i]);
13                }
14                result[i].index = indices[i];
15            }
16            return;
17        } else if(k-1 < pivotIndex) {
18            quickSelect(arr, indices, left, pivotIndex-1, k, result);
19        } else {
20            quickSelect(arr, indices, pivotIndex+1, right, k, result);
21        }
22    }
23 }
```

Listing 3: Quick Select algorithm implementation

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 × corpus) does not fit in memory. To combat this, we calculate matrix *D* in segments (chunks), use each chunk to find the *k*-NN of the corresponding queries, and then discard the chunk 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* chunk and running quickselect on it. We also use parallelism in order to calculate distances quicker. Chunk size is set to be small enough to be easily stored in memory and large enough to avoid unnecessary overhead. That way we also speed up the calculations significantly.

```
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     int const chunk_size = 300;
8     int num_chunks = (q + chunk_size-1) / chunk_size;
9
10    #pragma omp parallel for
11    for(int chunk = 0; chunk<num_chunks; chunk++) {
12
13        int start_idx = chunk*chunk_size;
14        int end_idx = (start_idx+chunk_size > q) ? q : start_idx+chunk_size;
15
16        long double* D = (long double*)malloc((end_idx-start_idx)*c*sizeof(long double));
17        memory_check(D);
18
19        calculate_distances(C, Q, start_idx, end_idx, D);
20
21        for(int i=start_idx; i<end_idx; i++) {
22
23            int query_idx = i - start_idx;
24
25            int* indices = (int*)malloc(c*sizeof(int));
26            memory_check(indices);
27
28            for(int j=0; j<c; j++) {
29                indices[j] = j;
30            }
31
32            quickSelect(D + query_idx*c, indices, 0, c-1, k, N + i*k);
33
34            free(indices);
35        }
36
37        free(D);
38    }
39 }
```

Listing 4: 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, with minimal error.

The Johnson–Lindenstrauss Lemma asserts that, given a set of points in a high-dimensional space, it is possible to project them onto a lower-dimensional space such, that the Euclidean distances between the points are approximately preserved. Specifically, for a set of n points in d -dimensional space, we can project them onto a subspace of t -dimensions where

$$t = O\left(\frac{\log n}{\epsilon^2}\right),$$

and the distances between the points will be distorted by no more than a factor of $1 \pm \epsilon$, with high probability. This dimensionality reduction is crucial for large datasets, as it allows us to reduce the size of the problem without a significant loss in the accuracy of distance calculations.

```
1 void random_projection(Mat* C, Mat* Q, int t, Mat* C_RP, Mat* Q_RP) {
2
3     srand(time(NULL) + (uintptr_t)C);
4
5     int c = (int)C->rows;
6     int q = (int)Q->rows;
7     int d = (int)C->cols;
8
9     C_RP->rows = c;
10    C_RP->cols = t;
11    Q_RP->rows = q;
12    Q_RP->cols = t;
13
14    C_RP->data = (double*)malloc(c*t*sizeof(double));
15    Q_RP->data = (double*)malloc(q*t*sizeof(double));
16
17    double* R = (double*)malloc(d*t*sizeof(double));
18    for (int i=0; i<d*t; i++) {
19        // Rademacher distribution (-1 or +1)
20        R[i] = (rand()%2 == 0 ? -1 : 1) / sqrt((double)t);
21    }
22
23    cblas_dgemm(CblasRowMajor, CblasNoTrans, CblasNoTrans, c, t, d, 1.0, C->data, d, R, t,
24                0.0, C_RP->data, t);
25    cblas_dgemm(CblasRowMajor, CblasNoTrans, CblasNoTrans, q, t, d, 1.0, Q->data, d, R, t,
26                0.0, Q_RP->data, t);
27
28    free(R);
29 }
```

Listing 5: Routine for applying Random Projection

When $c \gg d$, we cannot reduce the problem with random projection, as the target dimension t would exceed the original dimensionality d . Instead, we truncate the matrix C , retaining a set of representative rows, using k-means++ initialization and k-means clustering.

k-means clustering: Clustering is a technique that aims to group a set of data points into clusters, such that points within a cluster are similar. The similarity is typically measured using a distance metric, in this case the Euclidean distance. The idea is to partition the dataset into clusters in a way that maximizes the cohesion within clusters and the separation between different clusters.

1. **Initialization:** Randomly select k points from the dataset to serve as initial cluster centroids.
2. **Assignment step:** Each data point is assigned to the nearest centroid based on Euclidean distance, thus forming k clusters.
3. **Update step:** The centroids are updated by calculating the mean of all points in each cluster. This new mean becomes the new centroid for that cluster.
4. **Repeat:** The assignment and update steps are repeated until convergence, which occurs when the centroids no longer change significantly.

k-means++ initialization: The main challenge in k-means is selecting good initial centroids, as poor initialization can lead to suboptimal clustering results. To address this, k-means++ was introduced as an enhancement to k-means to select better initial centroids and improve the algorithm’s convergence rate.

1. **First centroid selection:** The first centroid is chosen randomly from the dataset.
2. **Subsequent centroids selection:** For each remaining centroid, the probability of selecting a data point is proportional to its squared distance from the nearest already selected centroid. This ensures that the selected centroids are spread out across the data, capturing the overall structure of the dataset.
3. **Repeat:** This process is repeated until k centroids have been selected.

