

Chapter 40

Finite Metric Spaces and Partitions

By Sarel Har-Peled, April 26, 2022^①

40.1. Finite Metric Spaces

Definition 40.1.1. A *metric space* is a pair (\mathcal{X}, d) where \mathcal{X} is a set and $d : \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty)$ is a *metric*, satisfying the following axioms:

- (i) $d(x, y) = 0 \iff x = y$,
- (ii) $d(x, y) = d(y, x)$, and
- (iii) $d(x, y) + d(y, z) \geq d(x, z)$ (triangle inequality).

The plane, \mathbb{R}^2 , with the regular Euclidean distance is a metric space.

Of special interest is the finite case, where \mathcal{X} is an n -point set. Then, the function d can be specified by $\binom{n}{2}$ real numbers. Alternatively, one can think about (\mathcal{X}, d) as a weighted complete graph, where positive weights are specified on the edges, and these weights comply with the triangle inequality.

Finite metric spaces rise naturally from (sparse) graphs. Indeed, let $G = (\mathcal{X}, E)$ be an undirected weighted graph defined over \mathcal{X} , and let $d_G(x, y)$ be the length of the shortest path between x and y in G . It is easy to verify that (\mathcal{X}, d_G) is a finite metric space. As such if the graph G is sparse, it provides a compact representation to the finite space (\mathcal{X}, d_G) .

Definition 40.1.2. Let (\mathcal{X}, d) be an n -point metric space. We denote the *open ball* of radius r about $x \in \mathcal{X}$, by $b(x, r) = \{y \in \mathcal{X} \mid d(x, y) < r\}$.

Underling our discussion of metric spaces are algorithmic applications. The hardness of various computational problems depends heavily on the structure of the finite metric space. Thus, given a finite metric space, and a computational task, it is natural to try to map the given metric space into a new metric where the task at hand becomes easy.

Example 40.1.3. Computing the diameter of a point set is not trivial in two dimensions (if one wants near linear running time), but is easy in one dimension. Thus, if we could map points in two dimensions into points in one dimension, such that the diameter is preserved, then computing the diameter becomes easy. This approach yields an efficient approximation algorithm, see Exercise 40.7.3 below.

Of course, this mapping from one metric space to another, is going to introduce error. Naturally, one would like to minimize the error introduced by such a mapping.

Definition 40.1.4. Let (\mathcal{X}, d_X) and (\mathcal{Y}, d_Y) be two metric spaces. A mapping $f : \mathcal{X} \rightarrow \mathcal{Y}$ is an *embedding*, and it is *C-Lipschitz* if $d_Y(f(x), f(y)) \leq C \cdot d_X(x, y)$ for all $x, y \in \mathcal{X}$. The mapping f is *K-bi-Lipschitz* if there exists a $C > 0$ such that

$$CK^{-1} \cdot d_X(x, y) \leq d_Y(f(x), f(y)) \leq C \cdot d_X(x, y),$$

^①This work is licensed under the Creative Commons Attribution-Noncommercial 3.0 License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc/3.0/> or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, USA.

for all $x, y \in \mathcal{X}$.

The least K for which f is K -bi-Lipschitz is the *distortion* of f , and is denoted $\text{dist}(f)$. The least distortion with which \mathcal{X} may be embedded in \mathcal{Y} is denoted $c_{\mathcal{Y}}(\mathcal{X})$.

Informally, if $f : \mathcal{X} \rightarrow \mathcal{Y}$ has distortion K , then the distances in \mathcal{X} and $f(\mathcal{X}) \subseteq \mathcal{Y}$ are the same up to a factor of K (one might need to scale up the distances by some constant C).

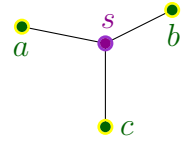
There are several powerful results about low distortion embeddings that would be presented:

- (I) **Probabilistic trees.** Every finite metric can be randomly embedded into a tree such that the “expected” distortion for a specific pair of points is $O(\log n)$.
- (II) **Bourgain embedding.** Any n -point metric space can be embedded into (finite dimensional) euclidean metric space with $O(\log n)$ distortion.
- (III) **Johnson-Lindenstrauss lemma.** Any n -point set in Euclidean space with the regular Euclidean distance can be embedded into \mathbb{R}^k with distortion $(1 + \varepsilon)$, where $k = O(\varepsilon^{-2} \log n)$.

40.2. Examples

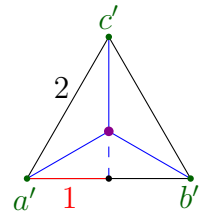
What is distortion? When considering a mapping $f : \mathcal{X} \rightarrow \mathbb{R}^d$ of a metric space $(\mathcal{X}, \mathbf{d})$ to \mathbb{R}^d , it would useful to observe that since \mathbb{R}^d can be scaled, we can consider f to be an expansion (i.e., no distances shrink). Furthermore, we can assume that there is at least one pair of points $x, y \in \mathcal{X}$, such that $\mathbf{d}(x, y) = \|x - y\|$. As such, we have $\text{dist}(f) = \max_{x,y} \frac{\|x-y\|}{\mathbf{d}(x,y)}$.

Why is distortion necessary? Consider the a graph $\mathbf{G} = (\mathbf{V}, \mathbf{E})$ with one vertex s connected to three other vertices a, b, c , where the weights on the edges are all one (i.e., \mathbf{G} is the star graph with three leaves). We claim that \mathbf{G} can not be embedded into Euclidean space with distortion $\leq \sqrt{2}$. Indeed, consider the associated metric space $(\mathbf{V}, \mathbf{d}_{\mathbf{G}})$ and an (expansive) embedding $f : \mathbf{V} \rightarrow \mathbb{R}^d$.



Consider the triangle formed by $\Delta = a'b'c'$, where $a' = f(a), b' = f(b)$ and $c' = f(c)$. Next, consider the following quantity $\max(\|a' - s'\|, \|b' - s'\|, \|c' - s'\|)$ which lower bounds the distortion of f . This quantity is minimized when $r = \|a' - s'\| = \|b' - s'\| = \|c' - s'\|$. Namely, s' is the center of the smallest enclosing circle of Δ . However, r is minimized when all the edges of Δ are of equal length, and are of length $\mathbf{d}_{\mathbf{G}}(a, b) = 2$. It follows that $\text{dist}(f) \geq r \geq 2/\sqrt{3}$.

