

Probabilistic analyses of closest pair algorithms

Mordecai Golin

► **To cite this version:**

Mordecai Golin. Probabilistic analyses of closest pair algorithms. [Research Report] RR-1471, INRIA. 1991. inria-00075091

HAL Id: inria-00075091

<https://hal.inria.fr/inria-00075091>

Submitted on 24 May 2006

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

INRIA

UNITÉ DE RECHERCHE
INRIA-ROCQUENCOURT

Institut National
de Recherche
en Informatique
et en Automatique

Domaine de Voluceau
Rocquencourt
B.P.105
78153 Le Chesnay Cedex
France
Tél.:(1) 39 63 55 11

Rapports de Recherche

N° 1471

Programme 2
Calcul Symbolique, Programmation
et Génie logiciel

PROBABILISTIC ANALYSES OF CLOSEST PAIR ALGORITHMS

Mordecai J. GOLIN

Juin 1991



* RR - 1 4 7 1 *

Probabilistic Analyses of Closest Pair Algorithms (Extended Abstract)

Mordecai J. Golin*
INRIA-Rocquencourt

November 1990

Abstract: In this paper we present three simple fast algorithms for finding the closest pair among a set of n points. We perform probabilistic analyses on the three algorithms under the assumption that the points are chosen I.I.D. from the uniform distribution over the unit square $[0, 1]^2$.

Analyses probabilistes d'algorithmes de paires les plus proches

Résumé : Dans cet article nous présentons trois algorithmes, simples et rapides pour trouver la paire la plus proche parmi n points. Nous faisons l'hypothèse que les n points sont I.I.D. dans le carré $[0, 1]^2$. Notre analyse est probabiliste.

*This work was supported in part by NSF Grant DCR-8605962 and Office of Naval Research Grant N00014-87-K-0460. It was later supported by a Chateaubriand fellowship from the French Ministère des Affaires Étrangères and by the European Community, Esprit II Basic Research Action Program Number 3075 (ALCOM). Author's current address: Algorithms' Project, Inria Rocquencourt, 78153, Le Chesnay Cedex, France.

Probabilistic Analyses of Closest Pair Algorithms

(Extended Abstract)

Mordecai Golin¹
Princeton University

November 1990

Abstract:

In this paper we present probabilistic analyses of some simple, practical algorithms for finding the closest pair among a set of n points. The first algorithm is a simplification of the sweep-line algorithm given in [HNS]. We show that, given n points chosen I.I.D. $U[0, 1]^2$, and then sorted by x -coordinate, the algorithm's i -th step requires an expected $\Theta(\sqrt{\frac{n}{i}})$ amount of time, using $\Theta(n)$ time overall. This analysis can be modified to show that [BP]'s projection method also has the same $\Theta(\sqrt{\frac{n}{i}})$, $\Theta(n)$ time bound. The final algorithm we examine is a modified version of the gridding algorithm, using the floor function, that appears in [BWY]. We show how some minor modifications to the gridding algorithm produces an algorithm which performs only an expected $\Theta(\log^2 n)$ distance calculations.

§1 Introduction

The *closest pair* problem, finding the closest pair among a set of points, is one of the more basic problems in computational geometry. A large number of algorithms, approaching the problem from many different perspectives, exist for its solution. There are “practical” $O(n \log n)$ time algorithms that are relatively simple and easy to program: Shamos and Hoey’s [Sh] divide-and-conquer algorithm and Hinrichs, Nievergelt, and Schorn’s [HNS] sweep-line algorithm are representative. There are two randomized algorithms, one due to Rabin [Ra], the other to Weide [We], that both run in expected linear time on all inputs. There is a “theoretical” algorithm due to Fortune and Hopcroft [FH], really a modification of [Ra], that manages to get rid of the randomization by adding complexity and raising the running time. There are algorithms that solve the closest pair problem as a side effect of solving a more complicated problem such as constructing Voronoi diagrams [SH]. There is also an algorithm, due to Bentley, Weide, and Yao [BWY], that finds the closest pair by using a “gridding method”, i.e. by using the floor function. When run on n points drawn I.I.D. $U[0, 1]^2$, (Independently, Identically Distributed from the uniform distribution over the unit square) this last algorithm requires $O(n)$ expected time.