```
1 void truncMat(Mat* C, int r, Mat* C_TR) {
2
3     int c = C->rows;
4     int d = C->cols;
5
6     C_TR->rows = r;
7     C_TR->cols = d;
8     C_TR->data = (double*)malloc(r*d*sizeof(double));
9
10    srand(time(NULL));
11
12    int first_idx = rand()%c;
13    memcpy(C_TR->data, C->data + first_idx*d, d*sizeof(double));
14
15    double* distances = (double*)malloc(c*sizeof(double));
16    #pragma omp parallel for
17    for(int i=0; i<c; i++) {
18        double dist = 0.0;
19        for(int j=0; j<d; j++) {
20            double diff = C->data[i*d + j] - C_TR->data[j];
21            dist += diff*diff;
22        }
23        distances[i] = dist;
24    }
25
26    for(int i=1; i<r; i++) {
27        double total_dist = 0.0;
28
29        #pragma omp parallel for reduction(+:total_dist)
30        for(int j=0; j<c; j++) {
31            total_dist += distances[j];
32        }
33
34        double rand_dist = ((double)rand()/RAND_MAX) * total_dist;
35        double cumulative_dist = 0.0;
36        int next_idx = 0;
37
38        for(int j=0; j<c; j++) {
39            cumulative_dist += distances[j];
40            if(cumulative_dist >= rand_dist) {
41                next_idx = j;
42                break;
43            }
44        }
45
46        memcpy(C_TR->data + i*d, C->data + next_idx*d, d*sizeof(double));
47
48        #pragma omp parallel for
49        for(int j=0; j<c; j++) {
50            double dist = 0.0;
51            for(int k=0; k<d; k++) {
52                double diff = C->data[j*d + k] - C_TR->data[i*d + k];
53                dist += diff*diff;
54            }
55            if(dist < distances[j]) {
56                distances[j] = dist;
57            }
58        }
59    }
60
61    free(distances);
62 }
```

Listing 6: Matrix Truncation

It is of great importance to choose the target dimensionality t carefully when using random projection. Based on the lemma, $t = \frac{\log(c)}{\epsilon^2}$, 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. Note that target dimensionality t does not depend on the actual dimensionality d . Moreover, the number of representative rows selected when truncating the matrix C affects the accuracy of our calculations. A heuristic metric that balances both speed and accuracy is given by $r = 100 \cdot \log(c) + 10d$.

Applying these minimization techniques is valuable only when the dataset is large. If the dataset is relatively small, we proceed with exact calculation.

```
1  double const e = 0.3;
2  int      const t = log((double)c) / (e*e);
3  if(t<d && (c>1000 && d>50)) {
4
5      Mat C_RP, Q_RP;
6      printf("Target dimension (t) for random projection: %d\n", t);
7      random_projection(&C, &Q, t, &C_RP, &Q_RP);
8
9      findKNN(&C_RP, &Q_RP, N, k);
10
11     free(C_RP.data);
12     free(Q_RP.data);
13 } else if(c>100000) {
14
15     Mat C_TR;
16     int const r = (int)(100*log(c)) + 10*d;
17     printf("Representative rows (r): %d\n", r);
18     truncMat(&C, r, &C_TR);
19
20     findKNN(&C_TR, &Q, N, k);
21
22     free(C_TR.data);
23 } else {
24
25     printf("Exact calculation\n");
26
27     findKNN(&C, &Q, N, k);
28     print_neighbors(N, q, k);
29 }
```

Listing 7: Approximation Check

3.3 Details

To improve computational efficiency, we avoid calculating the square root of distances in the ‘calculate-distances’ function. Instead, squared distances are compared directly within the ‘quickselect’ function. This is mathematically valid because, for any $a > 0$ and $b > 0$, where $a, b > 1$, the inequality $a > b$ is equivalent to $a^2 > b^2$. The square root is computed only for the final k -nearest neighbors. In the corner case where $distance < 1$ we compute the square root right away. This approach reduces the number of square root calculations from $c \times q$ (one for each distance) to $k \times q$, yielding a significant performance improvement.

While this method enhances runtime efficiency, it requires careful handling of numerical limits. The IEEE 754 double-precision format can reliably represent numbers up to approximately 10^{308} . However, squaring large elements in matrices C and Q may cause intermediate values to approach or exceed this limit, potentially resulting in overflow. To mitigate this, we use the ‘long double’ data type, which provides a wider range (up to approximately 10^{4932} in most implementations) to safely handle large squared norms.

4 Testing

To evaluate the performance and validity of our implementation, we employed three metrics: **Accuracy**, **Queries per Second**, and **Execution Time**. We also used Matlab’s kNN routine for further testing.

4.1 Accuracy

To assess the algorithm's accuracy, we selected a random sample of queries and calculated their exact *k*-nearest neighbors. The accuracy was determined by computing the deviation of approximate results from exact results using the formula:

$$\text{Error} = \frac{|d_{\text{approx}} - d_{\text{exact}}|}{|d_{\text{exact}}|}$$

The cumulative error for all samples was averaged, and the accuracy was computed as:

$$\text{Accuracy} = 1 - \text{Average Error}$$

We avoided recall-based metrics that rely on index comparisons because the minimization techniques (e.g., representative rows) modify the indexing, making recall unsuitable. This adjustment ensures accuracy is evaluated based solely on the deviation in distances.