This quantity is minimized when $r = \|a' - s'\| = \|b' - s'\| = \|c' - s'\|$. Namely, s' is the center of the smallest enclosing circle of Δ . However, r is minimized when all the edges of Δ are of equal length and are of length $\mathbf{d}_{\mathbf{G}}(a, b) = 2$. Observe that the height of the equilateral triangle with sidelength 2 is $h = \sqrt{3}$, and the radius of its inscribing circle is $r = (2/3)h = 2/\sqrt{3}$; see the figure on the right. As such, it follows that $\text{dist}(f) \geq r = 2/\sqrt{3}$.



Note that the above argument is independent of the target dimension d . A packing argument shows that embedding the star graph with n leaves into \mathbb{R}^d requires distortion $\Omega(n^{1/d})$; see Exercise ???. It is known that $\Omega(\log n)$ distortion is necessary in the worst case when embedding a graph into Euclidean space (this is shown using expanders). A proof of distortion $\Omega(\log n / \log \log n)$ is sketched in the bibliographical notes.

40.2.1. Hierarchical Tree Metrics

The following metric is quite useful in practice, and nicely demonstrate why algorithmically finite metric spaces are useful.

Definition 40.2.1. *Hierarchically well-separated tree* (HST) is a metric space defined on the leaves of a rooted tree T . To each vertex $u \in T$ there is associated a label $\Delta_u \geq 0$ such that $\Delta_u = 0$ if and only if u is a leaf of T . The labels are such that if a vertex u is a child of a vertex v then $\Delta_u \leq \Delta_v$. The distance between two leaves $x, y \in T$ is defined as $\Delta_{\text{lca}(x,y)}$, where $\text{lca}(x,y)$ is the least common ancestor of x and y in T .

A HST T is a *k-HST* if for a vertex $v \in T$, we have that $\Delta_v \leq \Delta_{\bar{p}(v)}/k$, where $\bar{p}(v)$ is the parent of v in T .

Note that a HST is a very limited metric. For example, consider the cycle $G = C_n$ of n vertices, with weight one on the edges, and consider an expansive embedding f of G into a HST HST. It is easy to verify, that there must be two consecutive nodes of the cycle, which are mapped to two different subtrees of the root r of HST. Since HST is expansive, it follows that $\Delta_r \geq n/2$. As such, $\text{dist}(f) \geq n/2$. Namely, HSTs fail to faithfully represent even very simple metrics.

40.2.2. Clustering

One natural problem we might want to solve on a graph (i.e., finite metric space) (\mathcal{X}, d) is to partition it into clusters. One such natural clustering is the *k-median clustering*, where we would like to choose a set $C \subseteq \mathcal{X}$ of k centers, such that $v_C(\mathcal{X}, d) = \sum_{u \in \mathcal{X}} d(u, C)$ is minimized, where $d(u, C) = \min_{c \in C} d(u, c)$ is the distance of u to its closest center in C .

It is known that finding the optimal k -median clustering in a (general weighted) graph is NP-complete. As such, the best we can hope for is an approximation algorithm. However, if the structure of the finite metric space (\mathcal{X}, d) is simple, then the problem can be solved efficiently. For example, if the points of \mathcal{X} are on the real line (and the distance between a and b is just $|a - b|$), then k -median can be solved using dynamic programming.

Another interesting case is when the metric space (\mathcal{X}, d) is a HST. Is not too hard to prove the following lemma. See Exercise 40.7.1.

Lemma 40.2.2. *Let (\mathcal{X}, d) be a HST defined over n points, and let $k > 0$ be an integer. One can compute the optimal k -median clustering of \mathcal{X} in $O(k^2n)$ time.*

Thus, if we can embed a general graph G into a HST HST, with low distortion, then we could approximate the k -median clustering on G by clustering the resulting HST, and “importing” the resulting partition to the original space. The quality of approximation, would be bounded by the distortion of the embedding of G into HST.

40.3. Random Partitions

Let (\mathcal{X}, d) be a finite metric space. Given a partition $P = \{C_1, \dots, C_m\}$ of \mathcal{X} , we refer to the sets C_i as *clusters*. We write $\mathcal{P}_{\mathcal{X}}$ for the set of all partitions of \mathcal{X} . For $x \in \mathcal{X}$ and a partition $P \in \mathcal{P}_{\mathcal{X}}$ we denote by $P(x)$ the unique cluster of P containing x . Finally, the set of all probability distributions on $\mathcal{P}_{\mathcal{X}}$ is denoted $\mathcal{D}_{\mathcal{X}}$.

The following partition scheme is due to [CKR04].