Our primary purpose in this paper is to examine practical algorithms for solving the closest pair problem in the plane, algorithms that are robust and easy to implement yet run well on average, algorithms that would be implemented to solve real problems for large inputs. The first algorithm that we discuss, a modification of

¹ This work was supported in part by NSF Grant DCR-8605962 and Office of Naval Research Grant N00014-87-K-0460. It was also supported by a Chateaubriand fellowship from the French Ministère des Affaires Étrangères and by the European Community, Esprit II Basic Research Action Program Number 3075 (ALCOM). Author’s current address: Algorithms’ Project, Inria Rocquencourt, 78153, Le Chesnay Cedex, France.

the sweep line algorithm of Hinrichs, Nievergalt and Schorn [HNS], has these characteristics. It consists of two steps. The first step is a sort on the x -coordinates of the input points. The second step is a simple sweep through the points that can be implemented using only nine lines of Pascal code. The algorithm pays a price for its simplicity: the second stage has a worst case running time of $\Theta(n^2)$. It's actual expected running time is much better: we will show that, if the n input points are chosen I.I.D. $\mathbf{U}[0, 1]^2$ then, for all but a negligible number of points, the i 'th step of the sweep requires only expected $\Theta(\sqrt{\frac{n}{i}})$ time. Summing over $i \leq n$ shows that the second stage of the sweep requires only $\sum_i \Theta(\sqrt{\frac{n}{i}}) = \Theta(n)$ expected time. Thus this is a very simple algorithm to implement physically that will run very fast on average. First sort the points by x -coordinate using whatever fine-tuned sorting algorithm you have on hand. Then run the sweep-stage over the sorted points. Our result shows that the dominant amount of work is performed in the sorting stage. This bound remains valid when distance is defined by any L_p metric.

The second algorithm we discuss, Bentley and Papadimitriou's [BP] projection method, is very similar to the first algorithm. This new algorithm is slightly more complicated than the one preceding it but the basic structure is the same: a first step sorting the points and a second step sweeping through them. It is possible to modify the analysis of the first algorithm to yield a $\Theta(\sqrt{\frac{n}{i}})$ expected time bound for the i 'th stage of the sweep step of this algorithm as well, again leading to a $\Theta(n)$ bound for the full sweep stage. This bound remains valid for any L_p metric. This algorithm's advantage over the preceding one is that in [BP] it was proved that it has a worst case running time of $O(n^{1.5})$ when closest is defined in the L_∞ sense as opposed to the $O(n^2)$ bound of the first algorithm. It's disadvantage is that it is slightly more complicated to program and has a higher overhead.

The third algorithm we discuss is a modification of the gridding algorithm presented by Bentley, Weide and Yao in [BWY]. In that paper they describe an algorithm that will find the closest pair (among other things) of n points chosen I.I.D. $\mathbf{U}[0, 1]^2$ in expected $\Omega(n)$ time. Their algorithm required the use of an expected $\Theta(n)$ number of distance calculations. We will show how to modify their algorithm by adding a linear time sort so that only an expected $2 \log^2 n + O(1)$ distance calculations will be performed.

There is a major qualitative difference between the sweep algorithms and the gridding algorithm: the use of the floor function. The sweep algorithms don't use the floor function while the gridding one does. From a complexity theory standpoint this is not a trivial distinction. For example, it is known that any *randomized* algorithm for solving the closest pair problem must use expected $\Omega(n \log n)$ time under the decision theory model of computation ([MT] and a reduction to element uniqueness). Yet [BWY]'s gridding algorithm uses only $O(n)$ expected time. This is because its use of the floor function means that it can no longer be analyzed under the decision tree model of computation. We conjecture that given n points chosen I.I.D. $\mathbf{U}[0, 1]^2$ any deterministic algorithm that finds their closest pair must use expected time $\Omega(n \log n)$ time under the decision tree model of computation.

The structure of this paper is as follows. In §2 we will examine the distribution of the distance between the closest pair and introduce notation. In §3 we will describe and analyze the sweep line algorithm. In §4 we will discuss how to extend this analysis to [BP]'s projection algorithm. In §5 we will describe and analyze the gridding algorithm.