```
1 double recall(Mat* C, Mat* Q, Neighbor* N, int k) {
2
3     double accuracy = 0.0;
4
5     Mat Q_TEST;
6     Q_TEST.rows = SAMPLE;
7     Q_TEST.cols = Q->cols;
8     Q_TEST.data = (double*)malloc(SAMPLE * Q->cols * sizeof(double));
9     memory_check(Q_TEST.data);
10    memcpy(Q_TEST.data, Q->data, SAMPLE * Q->cols * sizeof(double));
11
12    Neighbor* N_TEST = (Neighbor*)malloc(SAMPLE * k * sizeof(Neighbor));
13    memory_check(N_TEST);
14    findKNN(C, &Q_TEST, N_TEST, k);
15
16    for(int i=0; i<SAMPLE; i++) {
17        qsort(N + i*k, k, sizeof(Neighbor), compare);
18        qsort(N_TEST + i*k, k, sizeof(Neighbor), compare);
19    }
20
21    double error = 0.0;
22    for(int i=0; i<SAMPLE; i++) {
23        for(int j=0; j<k; j++) {
24            double approx = N[i*k + j].distance;
25            double exact = N_TEST[i*k + j].distance;
26
27            if(exact != 0) {
28                error += fabs(approx-exact) / fabs(exact);
29            }
30        }
31    }
32
33    double average = error / (SAMPLE*k);
34    accuracy = 1.0 - average;
35
36    free(Q_TEST.data);
37    free(N_TEST);
38
39    return accuracy;
40 }
```

Listing 8: Routine for Calculating Accuracy

4.2 Queries per Second

This metric measures the algorithm's throughput by dividing the total number of queries processed by the elapsed time. It reflects how many queries can be resolved per second.

```
1 double qps(struct timespec start, size_t q) {  
2  
3     double elapsed = duration(start);  
4     return q/elapsed;  
5 }
```

Listing 9: Routine for Calculating Queries per Second

4.3 Execution Time

Execution time is the total time taken by the algorithm to process the dataset, from the start of computation to the end of processing. It is a straightforward yet critical metric for benchmarking performance.

```
1 double duration(struct timespec start) {  
2  
3     struct timespec end;  
4     clock_gettime(CLOCK_MONOTONIC, &end);  
5  
6     double elapsed = (end.tv_sec-start.tv_sec) + (end.tv_nsec-start.tv_nsec)/1e9;  
7     return elapsed;  
8 }
```

Listing 10: Routine for Calculating Elapsed Time

5 Summary

This project implements an efficient k-NN algorithm by combining advanced matrix operations, sorting techniques, and approximation methods. By sacrificing a small degree of accuracy, the algorithm achieves substantial speedups, especially in high-dimensional datasets. Additionally, parallelism is leveraged to further improve scalability. This balanced approach demonstrates the practicality of neighbor searches in large-scale and high-dimensional scenarios, ensuring both computational efficiency and flexibility.

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, as well as Matlab.
4. The parallel programming libraries OpenMP, OpenCilk, and Pthreads.
5. The version control software Git.
6. The AI assistant Github Copilot.
7. The GNU/Linux OS.
8. The Neovim text editor.

7 Sources

C Data Types:

- https://en.wikipedia.org/wiki/Long_double
- https://www.tutorialspoint.com/cprogramming/c_data_types.htm

Random Projection:

- https://en.wikipedia.org/wiki/Random_projection

k-means Clustering:

- https://en.wikipedia.org/wiki/K-means_clustering

Johnson-Lindenstrauss Lemma:

- https://en.wikipedia.org/wiki/Johnson%E2%80%93Lindenstrauss_lemma
- <https://cs.stanford.edu/people/mmahoney/cs369m/Lectures/lecture1.pdf>
- <https://www.youtube.com/watch?v=j9qbuGSjzeE>

ANN Benchmarks:

- <https://github.com/erikbern/ann-benchmarks?tab=readme-ov-file>