

KD-Tree Coding Exercise

Uber Advanced Technologies Center is seeking **highly-skilled C++ developers** to join us in solving our ambitious objectives. This exercise is designed to assess your C++ skills by asking you to implement a data structure for spatial search: the KD-tree.

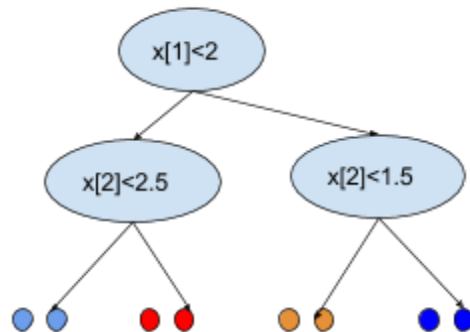
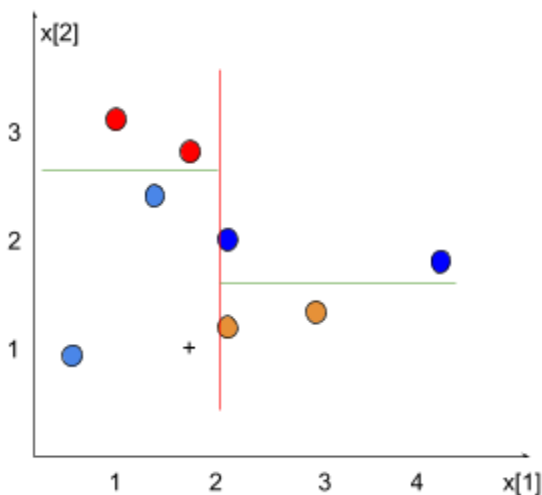
Background

A [KD-tree](#) defines a recursive binary space partitioning, designed for performing efficient spatial queries such as nearest neighbor search. Given a set of K-dimensional points $\{X_1, X_2, \dots, X_N\}$, each internal node in the binary tree defines a partitioning of the space – an axis-aligned hyperplane – that “splits” the points associated with that node, such that points in the left subtree are less than the splitting plane, and points in the right subtree are greater than the splitting plane. This splitting procedure continues recursively, until some minimal number of points remain (possibly one), which are stored explicitly, by reference, or by index in the leaves of the tree.

A key consideration in producing a well-balanced tree is the choice of splitting axis and position at each node. One heuristic for selecting the splitting axis is to simply cycle through the axes in order. Another is to select the axis that has the greatest range or variance. Using either heuristic, the split position can then be chosen as the median (or median-of-medians) of the point values along that axis.

Given a query point, Y , finding the nearest neighbor in the tree requires a depth-first search. At each node in the traversal, we test whether the value of the query Y in the node’s splitting axis (e.g. $Y[j]$) is less or greater than the split position, and descend into the left or right subtree respectively until a leaf node is reached.

A greedy search algorithm would terminate at this point, returning the X stored in the leaf node. This is suboptimal, however, since a closer point *may in fact* exist on the other side of the splitting hyperplanes. Consider the query in the figure below, shown as the ‘+’ point. Greedily descending through the tree would yield the light blue points, where in fact the orange point is closer.



Therefore, during traversal it is necessary to maintain a bound on how far the nearest neighbor could be from the query point. At internal nodes, if the bound is larger than the distance from the query point to the splitting point (indicating that points on the other side of the hyperplane may present viable candidates), it is necessary to search both subtrees. By keeping track of the best leaf node found to date, branches of the tree that are too far from the query point (i.e., the distance to the hyperplane is larger than the best found node to date) can be pruned/ignored entirely.

The Task

The exercise below is your opportunity to show us your C++ design and development skills. We are looking for candidates who are able to take the KD-tree description and turn it into **production quality** code: high performance, well written, well tested, robust pieces of C++.

We will evaluate your solution based on the following guidelines:

- ☐ Your solution should be your own code, or clearly attributed if you use other code sources.
- ☐ Use of the C++ standard library where appropriate.
- ☐ Evidence of C++11 constructs, design patterns, RAII, const correctness and good programming practices to make your code production quality and readily understandable.
- ☐ Evidence of established correctness of components within your program (e.g. unit tests).
- ☐ Understanding of how to use class and/or function templates.
- ☐ A clean, usable, and extensible API.

With these considerations in mind, your task is to develop a KD-tree library that provides at least the following interface/capabilities:

- ☐ Create a tree from a set of K-dimensional points. The tree should be created by recursively choosing a good axis (dimension) and position to split on. Choose an appropriate measure of “goodness” (e.g. largest range, variance, cycle) and justify your decision. Likewise for the split position (e.g. median, median-of-medians).
- ☐ Support efficient *exact* query for the nearest stored point to a query point. That is, search through the tree efficiently and return an iterator or reference to the nearest point.
- ☐ Support I/O to save/load the tree from disk.
- ☐ Your implementation should be templated on the scalar type of the point such that it supports at least `float` and `double` precision types.

Bonus Points:

- ☐ Allow the split axis and split position selection algorithms to be chosen at compile-time or runtime.

Using your library, you should develop two applications:

- ☐ `build_kdtree`: An application to build the KD-tree from a sample dataset (a simple CSV file – details below). It should build the tree and save it to a user-specifiable location.
- ☐ `query_kdtree`: An application to load a KD-tree generated by `build_kdtree` and a CSV file containing query points, and compute the nearest neighbors of each of the points. The application should output for each query point the exact nearest index (0-based) in the sample dataset, and the corresponding Euclidean distance from the query point to the nearest neighbor.

File Formats

We provide a set of test points and query points attached. The file format is CSV with a comma delimiter:

```
x00, x01, x02, ... x0N
x10, x11, x12, ... x1N
...
```

Each line contains data for a single point. You can safely assume the file is formatted correctly and do not need to develop an extensive parsing system. The point ID is the 0-based row number.

Your application should take in a set of query points and produce a file output that consists for each query point, the corresponding closest index of the sample data and the Euclidean distance from the query point to the nearest neighbor:

```
query0_closest_index, query0_minimum_distance
query1_closest_index, query1_minimum_distance
...
```

Parameters

This exercise is not intended to take more than 8 hours, so try to size your plans accordingly.

Please submit your source code, any associated build framework (e.g., cmake), testing, and a README clearly listing any other dependencies. Your implementation must be written in C++ (C++11 preferred) and compile with GNU g++. Your software should compile and run on a vanilla Ubuntu 14.04 LTS 64-bit Linux system. You may use additional dependencies, but clearly document them and make sure they are the vanilla install packages or provide source and make sure it builds as part of your build system. Your README should also detail how to build and how to run the software.

Additionally, please submit a brief pdf document that describes your implementation and design decisions undertaken. Detail the computational efficiency of your implementation and explain how you could make it better given more time. Good luck!