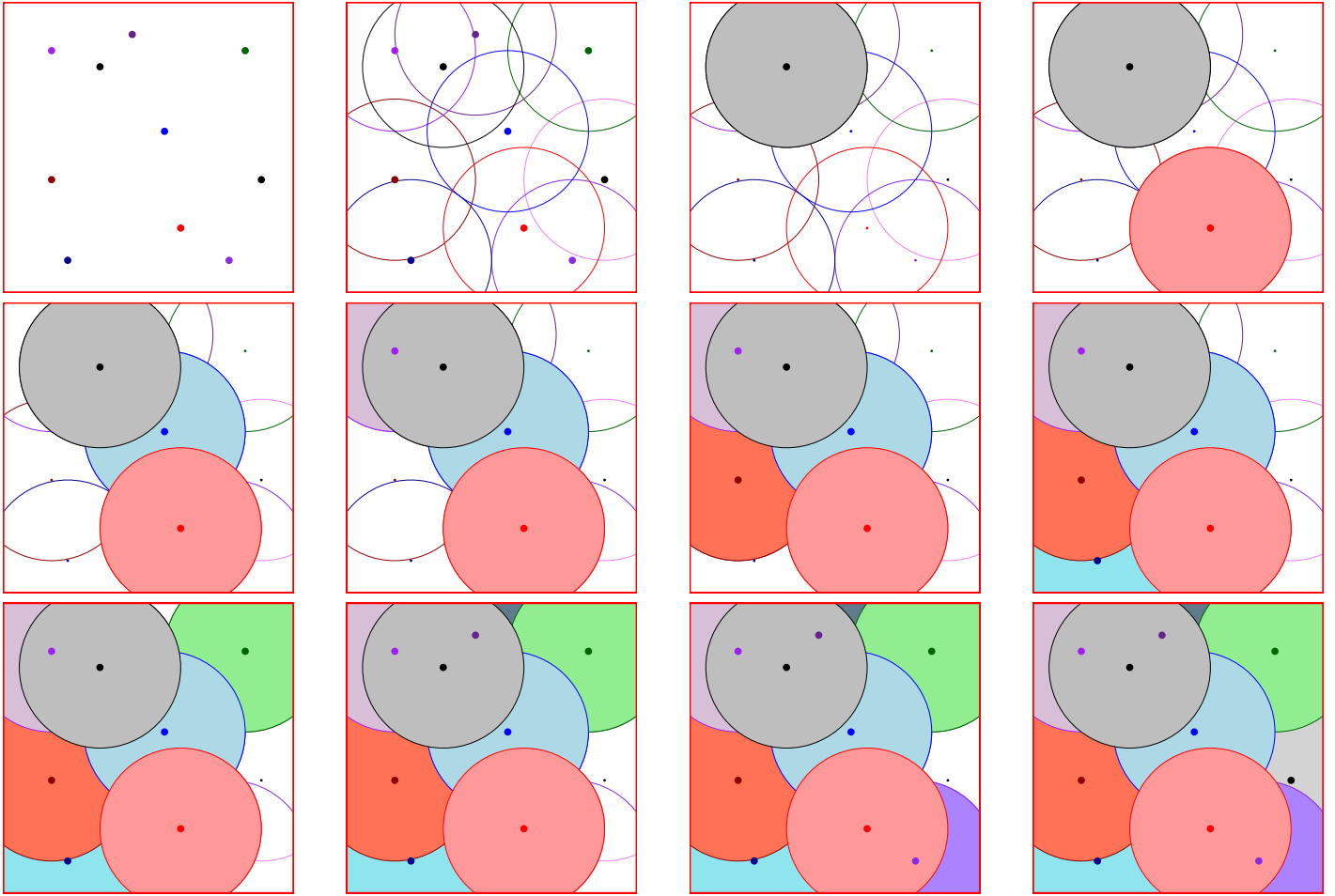


Figure 40.1: An example of the partition of a square (induced by a set of points) as described in Section 40.3.1.

40.3.1. Constructing the partition

Consider a given metric space (\mathcal{X}, d) , where \mathcal{X} is a set of n points.

Let $\Delta = 2^u$ be a prescribed parameter, which is the required diameter of the resulting clusters. Choose, uniformly at random, a permutation π of \mathcal{X} and a random value $\alpha \in [1/4, 1/2]$. Let $R = \alpha\Delta$, and observe that it is uniformly distributed in the interval $[\Delta/4, \Delta/2]$.

The partition is now defined as follows: A point $x \in \mathcal{X}$ is assigned to the cluster C_y of y , where y is the first point in the permutation in distance $\leq R$ from x . Formally,

$$C_y = \{x \in \mathcal{X} \mid x \in \mathbf{b}(y, R) \text{ and } \pi(y) \leq \pi(z) \text{ for all } z \in \mathcal{X} \text{ with } x \in \mathbf{b}(z, R)\}.$$

Let $P = \{C_y\}_{y \in \mathcal{X}}$ denote the resulting partition.

Here is a somewhat more intuitive explanation: Once we fix the radius of the clusters R , we start scooping out balls of radius R centered at the points of the random permutation π . At the i th stage, we scoop out only the remaining mass at the ball centered at x_i of radius r , where x_i is the i th point in the random permutation.

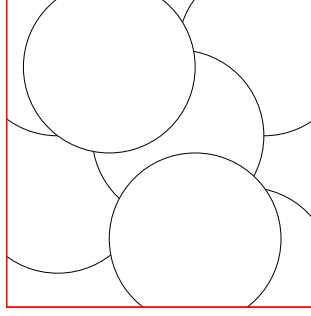


Figure 40.2: The resulting partition.

40.3.2. Properties

The following lemma quantifies the probability of a (crystal) ball of radius t centered at a point x is fully contained in one of the clusters of the partition? (Otherwise, the crystal ball is of course broken.)

Lemma 40.3.1. *Let (\mathcal{X}, d) be a finite metric space, $\Delta = 2^u$ a prescribed parameter, and let P be the partition of \mathcal{X} generated by the above random partition. Then the following holds:*

(i) *For any $C \in P$, we have $\text{diam}(C) \leq \Delta$.*

(ii) *Let x be any point of \mathcal{X} , and t a parameter $\leq \Delta/8$. Then,*

$$\mathbb{P}[\mathbf{b}(x, t) \not\subseteq P(x)] \leq \frac{8t}{\Delta} \ln \frac{b}{a},$$

where $a = |\mathbf{b}(x, \Delta/8)|$, and $b = |\mathbf{b}(x, \Delta)|$.

Proof: Since $C_y \subseteq \mathbf{b}(y, R)$, we have that $\text{diam}(C_y) \leq \Delta$, and thus the first claim holds.

Let U be the set of points of $\mathbf{b}(x, \Delta)$, such that $w \in U$ iff $\mathbf{b}(w, R) \cap \mathbf{b}(x, t) \neq \emptyset$. Arrange the points of U in increasing distance from x , and let $w_1, \dots, w_{b'}$ denote the resulting order, where $b' = |U|$. Let $I_k = [d(x, w_k) - t, d(x, w_k) + t]$ and write \mathcal{E}_k for the event that w_k is the first point in π such that $\mathbf{b}(x, t) \cap C_{w_k} \neq \emptyset$, and yet $\mathbf{b}(x, t) \not\subseteq C_{w_k}$. Note that if $w_k \in \mathbf{b}(x, \Delta/8)$, then $\mathbb{P}[\mathcal{E}_k] = 0$ since $\mathbf{b}(x, t) \subseteq \mathbf{b}(x, \Delta/8) \subseteq \mathbf{b}(w_k, \Delta/4) \subseteq \mathbf{b}(w_k, R)$.

In particular, $w_1, \dots, w_a \in \mathbf{b}(x, \Delta/8)$ and as such $\mathbb{P}[\mathcal{E}_1] = \dots = \mathbb{P}[\mathcal{E}_a] = 0$. Also, note that if $d(x, w_k) < R - t$ then $\mathbf{b}(w_k, R)$ contains $\mathbf{b}(x, t)$ and as such \mathcal{E}_k can not happen. Similarly, if $d(x, w_k) > R + t$ then $\mathbf{b}(w_k, R) \cap \mathbf{b}(x, t) = \emptyset$ and \mathcal{E}_k can not happen. As such, if \mathcal{E}_k happen then $R - t \leq d(x, w_k) \leq R + t$. Namely, if \mathcal{E}_k happen then $R \in I_k$. Namely, $\mathbb{P}[\mathcal{E}_k] = \mathbb{P}[\mathcal{E}_k \cap (R \in I_k)] = \mathbb{P}[R \in I_k] \cdot \mathbb{P}[\mathcal{E}_k | R \in I_k]$. Now, R is uniformly distributed in the interval $[\Delta/4, \Delta/2]$, and I_k is an interval of length $2t$. Thus, $\mathbb{P}[R \in I_k] \leq 2t/(\Delta/4) = 8t/\Delta$.

