

Rank Cover Trees for Nearest Neighbor Search

Michael E. HOULE
NII

Michael NETT
RWTH Aachen University, GERMANY

Why

Text, images, market data, biological data, scientific data, and other forms of information are currently being gathered in large data repositories at a rate that greatly outstrips our ability to analyze and to interpret. Together with this explosion of information, the demand for effective methods for searching, clustering, categorizing, summarizing and matching within data sets continues to grow. For such applications, solutions based on similarity search are among the most effective proposed in statistics, pattern recognition, and machine learning. The design and analysis of effective similarity search structures has consequently been the subject of intensive research for many decades.

What

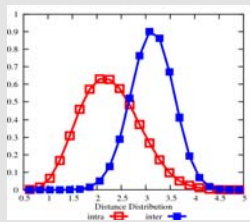
We propose the *rank cover tree* (RCT), a probabilistic data structure for similarity search in general metric spaces. The RCT reinterprets the design and analysis of the Cover Tree in terms of neighbor ranks as measured from the query point, rather than explicit distances. The ranked-based analysis results in a significantly smaller dependence on the intrinsic dimensionality over practical data set sizes. This allows the RCT to find approximations of very high qualities, orders of magnitudes faster than structures for exact similarity search. Moreover, the RCT is highly competitive with the SASH heuristic in terms of its speed-up accuracy trade-off.

Similarity Search

Curse of Dimensionality

Observations

- Sequential scans outperform classical search structures.
- Similarity values concentrate around their mean.
- Similar and dissimilar points are hardly distinguishable.
- Spatial intuitions from 3D spaces are invalid.

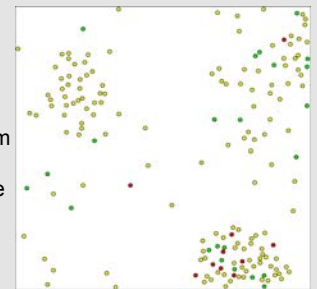


Consequences

- Use approximate similarity search.

Idea of Sampling

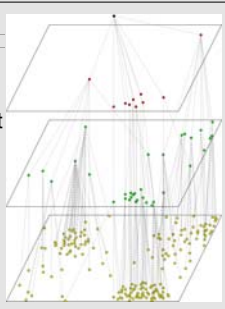
- We try to find items similar to a query q with respect to some data set T .
- Suppose for a subset $S \subseteq T$ we already know an item x that is similar to q .
- An item $y \in T \setminus S$ that is similar to x is likely to be similar to q , and the probability of this is known!



Rank Cover Trees

Construction

- For each item $x \in T$, introduce x into levels $0, \dots, h$. The x 's are geometrically distributed by $p = n^{-1/h}$ with tree height h .
- Build a partial RCT on the highest level by connecting all items in that level to an artificial root node.
- Connect next level by selecting approximate nearest-neighbors found in the partial RCT.
- Errors are amplified, but we can control them!



Search

Find the k items most similar to a query q .

- For each level i , find a cover set C_i .
- Start with C_h containing the artificial root.
- C_i is constructed from C_{i+1} by keeping the k_i children of all elements in C_{i+1} , which are most similar to the query q .

• $k_i = \lfloor \max\{k n^{-1/h}, 1\} \rfloor$, where α is parameter allowing to trade-off between accuracy and query time.

Performance on real data sets

The following figures display the trade-off between approximation quality and query time. Query times are measured as factors of the time consumed by a sequential scan. The query time of the Cover Tree structure for exact similarity searches is provided as a reference point (★). The plots include the RCT for heights 3, 4, 5 and 8 (●, ■, ◆ and ▲) and the SASH heuristic (X).

