# k-NN and range search with kd-trees

Dominika Kubániová

November 11, 2019

# kd-tree

- Binary tree
- Data points $\in \mathbb{R}^k$ (my case k is max 4)
- Root wraps data into hyperrectange (cell)
- Hyperrectangle iteratively splitted by orthogonal planes

# Non-optimized structures

```
1  struct kd_node {
2    kd_node * left , * right ;
3    byte split_dim ;
4    float split_value ;
5    bool is_leaf ;
6    float * p ;
7  }
```

Note: Some implementations store points in nodes and leaves, some only in leaves

# Search problems

- k-NN search
- Range search (rectangular and spherical queries)

# Construction of a kd-tree

```
1  fun construct_kdtree(dataset of size N in dimension k):
2      if dataset size < M:
3          return nullptr
4      kd_node * node
5      [split_dim, split_value] = choose_split()
6      // ...
7      data_left = {p from dataset | p[dim] < value}
8      data_left = {p from dataset | p[dim] > value}
9      node->left = construct_kdtree(data_left )
10     node->right = construct_kdtree(data_right)
11     return node
12
```

Data should be sorted in every dimension

# How to choose splitting dimension?

- Does not affect correctness of the algorithm
- Can increase complexity of the search
- Balanced tree
- We want ratio of the longest to shortest side to be bounded
- The number of visited leaves in search should be small
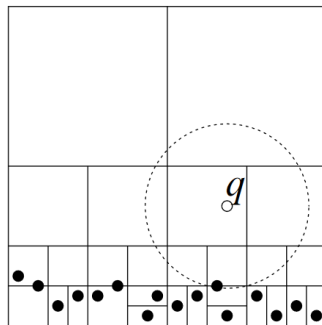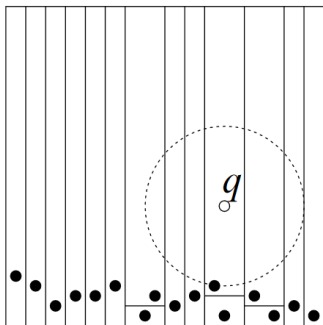- Bad dataset distribution can ruin everything

# Dimension with maximal variance

- Splitting value = median
- Pros:
  - O(N) leaves and O(log N) depth
  - Probably every next chosen dimension will be different (not true when points have bad distribution)
- Cons:
  - Unbounded side length ratio (median does not need to be in the middle of the hyperrectangle)

# Midpoint split

- Splitting value = half of the longest side of the hyperrectangle
- Balanced
- Pros:
  - Side length ratio bounded
- Cons:
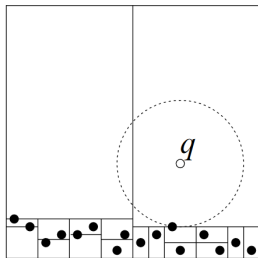  - Possibly a lot of empty nodes

# Max variance and Midpoint



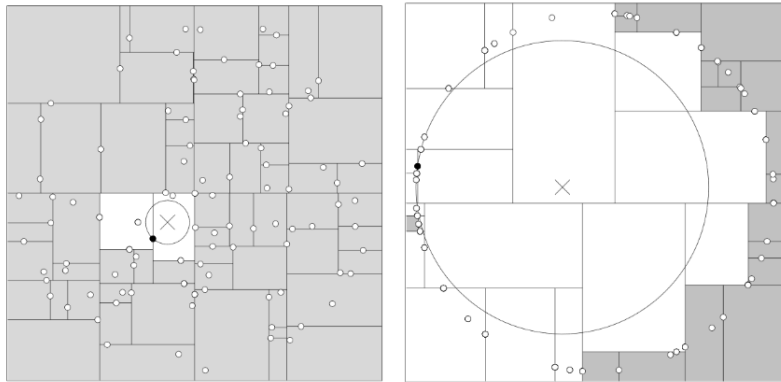Source [2]
Large number of leaves is visited in k-NN search!

# Sliding midpoint

- First step as midpoint
- If points lie on both sides of this split $\rightarrow$ continue
- Else $\rightarrow$ slide towards the closest point creating one NON-EMPTY leaf from it
- Pros:
  - O(N) leaves
- Cons:
  - What about its side length ratio? Is it balanced?