Note: Most of the proofs in this extended abstract require straightforward but long calculations. In the interest of brevity we will sometimes only provide a sketch of the full proof.

§2 Distribution of the Nearest Neighbor Distance

In this section we will derive upper and lower bounds on the distribution function of the nearest neighbor distance for n points Independently, Identically Distributed (I.I.D.) from the uniform distribution from the rectangle $\mathbf{R} = [0, x] \times [0, 1]$, i.e. the rectangle with width x and height 1. We will need these bounds for

our later algorithmic analyses. As a corollary we will show that the expected nearest neighbor distance for n points chosen I.I.D. uniformly from $[0, 1]^2$ is $\Theta(1/n)$. This should be contrasted to the $\Theta(1/n)$ expected distance between an arbitrary point and its nearest neighbor.

Formally, let q_1, \dots, q_n be chosen I.I.D. from the uniform distribution over $[0, x] \times [0, 1]$. The random variable that we would like to analyze is

$$Z(q_1, \dots, q_n) = \min_{1 \leq i < j \leq n} d(q_i, q_j)$$

where $d((x, y), (x', y')) = \sqrt{(x - x')^2 + (y - y')^2}$. We will prove the following:

Lemma 1: Let $x \in [0, 1]$ and $\alpha \leq x$. Choose n points q_1, \dots, q_n , I.I.D. from the uniform distribution over $\mathbb{R} = [0, x] \times [0, 1]$. Then

$$e^{-\pi n^2 \alpha^2 / 2x} [1 + O(n\alpha^2/x + n^3 \alpha^4/x^2)] \leq \Pr(Z(q_1, \dots, q_n) \geq \alpha) \leq e^{-\pi n^2 \alpha^2 / 8x} [1 + O(n\alpha^2/x + n^3 \alpha^4/x^2)] \quad (2.1)$$

Proof: We will derive the lower bound first. The event $Z \geq \alpha$ can be rewritten as

$$Z(q_1, \dots, q_n) \geq \alpha \Leftrightarrow \forall k : q_k \notin \bigcup_{j < k} B(q_j, \alpha) \quad (2.2)$$

where $B(q_j, \alpha)$ is the Euclidean ball of radius α around q_j . We can think of this characterization as defining a point placement process. The points are placed in \mathbb{R} in order, each under the constraint that it is outside the α -balls surrounding its predecessors. Viewing the situation from this perspective lets us write

$$\begin{aligned} \Pr(Z(q_1, \dots, q_n) \geq \alpha) &= \Pr(Z(q_1, q_2) \geq \alpha) \cdot \Pr(Z(q_1, q_2, q_3) \geq \alpha \mid Z(q_1, q_2) \geq \alpha) \cdots \\ &= \prod_{k \leq n} \Pr(Z(q_1, \dots, q_k) \geq \alpha \mid Z(q_1, \dots, q_{k-1}) \geq \alpha) \\ &= \prod_{k \leq n} \Pr\left(q_k \notin \bigcup_{j < k} B(q_j, \alpha) \mid Z(q_1, \dots, q_{k-1}) \geq \alpha\right). \end{aligned} \quad (2.3)$$

Because the q_k are uniformly distributed in the rectangle \mathbb{R}

$$\Pr\left(q_k \notin \bigcup_{j < k} B(q_j, \alpha) \mid Z(q_1, \dots, q_{k-1}) \geq \alpha\right) = 1 - \frac{\mathbf{E}\left(\text{Area}\left[\bigcup_{j < k} B(q_j, \alpha) \cap \mathbb{R}\right] \mid Z(q_1, \dots, q_{k-1}) \geq \alpha\right)}{x}. \quad (2.4)$$

No matter where the q_j are placed we can always bound $\text{Area}(B(q_k, \alpha) \cap \mathbb{R}) \leq \pi\alpha^2$ so we can upper bound the expected area in (2.4) and lower bound the product in (2.3) by

$$\prod_{k < m} \left[1 - \frac{k\pi\alpha^2}{x}\right] \leq \Pr(Z(q_1, \dots, q_m) \geq \alpha). \quad (2.5)$$

Substituting the asymptotic identity

$$\left[1 - \frac{\pi k \alpha^2}{x}\right] = \exp(-\pi k \alpha^2 / x + O(k^2 \alpha^4 / x^2)) \quad (2.6)$$

into this inequality yields the lower bound of the lemma.

