

A quick note on $B_2(\cdot)$

Ryan R. Curtin

May 25, 2016

Our goal here is to show that $B_2(\cdot)$ from [1] is a correct bound for nearest neighbor search. We can prove this, but first let's rewrite the bound function itself:

$$B_2(\mathcal{N}_q) = \min\left\{\min_{p \in \mathcal{P}_q} (D_p[k] + \rho(\mathcal{N}_q) + \lambda(\mathcal{N}_q)), \min_{\mathcal{N}_c \in \mathcal{C}_q} (B_2(\mathcal{N}_c) + 2(\lambda(\mathcal{N}_q) - \lambda(\mathcal{N}_c)))\right\}. \quad (1)$$

Theorem 1. $B_2(\mathcal{N}_q)$ gives, for any \mathcal{N}_q , an upper bound on the distance between any descendant point of \mathcal{N}_q and its k -nearest neighbor.

Proof. To prove the correctness of $B_2(\mathcal{N}_q)$, we have to consider two cases: when \mathcal{N}_q is a leaf (has no children), and when \mathcal{N}_q is not a leaf. This strategy resembles induction, where the base case is a leaf.

First, consider when \mathcal{N}_q is a leaf. In this setting, the second min in Equation 1 does not evaluate since $|\mathcal{C}(\mathcal{N}_q)| = 0$. So we only need to consider the first term. Also, when \mathcal{N}_q is a leaf, $\lambda(\mathcal{N}_q) = \rho(\mathcal{N}_q)$ because $\mathcal{P}_q = \mathcal{D}_q^p$ (that is, the set of points held in \mathcal{N}_q is the same as the set of descendant points of \mathcal{N}_q). Thus in this case,

$$B_2(\mathcal{N}_q) = \min_{p \in \mathcal{P}_q} D_p[k] + 2\lambda(\mathcal{N}_q). \quad (2)$$

We can show the correctness here using the triangle inequality. Any point in \mathcal{P}_q is separated from any other points in \mathcal{P}_q by a maximum of $2\lambda(\mathcal{N}_q)$. Thus, if there exists some point p with k -furthest neighbor candidate distance $D_p[k]$, then for any other point p_i in \mathcal{P}_q , then

$$D_{p_i}[k] \leq D_p[k] + d(p, p_i) \quad (3)$$

$$\leq D_p[k] + 2\lambda(\mathcal{N}_q). \quad (4)$$

Thus, $B_2(\mathcal{N}_q)$ is correct when \mathcal{N}_q is a leaf. Now, let us consider the other case, where \mathcal{N}_q is not a leaf. Here we must prove that both sides of Equation 1 are correct. We will consider the first side first, with a similar argument.

Since \mathcal{N}_q is not a leaf, then $\rho(\mathcal{N}_q) \leq \lambda(\mathcal{N}_q)$ (that is, we do not have strict equality). We know that any point in \mathcal{P}_q (any point held in \mathcal{N}_q) is separated from \mathcal{D}_q^p (any descendant point of \mathcal{N}_q) by at most $\rho(\mathcal{N}_q) + \lambda(\mathcal{N}_q)$. Thus, if there exists some point $p \in \mathcal{P}_q$ with k -furthest neighbor candidate distance $D_p[k]$, then for any descendant point $p_i \in \mathcal{D}_q^p$, then

$$D_{p_i}[k] \leq D_p[k] + d(p, p_i) \quad (5)$$

$$\leq D_p[k] + \rho(\mathcal{N}_q) + \lambda(\mathcal{N}_q). \quad (6)$$

Now we may turn to proving the correctness of the second side of Equation 1. Assume that $\mathcal{B}_2(\mathcal{N}_c)$ is valid for each child \mathcal{N}_c of \mathcal{N}_q (that is, it satisfies the statement of the theorem). This means that $\mathcal{B}_2(\mathcal{N}_c)$ is a valid upper bound on the distance between any descendant point of \mathcal{N}_c and its k -nearest neighbor. But we can actually say something slightly stricter due to the way $\mathcal{B}_2(\mathcal{N}_c)$ is constructed: $\mathcal{B}_2(\mathcal{N}_c)$ is a valid upper bound on the distance between *any point that falls into the ball of radius $\lambda(\mathcal{N}_c)$ centered at the center of the node \mathcal{N}_c* and its k -nearest neighbor.

The ball of radius $\lambda(\mathcal{N}_c)$ centered at the center of the node \mathcal{N}_c lies entirely within the ball of radius $\lambda(\mathcal{N}_q)$ centered at the center of the node \mathcal{N}_q . For simplicity for what I'm about to write, call B_i the ball of radius λ_i centered at the center of node \mathcal{N}_i .

Then, for any point $p_q \in B_q$ and any point $p_c \in B_c$, we may construct a valid upper bound u_q on the k -nearest neighbor of p_q :

$$u_q = D_{p_c}[k] + d(p_q, p_c). \tag{7}$$

If $p_q \in B_c$ (that is, p_q not only is contained in the ball B_q but also in B_c) then we may simply pick $p_q = p_c$ so $d(p_q, p_c) = 0$. And if $p_q \notin B_c$, we can pick the closest point in B_c to p_q . The furthest possible distance between any $p_q \in B_q$ and the closest $p_c \in B_c$ is $2\lambda(\mathcal{N}_q) - 2\lambda(\mathcal{N}_c)$. (Maybe it is easiest to see this geometrically, but I don't feel like drawing out the figure for this 'short' response.)

Thus we can conclude that in any situation, $d(p_q, p_c) \leq 2(\lambda(\mathcal{N}_q) - \lambda(\mathcal{N}_c))$. Therefore

$$u_q = D_{p_c}[k] + 2(\lambda(\mathcal{N}_q) + \lambda(\mathcal{N}_c)) \tag{8}$$

and since $\mathcal{B}_2(\mathcal{N}_q)$ is a valid upper bound for any point $p_c \in B_c$, we may simplify to

$$u_q = \mathcal{B}_2(\mathcal{N}_q) + 2(\lambda(\mathcal{N}_q) + \lambda(\mathcal{N}_c)). \tag{9}$$

We know that u_q is valid for any $p_q \in B_q$; thus, we can conclude that the second term in Equation 1 is a valid upper bound on the k -nearest neighbor for any p_q that is a descendant point of \mathcal{N}_q .

Combining upper bounds via min still gives valid upper bounds, so the statement of the theorem holds. \square

References

- [1] Ryan R. Curtin, William B. March, Parikshit Ram, David V. Anderson, Alexander G. Gray, and Charles L. Isbell Jr. Tree-independent dual-tree algorithms. In *Proceedings of The 30th International Conference on Machine Learning (ICML '13)*, pages 1435–1443, 2013.