Next, to bound $\mathbb{P}[\mathcal{E}_k | R \in I_k]$, we observe that w_1, \dots, w_{k-1} are closer to x than w_k and their distance to $\mathbf{b}(x, t)$ is smaller than R . Thus, if any of them appear before w_k in π then \mathcal{E}_k does not happen. Thus, $\mathbb{P}[\mathcal{E}_k | R \in I_k]$ is bounded by the probability that w_k is the first to appear in π out of w_1, \dots, w_k . But this probability is $1/k$, and thus $\mathbb{P}[\mathcal{E}_k | R \in I_k] \leq 1/k$.

We are now ready for the kill. Indeed,

$$\begin{aligned} \mathbb{P}[\mathbf{b}(x, t) \not\subseteq P(x)] &= \sum_{k=1}^{b'} \mathbb{P}[\mathcal{E}_k] = \sum_{k=a+1}^{b'} \mathbb{P}[\mathcal{E}_k] = \sum_{k=a+1}^{b'} \mathbb{P}[R \in I_k] \cdot \mathbb{P}[\mathcal{E}_k | R \in I_k] \\ &\leq \sum_{k=a+1}^{b'} \frac{8t}{\Delta} \cdot \frac{1}{k} \leq \frac{8t}{\Delta} \ln \frac{b'}{a} \leq \frac{8t}{\Delta} \ln \frac{b}{a}, \end{aligned}$$

since $\sum_{k=a+1}^b \frac{1}{k} \leq \int_a^b \frac{dx}{x} = \ln \frac{b}{a}$ and $b' \leq b$. ■

40.4. Probabilistic embedding into trees

In this section, given n -point finite metric $(\mathcal{X}, \mathbf{d})$, we would like to embed it into a HST. As mentioned above, one can verify that for any embedding into HST, the distortion in the worst case is $\Omega(n)$. Thus, we define a randomized algorithm that embed $(\mathcal{X}, \mathbf{d})$ into a tree. Let T be the resulting tree, and consider two points $x, y \in \mathcal{X}$. Consider the *random variable* $\mathbf{d}_T(x, y)$. We constructed the tree T such that distances never shrink; i.e. $\mathbf{d}(x, y) \leq \mathbf{d}_T(x, y)$. The *probabilistic distortion* of this embedding is $\max_{x,y} \mathbb{E} \left[\frac{\mathbf{d}_T(x,y)}{\mathbf{d}(x,y)} \right]$. Somewhat surprisingly, one can find such an embedding with logarithmic probabilistic distortion.

Theorem 40.4.1. *Given n -point metric $(\mathcal{X}, \mathbf{d})$ one can randomly embed it into a 2-HST with probabilistic distortion $\leq 24 \ln n$.*

Proof: The construction is recursive. Let $\text{diam}(P)$, and compute a random partition of \mathcal{X} with cluster diameter $\text{diam}(P)/2$, using the construction of Section 40.3.1. We recursively construct a 2-HST for each cluster, and hang the resulting clusters on the root node v , which is marked by $\Delta_v = \text{diam}(P)$. Clearly, the resulting tree is a 2-HST.

For a node $v \in T$, let $\mathcal{X}(v)$ be the set of points of \mathcal{X} contained in the subtree of v .

For the analysis, assume $\text{diam}(P) = 1$, and consider two points $x, y \in \mathcal{X}$. We consider a node $v \in T$ to be in level i if $\text{level}(v) = \lceil \lg \Delta_v \rceil = i$. The two points x and y correspond to two leaves in T , and let \hat{u} be the least common ancestor of x and y in T . We have $\mathbf{d}_T(x, y) \leq 2^{\text{level}(v)}$. Furthermore, note that along a path the levels are strictly monotonically increasing.

Being more conservative, let w be the first ancestor of x , such that $\mathbf{b} = \mathbf{b}(x, \mathbf{d}(x, y))$ is not completely contained in $\mathcal{X}(u_1), \dots, \mathcal{X}(u_m)$, where u_1, \dots, u_m are the children of w . Clearly, $\text{level}(w) > \text{level}(\hat{u})$. Thus, $\mathbf{d}_T(x, y) \leq 2^{\text{level}(w)}$.

Consider the path σ from the root of T to x , and let \mathcal{E}_i be the event that \mathbf{b} is not fully contained in $\mathcal{X}(v_i)$, where v_i is the node of σ of level i (if such a node exists). Furthermore, let Y_i be the indicator variable which is 1 if \mathcal{E}_i is the first to happen out of the sequence of events $\mathcal{E}_0, \mathcal{E}_1, \dots$. Clearly, $\mathbf{d}_T(x, y) \leq \sum Y_i 2^i$.

Let $t = \mathbf{d}(x, y)$ and $j = \lfloor \lg \mathbf{d}(x, y) \rfloor$, and $n_i = |\mathbf{b}(x, 2^i)|$ for $i = 0, \dots, -\infty$. We have

$$\mathbb{E}[\mathbf{d}_T(x, y)] \leq \sum_{i=j}^0 \mathbb{E}[Y_i] 2^i \leq \sum_{i=j}^0 2^i \mathbb{P} \left[\mathcal{E}_i \cap \overline{\mathcal{E}_{i-1}} \cap \overline{\mathcal{E}_{i-2}} \cdots \overline{\mathcal{E}_0} \right] \leq \sum_{i=j}^0 2^i \cdot \frac{8t}{2^i} \ln \frac{n_i}{n_{i-3}},$$

by Lemma 40.3.1. Thus,

$$\mathbb{E}[\mathbf{d}_T(x, y)] \leq 8t \ln \left(\prod_{i=j}^0 \frac{n_i}{n_{i-3}} \right) \leq 8t \ln(n_0 \cdot n_1 \cdot n_2) \leq 24t \ln n.$$

It thus follows, that the expected distortion for x and y is $\leq 24 \ln n$. ■

40.4.1. Application: approximation algorithm for k -median clustering

Let $(\mathcal{X}, \mathbf{d})$ be a n -point metric space, and let k be an integer number. We would like to compute the optimal k -median clustering. Number, find a subset $C_{\text{opt}} \subseteq \mathcal{X}$, such that $\nu_{C_{\text{opt}}}(\mathcal{X}, \mathbf{d})$ is minimized, see Section 40.2.2. To this end, we randomly embed $(\mathcal{X}, \mathbf{d})$ into a HST using Theorem 40.4.1. Next, using Lemma 40.2.2, we compute the optimal k -median clustering of HST. Let C be the set of centers computed. We return C together with the partition of \mathcal{X} it induces as the required clustering.