Now we derive the upper bound in (2.2). For this we use a second characterization of $Z \geq \alpha$:

$$Z(q_1, \dots, q_n) \geq \alpha \Leftrightarrow \forall j \neq k, B(q_j, \alpha/2) \cap B(q_k, \alpha/2) = \emptyset. \quad (2.7)$$

Since we've assumed that $\alpha \leq x$ it follows that for every point, q_j , at least one full quadrant (upper right, lower right, etc.) of $B(q_j, \alpha/2)$ is totally contained in \mathbb{R} . If $Z(q_1, \dots, q_{k-1}) \geq \alpha$ then combining this last fact with (2.7) gives

$$\text{Area} \left[\bigcup_{j < k} B(q_j, \alpha) \cap \mathbb{R} \right] \geq \text{Area} \left[\bigcup_{j < k} B(q_j, \alpha/2) \cap \mathbb{R} \right] \geq (k-1)\pi\alpha^2/4. \quad (2.8)$$

This lower bounds the expected size of the area.

$$\mathbb{E} \left(\text{Area} \left[\bigcup_{j < k} B(q_j, \alpha) \cap \mathbb{R} \right] \mid Z(q_1, \dots, q_{k-1}) \geq \alpha \right) \geq (k-1)\pi\alpha^2/4.$$

Inserting this inequality back into (2.3) and evaluating using (2.6) yields the lower bound of the lemma.

Q.E.D.

As an immediate consequence Lemma 1 easily lets us find the expected value of Z when the input points are chosen I.I.D. from the uniform distribution on the unit square (from now on we will write that the points are I.I.D. $U[0, 1]^2$).

Corollary 2: Let q_1, \dots, q_n be I.I.D. $U[0, 1]^2$ and $Z(q_1, \dots, q_n)$ be the nearest neighbor distance. Then

$$\mathbb{E}(Z(q_1, \dots, q_n)) = \Theta(1/n).$$

Proof: Apply Lemma 1 with $x = 1$ and integrate $\mathbb{E}(Z) = \int_0^{\sqrt{2}} \Pr(Z \geq \alpha) d\alpha$.

Q.E.D.

§3 The Sweep Line algorithm

§3.1 Description of the algorithm

The algorithm we consider here is the sweep line algorithm for computing the closest distance among points in the plane. Before presenting our analysis we give a quick description of the algorithm. For a more complete description as well as a correctness proof and implementation details see [HNS].

Suppose we are given n points $p_i = (x_i, y_i)$, $i = 1 \dots n$, in the plane. Sort them in increasing x order (arbitrarily breaking ties). We use the standard notation for order statistics. The i 'th sorted point is denoted by $p_{(i)} = (x_{(i)}, y_{(i)})$ where

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n-1)} \leq x_{(n)}.$$

The x -coordinate of p_i is the i -th smallest x -value. The y -coordinate of $p_{(i)}$ is the y coordinate that was associated with $x_{(i)}$. For example, if the original points were $p_1 = (5, 3)$, $p_2 = (3, 7)$, and $p_3 = (18, 6)$, then the relabeled points are $p_{(1)} = (3, 7)$, $p_{(2)} = (5, 3)$, and $p_{(3)} = (18, 6)$.

Let δ_i be the distance between between the closest pair among the first i points, i.e.

$$\delta_i = \min_{1 \leq k < l \leq i} d(p_{(k)}, p_{(l)}) = Z(p_{(1)}, \dots, p_{(i)}).$$

The closest pair distance problem is then seen to be the problem of determining δ_n . It follows from the definition of δ_i that

$$\delta_{i+1} = \min \left(\delta_i, \min_{1 \leq k \leq i} d(p_{(i+1)}, p_{(k)}) \right). \quad (3.1)$$

Updating δ_i to δ_{i+1} only requires calculating the minimum among $d(p_{(i+1)}, p_{(k)})$. Furthermore, if there is some $k \leq i$ such that $d(p_{(i+1)}, p_{(k)}) < \delta_i$ then for that k , $x_{(i+1)} - x_{(k)} < \delta_i$. Thus we only have to check the distance between $p_{(i+1)}$ and those points whose x -coordinates are in the interval $(x_{(i+1)} - \delta_i, x_{(i+1)})$, a region which we will call the δ_i -interval (see Figure 1(a)).