Source: [2]

It is observed that the depth of the tree is not that dominant factor than the number of visited leaves. Example:



Source: [2]

# k-NN search (with PQ)

```
1  Point query
2  float knn_dist = INF
3  PQ kNNq
4  PQ q
5
6  fun process_leaf(n):
7    for p in n->points:
8      if (kNNq.size) < k:
9        kNNq.push(p)
10     else if dist(p, query) < knn_dist:
11       kNNq.push(p)
12       if kNNq.size() > k:
13         kNNq.pop_last()
14         knn_dist = dist(kNNq.last, query)
15
```

# k-NN search (with PQ) cont.

```
1  q.push([root, 0.0])
2  while !q.empty():
3    [n, dist] = q.pop()
4    if dis < knn_dist:
5      while (!n->is_leaf):
6        off_split_plane = query[n->split_dim] - n->
       split_value
7        if off_split_plane < 0:
8          off_box_boundary = query[n->split_dim] - n->
       bound
9          if off_box_boundary > 0:
10           off_box_boundary = 0
11         dist_right = dist - off_box_boundary^2 +
       off_split_plane^2
12         q.push([n->right, dist_right])
13         n = n->left
14       else:
15         // analogously
16     process_leaf(n)
17
```

# Range search (spherical query)

```
1  Point query
2  float r_sq
3  float offs[k] // initialized 0
4  vector<Point> in_range
5  stack s
6
7  fun process_leaf(n):
8    for p in n->points:
9      if dist(p, query) < r_sq:
10       in_range.push(p)
11
```

# Range search (spherical query) cont.

```
1  s.push([root, dist, offs])
2  while !s.empty():
3    [n, dist, off] = s.pop()
4    if n->is_leaf:
5      process_leaf(n)
6    else:
7      old_off_split_plane = off[n->dim]
8      new_off_split_plane = query[n->dim] - n->
       split_value
9      if new_off_split_plane < 0:
10       new_dist = dist - old_off_split_plane^2 +
       new_off_split_plane^2
11       if new_dist < r_sq:
12         new_off = off
13         new_off[n->dim] = new_off_split_plane
14           s.push([n->right, dist, new_off])
15         s.push([n->left, dist, off])
16       else:
17         // analogously
18
```

# Range search (rectangular query)

```
1  Point query
2  float size[k] // only halves stored
3  float offs[k] // initialized 0
4  vector<Point> in_range
5  stack s
6
7  fun process_leaf(n):
8    for p in n->points:
9      dim = n->dim
10     for i in range(k):
11       if query[dim] - size[dim] < p[dim] < query[dim] +
        size[dim]
12         in_range.push(p)
13       else:
14         break
15
```

# Range search (rectangular query) cont.
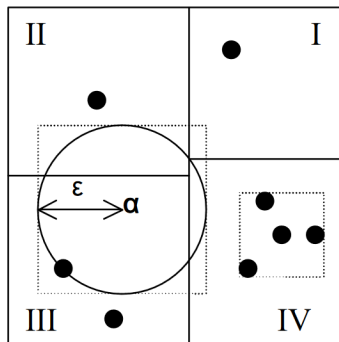
```
1  s.push(root)
2  while !s.empty():
3    n = s.pop()
4
5    if n->is_leaf:
6      process_leaf(n)
7    else:
8      if n->value < query[n->dim] - size[n->dim]:
9        s.push(n->right)
10
11     else if n->value > query[n->dim] + size[n->dim]:
12        s.push(n->left)
13
14     else:
15        s.push(n->left)
16        s.push(n->right)
17
```

# Optimized structures

```
1  struct kd_node {
2      union {
3          kd_node * left;
4          float * p;
5      }
6      byte is_leaf_split_dim;
7      float split_value;
8  }
9
```

# Other simple heuristics

▶ Bounds-Overlap-Ball (BOB) test
Effective when number of points stored in a leaf is high



Source: [2]

# What next needs to be done

- Dataset generation (various distributions) in 2D-4D
- Dataset size: $10^3 - 10^8$
- Different implementation of priority queue
- Comparison (mainly with naïve algorithm)

# References

- 1 Moore, Andrew. (2004). An Intoductory Tutorial on Kd-Trees.
- 2 Maneewongvatana, Songrit & Mount, David. (2000). It's Okay to Be Skinny, If Your Friends Are Fat.
- 3 Sample, Neal & Haines, Matthew & Arnold, Mark & Purcell, Timothy. (2001). Optimizing Search Strategies in k-d Trees.

Thank you for your attention