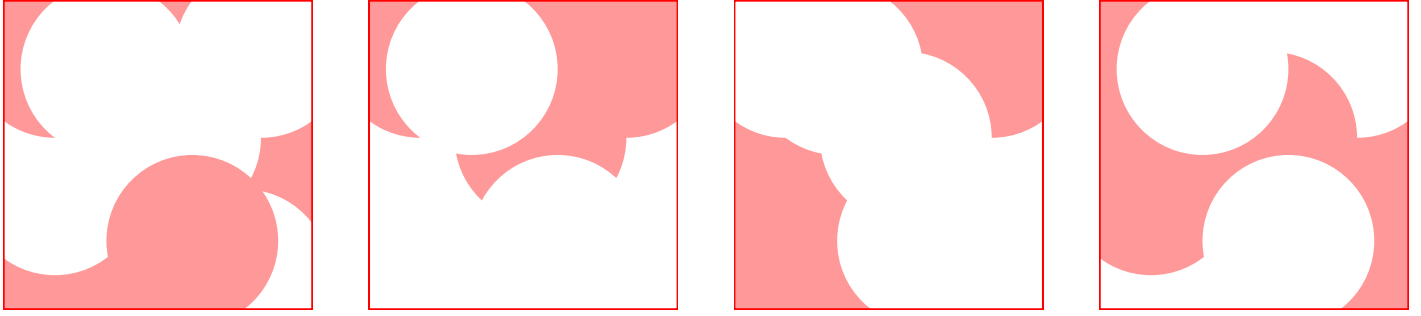


Figure 40.3: Examples of the sets resulting from the partition of Figure 40.1 and taking clusters into a set with probability $1/2$.

Theorem 40.4.2. *Let $(\mathcal{X}, \mathbf{d})$ be a n -point metric space. One can compute in polynomial time a k -median clustering of \mathcal{X} which has expected price $O(\alpha \log n)$, where α is the price of the optimal k -median clustering of $(\mathcal{X}, \mathbf{d})$.*

Proof: The algorithm is described above, and the fact that its running time is polynomial can be easily be verified. To prove the bound on the quality of the clustering, for any point $p \in \mathcal{X}$, let $\text{cen}(p)$ denote the closest point in C_{opt} to p according to \mathbf{d} , where C_{opt} is the set of k -medians in the optimal clustering. Let C be the set of k -medians returned by the algorithm, and let HST be the HST used by the algorithm. We have

$$\beta = \nu_C(\mathcal{X}, \mathbf{d}) \leq \nu_C(\mathcal{X}, \mathbf{d}_{\text{HST}}) \leq \nu_{C_{\text{opt}}}(\mathcal{X}, \mathbf{d}_{\text{HST}}) \leq \sum_{p \in \mathcal{X}} \mathbf{d}_{\text{HST}}(p, C_{\text{opt}}) \leq \sum_{p \in \mathcal{X}} \mathbf{d}_{\text{HST}}(p, \text{cen}(p)).$$

Thus, in expectation we have

$$\begin{aligned} \mathbb{E}[\beta] &= \mathbb{E}\left[\sum_{p \in \mathcal{X}} \mathbf{d}_{\text{HST}}(p, \text{cen}(p))\right] = \sum_{p \in \mathcal{X}} \mathbb{E}[\mathbf{d}_{\text{HST}}(p, \text{cen}(p))] = \sum_{p \in \mathcal{X}} O(\mathbf{d}(p, \text{cen}(p)) \log n) \\ &= O\left((\log n) \sum_{p \in \mathcal{X}} \mathbf{d}(p, \text{cen}(p))\right) = O\left(\nu_{C_{\text{opt}}}(\mathcal{X}, \mathbf{d}) \log n\right), \end{aligned}$$

by linearity of expectation and Theorem 40.4.1. ■

40.5. Embedding any metric space into Euclidean space

Lemma 40.5.1. *Let $(\mathcal{X}, \mathbf{d})$ be a metric, and let $Y \subset \mathcal{X}$. Consider the mapping $f : \mathcal{X} \rightarrow \mathbb{R}$, where $f(x) = \mathbf{d}(x, Y) = \min_{y \in Y} \mathbf{d}(x, y)$. Then for any $x, y \in \mathcal{X}$, we have $|f(x) - f(y)| \leq \mathbf{d}(x, y)$. Namely f is nonexpansive.*

Proof: Indeed, let x' and y' be the closet points of Y , to x and y , respectively. Observe that

$$f(x) = \mathbf{d}(x, x') \leq \mathbf{d}(x, y') \leq \mathbf{d}(x, y) + \mathbf{d}(y, y') = \mathbf{d}(x, y) + f(y)$$

by the triangle inequality. Thus, $f(x) - f(y) \leq \mathbf{d}(x, y)$. By symmetry, we have $f(y) - f(x) \leq \mathbf{d}(x, y)$. Thus, $|f(x) - f(y)| \leq \mathbf{d}(x, y)$. ■

40.5.1. The bounded spread case

Let $(\mathcal{X}, \mathbf{d})$ be a n -point metric. The *spread* of \mathcal{X} , denoted by

$$\Phi(\mathcal{X}) = \frac{\text{diam}(\mathcal{X})}{\min_{x,y \in \mathcal{X}, x \neq y} \mathbf{d}(x,y)},$$

is the ratio between the diameter of \mathcal{X} and the distance between the closest pair of points.

Theorem 40.5.2. *Given a n -point metric $\mathcal{Y} = (\mathcal{X}, d)$, with spread Φ , one can embed it into Euclidean space \mathbb{R}^k with distortion $O(\sqrt{\ln \Phi \ln n})$, where $k = O(\ln \Phi \ln n)$.*

Proof: Assume that $\text{diam}(\mathcal{Y}) = \Phi$ (i.e., the smallest distance in \mathcal{Y} is 1), and let $r_i = 2^{i-2}$, for $i = 1, \dots, \alpha$, where $\alpha = \lceil \lg \Phi \rceil$. Let $P_{i,j}$ be a random partition of P with diameter r_i , using [Theorem 40.4.1](#), for $i = 1, \dots, \alpha$ and $j = 1, \dots, \beta$, where $\beta = \lceil c \log n \rceil$ and c is a large enough constant to be determined shortly.