There is a geometric fact (proven in [HNS]) that is essential to understanding how to insure the algorithm $O(n \log n)$ worst case behavior. Suppose we have sorted the points in the δ_i -interval by their y -coordinates. Furthermore suppose that

$$\exists k \leq i \text{ such that } d(p_{(i+1)}, p_{(k)}) < \delta_i.$$

Then $p_{(k)}$ has to be one of the four points immediately above or one of the four points immediately below the point $p_{(i+1)}$ in the δ_i -interval. Thus, if the points in the δ_i -interval are stored in a balanced tree sorted by y -coordinate, we can calculate δ_{i+1} in $O(\log n)$ time by retrieving these eight points. All that remains is to describe how to update the δ_i -interval, i.e. update the y -sorted balanced tree so that it contains the points in the δ_{i+1} -interval. This is not difficult. First we insert $p_{(i+1)}$ into the balanced tree in $O(\log n)$ time. Then, using the fact that we already have the points sorted in increasing x -order (and have stored somewhere the index of the leftmost point in the δ_i -interval) we scan rightward deleting points from the tree until we enter the δ_{i+1} -interval (see Figure 1(b)). Since each point is deleted at most once during the execution of the algorithm the total time for deletion is $O(n \log n)$. Furthermore, since all other operations are $O(\log n)$ time per step we find that the entire second stage of the algorithm can be implemented in $O(n \log n)$ time.

The main purpose of this section is to show that, probabilistically, the performance of the second stage of this algorithm – scanning through the sorted points while updating δ_i – will be much better than the worst case analysis suggests. In fact, we will show that it is possible to dispense totally with complicated balanced tree data structures and *still* get very good behavior: if we compare $p_{(i+1)}$ to *all* of the points in the δ_i -interval the algorithm will still run very fast on the average. We give a short but full listing of code for this version of the algorithm in Figure 2. It is not hard to see that this simplified version of the algorithm will require $\Omega(n^2)$ time in the worst case when run on n , x -sorted points, e.g. if the points are all on the same vertical line. We will show that, on average, it will be much faster than that.

Mathematically, given n points p_1, \dots, p_n , chosen I.I.D. $\mathbf{U}[0, 1]^2$, we set

$$N_i = \text{number of points in the } \delta_i\text{-interval}.$$

The amount of work needed to update δ_i to δ_{i+1} is N_i so the total amount of work performed by the algorithm is $\sum N_i$. To prove that the algorithm runs in expected linear time we must show that

$$\mathbf{E}\left(\sum_{i=1}^n N_i\right) = \sum_{i=1}^n \mathbf{E}(N_i) = O(n). \quad (3.2)$$

This, along with a matching lower bound, will be an immediate consequence of the following theorem.

Theorem 3: Choose n points I.I.D. $\mathbf{U}[0, 1]^2$. Let N_i be the number of points in the δ_i -interval. Then

$$\mathbf{E}(N_i) = \Theta\left(\sqrt{\frac{n}{i}}\right), \quad i > n^{1/3} \lg n. \quad (3.3)$$

More specifically

$$\sqrt{\frac{2}{3\pi}} \sqrt{\frac{n}{i}} + O\left(\sqrt{\frac{n}{i^3}}\right) \leq \mathbf{E}(N_i) \leq \sqrt{\frac{16}{\pi}} \sqrt{\frac{n}{i}} + 4 + O\left(\sqrt{\frac{n}{i^3}}\right).$$

Combining this theorem with the fact that $\sum_{i \leq n} \sqrt{\frac{n}{i}} \sim 2n$ and $\sum_{i \leq n} \sqrt{\frac{n}{i^3}} = O(\sqrt{n})$ proves what we want:

n	$\sum_i N_i$	$\frac{1}{n} \sum_i N_i$
1000	1287.9 ± 437.6	1.287900 ± 0.4376
5000	6665.0 ± 2179.6	1.332996 ± 0.4359
10000	14254.9 ± 4981.7	1.425485 ± 0.4982
20000	27536.2 ± 9550.5	1.376812 ± 0.4775
30000	38866.0 ± 12462.7	1.295532 ± 0.4154
40000	56985.9 ± 17986.8	1.424649 ± 0.4497
50000	67907.2 ± 19342.2	1.358143 ± 0.3868

Table 1. For each value of n we generated 100 sets of n random points I.I.D. $U[0, 1]^2$. Each set was sorted and $\sum N_i$ was calculated. The value in the second column is the average value of $\sum N_i$ over all 100 point sets. The value in the third column is the value in the second normalized by division by n .

Corollary 4: The total amount of expected work, $\sum E(N_i)$, performed by the second stage of the sweep-line algorithm is $\Theta(n)$. More specifically

$$0.6514 \dots n + O(n^{2/3} \log^2 n) \leq E(N_i) \leq 6.2567 \dots n + O(n^{2/3} \log^2 n).$$

The constant on the right hand side is small so the algorithm should run relatively fast. How well does this work in practice? In Table 1 we provide statistics on $E(\sum N_i)$ collected by running the algorithm on random point sets. For each value of n we ran 100 random trials and averaged the results. The statistics seem to show that the running time of the algorithm is much closer to the lower bound than to the upper.

In §3.2 we will introduce some probabilistic notation and facts and in §3.3 we will prove the Theorem 4.

§3.2 Probabilistic Assumptions and Notation

We will assume that we are given n points $p_i = (x_i, y_i)$, $i = 1, \dots, n$, chosen I.I.D. $U[0, 1]^2$. The p_i -s are then sorted by x -coordinate² and renamed $p_{(i)} = (x_{(i)}, y_{(i)})$ where

$$x_{(1)} < x_{(2)} < \dots < x_{(n-1)} < x_{(n)}. \quad (3.4)$$

We define the random variables $X_{(i)} = x_{(i)}$ and let $\delta_i = \min_{1 \leq k < l \leq i} d(p_{(k)}, p_{(l)})$ as before. Given a real value X we say that the α -interval is the strip of width α to the left of the line $x = X$. The number of points in an α -interval will be denoted by

$$\begin{aligned} L(X, \alpha) &= \text{number of } p_i \text{ in the } \alpha \text{ strip to the left of } X \\ &= |\{p_j \mid X - \alpha \leq x_j \leq X\}|. \end{aligned}$$

In this notation the number of points in the δ_i -interval is

$$N_i = L(X_{(i+1)}, \delta_i).$$

Directly from the definitions we see that $N_i \leq i$. Our goal is to analyze $E(N_i) = E(L(X_{(i+1)}, \delta_i))$. If α is a constant then it is not hard to see that $L(X_{(i+1)}, \alpha) \mid X_{(i+1)} = x$ is a binomial³ $\mathcal{B}(i, \alpha/x)$ random variable

² A technical quibble: this definition is not valid if there are two points with the same x -coordinate. We don't have to worry about this happening. If the points are chosen I.I.D. $U[0, 1]^2$ then two points sharing the same x -coordinate is a zero probability event which we can safely ignore.

³ A random variable Y is binomial $\mathcal{B}(m, p)$ if $\Pr(Y = k) = \binom{m}{k} p^k (1-p)^{m-k}$ for $0 \leq k \leq m$.

and therefore easy to analyze. The random variable $L(\mathbf{X}_{(i+1)}, \delta_i) | \mathbf{X}_{(i+1)} = \mathbf{x}, \delta_i = \alpha$ is not binomial though because conditioning on $\delta_i = \alpha$ totally changes the distribution of the points. It restricts all pairs of points from being less than a distance α from each other and also requires that there is some pair that is *exactly* a distance α apart. A large portion of our analysis will be devoted to dealing with this conditionality.