For each cluster of $P_{i,j}$ randomly toss a coin, and let $V_{i,j}$ be the all the points of \mathcal{X} that belong to clusters in $P_{i,j}$ that got 'T' in their coin toss. For a point $u \in x$, let

$$f_{i,j}(x) = \mathbf{d}(x, \mathcal{X} \setminus V_{i,j}) = \min_{v \in \mathcal{X} \setminus V_{i,j}} \mathbf{d}(x, v),$$

for $i = 0, \dots, m$ and $j = 1, \dots, \beta$. Let $F : \mathcal{X} \rightarrow \mathbb{R}^{(m+1) \cdot \beta}$ be the embedding, such that

$$F(x) = \underbrace{(f_{0,1}(x), f_{0,2}(x), \dots, f_{0,\beta}(x))}_{\text{first } n \text{ resolution block}}, f_{1,1}(x), f_{1,2}(x), \dots, f_{1,\beta}(x), \dots, f_{m,1}(x), f_{m,2}(x), \dots, f_{m,\beta}(x)).$$

Next, consider two points $x, y \in \mathcal{X}$, with distance $\phi = \mathbf{d}(x, y)$. Let k be an integer such that $r_u \leq \phi/2 \leq r_{u+1}$. Clearly, in any partition of $P_{u,1}, \dots, P_{u,\beta}$ the points x and y belong to different clusters. Furthermore, with probability half $x \in V_{u,j}$ and $y \notin V_{u,j}$ or $x \notin V_{u,j}$ and $y \in V_{u,j}$, for $1 \leq j \leq \beta$.

Let \mathcal{E}_j denote the event that $\mathbf{b}(x, \rho) \subseteq V_{u,j}$ and $y \notin V_{u,j}$, for $j = 1, \dots, \beta$, where $\rho = \phi/(64 \ln n)$. By [Lemma 40.3.1](#), we have

$$\mathbb{P}[\mathbf{b}(x, \rho) \not\subseteq P_{u,j}(x)] \leq \frac{8\rho}{r_u} \ln n \leq \frac{\phi}{8r_u} \leq 1/2.$$

Thus,

$$\begin{aligned} \mathbb{P}[\mathcal{E}_j] &= \mathbb{P}[(\mathbf{b}(x, \rho) \subseteq P_{u,j}(x)) \cap (x \in V_{u,j}) \cap (y \notin V_{u,j})] \\ &= \mathbb{P}[\mathbf{b}(x, \rho) \subseteq P_{u,j}(x)] \cdot \mathbb{P}[x \in V_{u,j}] \cdot \mathbb{P}[y \notin V_{u,j}] \geq 1/8, \end{aligned}$$

since those three events are independent. Notice, that if \mathcal{E}_j happens, then $f_{u,j}(x) \geq \rho$ and $f_{u,j}(y) = 0$.

Let X_j be an indicator variable which is 1 if \mathcal{E}_j happens, for $j = 1, \dots, \beta$. Let $Z = \sum_j X_j$, and we have $\mu = \mathbb{E}[Z] = \mathbb{E}[\sum_j X_j] \geq \beta/8$. Thus, the probability that only $\beta/16$ of $\mathcal{E}_1, \dots, \mathcal{E}_\beta$ happens, is $\mathbb{P}[Z < (1 - 1/2) \mathbb{E}[Z]]$. By the Chernoff inequality, we have $\mathbb{P}[Z < (1 - 1/2) \mathbb{E}[Z]] \leq \exp(-\mu/2) \leq \exp(-\beta/16) \leq 1/n^{10}$, if we set $c = 640$.

Thus, with high probability

$$\|F(x) - F(y)\| \geq \sqrt{\sum_{j=1}^{\beta} (f_{u,j}(x) - f_{u,j}(y))^2} \geq \sqrt{\rho^2 \frac{\beta}{16}} = \sqrt{\beta} \frac{\rho}{4} = \phi \cdot \frac{\sqrt{\beta}}{256 \ln n}.$$

On the other hand, $|f_{i,j}(x) - f_{i,j}(y)| \leq \mathbf{d}(x, y) = \phi \leq 64\rho \ln n$. Thus,

$$\|F(x) - F(y)\| \leq \sqrt{\alpha\beta(64\rho \ln n)^2} \leq 64\sqrt{\alpha\beta}\rho \ln n = \sqrt{\alpha\beta} \cdot \phi.$$

Thus, setting $G(x) = F(x) \frac{256 \ln n}{\sqrt{\beta}}$, we get a mapping that maps two points of distance ϕ from each other to two points with distance in the range $\left[\phi, \phi \cdot \sqrt{\alpha\beta} \cdot \frac{256 \ln n}{\sqrt{\beta}}\right]$. Namely, $G(\cdot)$ is an embedding with distortion $O(\sqrt{\alpha} \ln n) = O(\sqrt{\ln \Phi} \ln n)$.

The probability that \mathbf{G} fails on one of the pairs, is smaller than $(1/n^{10}) \cdot \binom{n}{2} < 1/n^8$. In particular, we can check the distortion of \mathbf{G} for all $\binom{n}{2}$ pairs, and if any of them fail (i.e., the distortion is too big), we restart the process. \blacksquare

40.5.2. The unbounded spread case

Our next task, is to extend [Theorem 40.5.2](#) to the case of unbounded spread. Indeed, let (\mathcal{X}, d) be a n -point metric, such that $\text{diam}(\mathcal{X}) \leq 1/2$. Again, we look on the different resolutions r_1, r_2, \dots , where $r_i = 1/2^{i-1}$. For each one of those resolutions r_i , we can embed this resolution into β coordinates, as done for the bounded case. Then we concatenate the coordinates together.

There are two problems with this approach: (i) the number of resulting coordinates is infinite, and (ii) a pair x, y , might be distorted a ‘‘lot’’ because it contributes to all resolutions, not only to its ‘‘relevant’’ resolutions.

Both problems can be overcome with careful tinkering. Indeed, for a resolution r_i , we are going to modify the metric, so that it ignores short distances (i.e., distances $\leq r_i/n^2$). Formally, for each resolution r_i , let $G_i = (\mathcal{X}, \widehat{E}_i)$ be the graph where two points x and y are connected if $\mathbf{d}(x, y) \leq r_i/n^2$. Consider a connected component $C \in G_i$. For any two points $x, y \in C$, we have $\mathbf{d}(x, y) \leq n(r_i/n^2) \leq r_i/n$. Let \mathcal{X}_i be the set of connected components of G_i , and define the distances between two connected components $C, C' \in \mathcal{X}_i$, to be $\mathbf{d}_i(C, C') = \mathbf{d}(C, C') = \min_{c \in C, c' \in C'} \mathbf{d}(c, c')$.

It is easy to verify that $(\mathcal{X}_i, \mathbf{d}_i)$ is a metric space (see [Exercise 40.7.2](#)). Furthermore, we can naturally embed (\mathcal{X}, d) into $(\mathcal{X}_i, \mathbf{d}_i)$ by mapping a point $x \in \mathcal{X}$ to its connected components in \mathcal{X}_i . Essentially $(\mathcal{X}_i, \mathbf{d}_i)$ is a snapped version of the metric (\mathcal{X}, d) , with the advantage that $\Phi((\mathcal{X}_i, \mathbf{d}_i)) = O(n^2)$. We now embed \mathcal{X}_i into $\beta = O(\log n)$ coordinates. Next, for any point of \mathcal{X} we embed it into those β coordinates, by using the embedding of its connected component in \mathcal{X}_i . Let E_i be the embedding for resolution r_i . Namely, $E_i(x) = (f_{i,1}(x), f_{i,2}(x), \dots, f_{i,\beta}(x))$, where $f_{i,j}(x) = \min(\mathbf{d}_i(x, \mathcal{X} \setminus V_{i,j}), 2r_i)$. The resulting embedding is $F(x) = \oplus E_i(x) = (E_1(x), E_2(x), \dots)$.

Since we slightly modified the definition of $f_{i,j}(\cdot)$, we have to show that $f_{i,j}(\cdot)$ is nonexpansive. Indeed, consider two points $x, y \in \mathcal{X}_i$, and observe that

$$|f_{i,j}(x) - f_{i,j}(y)| \leq |\mathbf{d}_i(x, V_{i,j}) - \mathbf{d}_i(y, V_{i,j})| \leq \mathbf{d}_i(x, y) \leq \mathbf{d}(x, y),$$

as a simple case analysis^② shows.

For a pair $x, y \in \mathcal{X}$, and let $\phi = \mathbf{d}(x, y)$. To see that $F(\cdot)$ is the required embedding (up to scaling), observe that, by the same argumentation of [Theorem 40.5.2](#), we have that with high probability

$$\|F(x) - F(y)\| \geq \phi \cdot \frac{\sqrt{\beta}}{256 \ln n}.$$

^②Indeed, if $f_{i,j}(x) < \mathbf{d}_i(x, V_{i,j})$ and $f_{i,j}(y) < \mathbf{d}_i(x, V_{i,j})$ then $f_{i,j}(x) = 2r_i$ and $f_{i,j}(y) = 2r_i$, which implies the above inequality. If $f_{i,j}(x) = \mathbf{d}_i(x, V_{i,j})$ and $f_{i,j}(y) = \mathbf{d}_i(x, V_{i,j})$ then the inequality trivially holds. The other option is handled in a similar fashion.

To get an upper bound on this distance, observe that for i such that $r_i > \phi n^2$, we have $E_i(x) = E_i(y)$. Thus,

$$\begin{aligned} \|F(x) - F(y)\|^2 &= \sum_i \|E_i(x) - E_i(y)\|^2 = \sum_{i, r_i < \phi n^2} \|E_i(x) - E_i(y)\|^2 \\ &= \sum_{i, \phi/n^2 < r_i < \phi n^2} \|E_i(x) - E_i(y)\|^2 + \sum_{i, r_i < \phi/n^2} \|E_i(x) - E_i(y)\|^2 \\ &= \beta \phi^2 \lg(n^4) + \sum_{i, r_i < \phi/n^2} (2r_i)^2 \beta \leq 4\beta \phi^2 \lg n + \frac{4\phi^2 \beta}{n^4} \leq 5\beta \phi^2 \lg n. \end{aligned}$$

Thus, $\|F(x) - F(y)\| \leq \phi \sqrt{5\beta \lg n}$. We conclude, that with high probability, $F(\cdot)$ is an embedding of \mathcal{X} into Euclidean space with distortion $\left(\phi \sqrt{5\beta \lg n}\right) / \left(\phi \cdot \frac{\sqrt{\beta}}{256 \ln n}\right) = O(\log^{3/2} n)$.

We still have to handle the infinite number of coordinates problem. However, the above proof shows that we care about a resolution r_i (i.e., it contributes to the estimates in the above proof) only if there is a pair x and y such that $r_i/n^2 \leq \mathbf{d}(x, y) \leq r_i n^2$. Thus, for every pair of distances there are $O(\log n)$ relevant resolutions. Thus, there are at most $\eta = O(n^2 \beta \log n) = O(n^2 \log^2 n)$ relevant coordinates, and we can ignore all the other coordinates. Next, consider the affine subspace h that spans $F(P)$. Clearly, it is $n - 1$ dimensional, and consider the projection $G : \mathbb{R}^\eta \rightarrow \mathbb{R}^{n-1}$ that projects a point to its closest point in h . Clearly, $G(F(\cdot))$ is an embedding with the same distortion for P , and the target space is of dimension $n - 1$.

Note, that all this process succeeds with high probability. If it fails, we try again. We conclude:

Theorem 40.5.3 (Low quality Bourgain theorem). *Given a n -point metric M , one can embed it into Euclidean space of dimension $n - 1$, such that the distortion of the embedding is at most $O(\log^{3/2} n)$.*

Using the Johnson-Lindenstrauss lemma, the dimension can be further reduced to $O(\log n)$. Being more careful in the proof, it is possible to reduce the dimension to $O(\log n)$ directly.

40.6. Bibliographical notes

The partitions we use are due to Calinescu *et al.* [CKR04]. The idea of embedding into spanning trees is due to Alon *et al.* [AKPW95], which showed that one can get a probabilistic distortion of $2^{O(\sqrt{\log n \log \log n})}$. Yair Bartal realized that by allowing trees with additional vertices, one can get a considerably better result. In particular, he showed [Bar96] that probabilistic embedding into trees can be done with polylogarithmic average distortion. He later improved the distortion to $O(\log n \log \log n)$ in [Bar98]. Improving this result was an open question, culminating in the work of Fakcharoenphol *et al.* [FRT04] which achieve the optimal $O(\log n)$ distortion.

Our proof of Lemma 40.3.1 (which is originally from [FRT04]) is taken from [KLMN05]. The proof of Theorem 40.5.3 is by Gupta [Gup00].

A good exposition of metric spaces is available in Matoušek [Mat02].

Embedding into spanning trees. The above embeds the graph into a Steiner tree. A more useful representation, would be a random embedding into a spanning tree. Surprisingly, this can be done, as shown by Emek *et al.* [EEST08]. This was improved to $O(\log n \cdot \log \log n \cdot (\log \log \log n)^3)$ ^③ by Abraham *et al.* [ABN08a, ABN08b].

^③Truely a polyglot of logs.

Alternative proof of the tree embedding result. Interestingly, if one does not care about the optimal distortion, one can get similar result (for embedding into probabilistic trees), by first embedding the metric into Euclidean space, then reduce the dimension by the Johnson-Lindenstrauss lemma, and finally, construct an HST by constructing a quadtree over the points. The “trick” is to randomly translate the quadtree. It is easy to verify that this yields $O(\log^4 n)$ distortion. See the survey by Indyk [Ind01] for more details. This random shifting of quadtrees is a powerful technique that was used in getting several result, and it is a crucial ingredient in Arora [Aro98] approximation algorithm for Euclidean TSP.

40.7. Exercises

Exercise 40.7.1 (Clustering for HST). Let $(\mathcal{X}, \mathbf{d})$ be a HST defined over n points, and let $k > 0$ be an integer. Provide an algorithm that computes the optimal k -median clustering of \mathcal{X} in $O(k^2 n)$ time.

[Transform the HST into a tree where every node has only two children. Next, run a dynamic programming algorithm on this tree.]

Exercise 40.7.2 (Partition induced metric).

- Give a counter example to the following claim: Let $(\mathcal{X}, \mathbf{d})$ be a metric space, and let P be a partition of \mathcal{X} . Then, the pair (P, \mathbf{d}') is a metric, where $\mathbf{d}'(C, C') = \mathbf{d}(C, C') = \min_{x \in C, y \in C'} \mathbf{d}(x, y)$ and $C, C' \in P$.
- Let $(\mathcal{X}, \mathbf{d})$ be a n -point metric space, and consider the set $U = \{i \mid 2^i \leq \mathbf{d}(x, y) \leq 2^{i+1}, \text{ for } x, y \in \mathcal{X}\}$. Prove that $|U| = O(n)$. Namely, there are only n different resolutions that “matter” for a finite metric space.

Exercise 40.7.3 (Computing the diameter via embeddings).

- (h:1) Let ℓ be a line in the plane, and consider the embedding $f : \mathbb{R}^2 \rightarrow \ell$, which is the projection of the plane into ℓ . Prove that f is 1-Lipschitz, but it is not K -bi-Lipschitz for any constant K .
- (h:3) Prove that one can find a family of projections \mathcal{F} of size $O(1/\sqrt{\varepsilon})$, such that for any two points $x, y \in \mathbb{R}^2$, for one of the projections $f \in \mathcal{F}$ we have $\mathbf{d}(f(x), f(y)) \geq (1 - \varepsilon)\mathbf{d}(x, y)$.
- (h:1) Given a set P of n in the plane, given a $O(n/\sqrt{\varepsilon})$ time algorithm that outputs two points $x, y \in P$, such that $\mathbf{d}(x, y) \geq (1 - \varepsilon)\text{diam}(P)$, where $\text{diam}(P) = \max_{z, w \in P} \mathbf{d}(z, w)$ is the diameter of P .
- (h:2) Given P , show how to extract, in $O(n)$ time, a set $Q \subseteq P$ of size $O(\varepsilon^{-2})$, such that $\text{diam}(Q) \geq (1 - \varepsilon/2)\text{diam}(P)$. (Hint: Construct a grid of appropriate resolution.)

In particular, give an $(1 - \varepsilon)$ -approximation algorithm to the diameter of P that works in $O(n + \varepsilon^{-2.5})$ time. (There are slightly faster approximation algorithms known for approximating the diameter.)

Acknowledgments

The presentation in this write-up follows closely the insightful suggestions of Manor Mendel.

References

- [ABN08a] I. Abraham, Y. Bartal, and O. Neiman. *Nearly tight low stretch spanning trees*. *Proc. 49th Annu. IEEE Sympos. Found. Comput. Sci. (FOCS)*, 781–790, 2008.

- [ABN08b] I. Abraham, Y. Bartal, and O. Neiman. *Nearly tight low stretch spanning trees*. *CoRR*, abs/0808.2017, 2008. arXiv: [0808.2017](https://arxiv.org/abs/0808.2017).
- [AKPW95] N. Alon, R. M. Karp, D. Peleg, and D. West. *A graph-theoretic game and its application to the k -server problem*. *SIAM J. Comput.*, 24(1): 78–100, 1995.
- [Aro98] S. Arora. *Polynomial time approximation schemes for Euclidean TSP and other geometric problems*. *J. Assoc. Comput. Mach.*, 45(5): 753–782, 1998.
- [Bar96] Y. Bartal. *Probabilistic approximations of metric space and its algorithmic application*. *Proc. 37th Annu. IEEE Sympos. Found. Comput. Sci. (FOCS)*, 183–193, 1996.
- [Bar98] Y. Bartal. *On approximating arbitrary metrics by tree metrics*. *Proc. 30th Annu. ACM Sympos. Theory Comput. (STOC)*, 161–168, 1998.
- [CKR04] G. Călinescu, H. J. Karloff, and Y. Rabani. *Approximation algorithms for the 0-extension problem*. *SIAM J. Comput.*, 34(2): 358–372, 2004.
- [EEST08] M. Elkin, Y. Emek, D. A. Spielman, and S. Teng. *Lower-stretch spanning trees*. *SIAM J. Comput.*, 38(2): 608–628, 2008.
- [FRT04] J. Fakcharoenphol, S. Rao, and K. Talwar. *A tight bound on approximating arbitrary metrics by tree metrics*. *J. Comput. Sys. Sci.*, 69(3): 485–497, 2004.
- [Gup00] A. Gupta. *Embeddings of Finite Metrics*. PhD thesis. University of California, Berkeley, 2000.
- [Ind01] P. Indyk. *Algorithmic applications of low-distortion geometric embeddings*. *Proc. 42nd Annu. IEEE Sympos. Found. Comput. Sci. (FOCS)*, Tutorial. 10–31, 2001.
- [KLMN05] R. Krauthgamer, J. R. Lee, M. Mendel, and A. Naor. *Measured descent: a new embedding method for finite metric spaces*. *Geom. funct. anal. (GAFA)*, 15(4): 839–858, 2005.
- [Mat02] J. Matoušek. *Lectures on discrete geometry*. Vol. 212. Grad. Text in Math. Springer, 2002.