First we examine the $\mathbf{X}_{(i)}$. We take advantage of the fact that the $\mathbf{X}_{(i)}$ have been extensively studied [De] because they are the *order statistics* of random variables. More specifically, since the $p_i = (x_i, y_i)$ are I.I.D. uniformly in the unit square the x_i have to be I.I.D. uniformly distributed in $[0, 1]$. The sorted x -coordinates, $\mathbf{X}_{(i)}$, are the order statistics of n uniformly distributed random variables on $[0, 1]$. The probability density $f_{(i)}(\mathbf{x})$, of $\mathbf{X}_{(i)}$ is well known (see [De], p. 17, for a derivation) and is given by the Beta distribution:

$$f_{(i)}(\mathbf{x}) = \frac{n!}{(i-1)!(n-i)!} x^{i-1}(1-x)^{n-i} \quad \mathbf{x} \in [0, 1]. \quad (3.5)$$

Integration gives us the expectation and the variance:

$$\mathbf{E}(\mathbf{X}_{(i)}) = \frac{i}{n+1} \quad \text{Var}(\mathbf{X}_{(i)}) = \frac{i(n+1-i)}{(n+2)(n+1)^2}. \quad (3.6)$$

We can now use Chebyshev's inequality to restrict the range of $\mathbf{X}_{(i+1)}$:

Corollary 6:

$$\Pr \left(\mathbf{X}_{(i+1)} \notin \left[\frac{1}{2} \frac{(i+1)}{(n+1)}, \frac{3}{2} \frac{(i+1)}{(n+1)} \right] \right) \leq \frac{4}{i+1}. \quad (3.7)$$

Proof. Chebyshev's inequality gives

$$\begin{aligned} \Pr \left(\mathbf{X}_{(i+1)} \notin \left[\frac{1}{2} \frac{(i+1)}{(n+1)}, \frac{3}{2} \frac{(i+1)}{(n+1)} \right] \right) &= \Pr \left(|\mathbf{X}_{(i+1)} - \mathbf{E}(\mathbf{X}_{(i+1)})| > \frac{(i+1)}{2(n+1)} \right) \\ &\leq \text{Var}(\mathbf{X}_{(i+1)}) \left(\frac{2(n+1)}{(i+1)} \right)^2 \leq \frac{4}{(i+1)}. \end{aligned}$$

§3.3 Proof of Theorem 4

Proof: We will prove that

$$\sqrt{\frac{2}{3\pi}} \sqrt{\frac{n}{i}} + O \left(\sqrt{\frac{n}{i^3}} \right) \leq \mathbf{E}(\mathbf{N}_i) \leq \sqrt{\frac{16}{\pi}} \sqrt{\frac{n}{i}} + 4 + O \left(\sqrt{\frac{n}{i^3}} \right), \quad i > n^{1/3} \lg n.$$

We start by splitting $\mathbf{E}(\mathbf{N}_i)$ into the sum of two integrals:

$$\mathbf{E}(\mathbf{N}_i) = \int_{\mathbf{x} \in D_i} \mathbf{E}(\mathbf{N}_i | \mathbf{X}_{(i+1)} = \mathbf{x}) f_{(i+1)}(\mathbf{x}) d\mathbf{x} + \int_{\mathbf{x} \notin D_i} \mathbf{E}(\mathbf{N}_i | \mathbf{X}_{(i+1)} = \mathbf{x}) f_{(i+1)}(\mathbf{x}) d\mathbf{x}. \quad (3.8)$$

From Corollary 6 and the fact that $\mathbf{N}_i \leq i$ we immediately see that

$$\int_{\mathbf{x} \notin D_i} \mathbf{E}(\mathbf{N}_i | \mathbf{X}_{(i+1)} = \mathbf{x}) f_{(i+1)}(\mathbf{x}) d\mathbf{x} \leq i \cdot \Pr(\mathbf{X}_{(i+1)} \notin D_i) < 4.$$

We can therefore rewrite (3.8) as

$$\mathbf{E}(\mathbf{N}_i) = \int_{\mathbf{x} \in D_i} \mathbf{E}(\mathbf{N}_i | \mathbf{X}_{(i+1)} = \mathbf{x}) f_{(i+1)}(\mathbf{x}) d\mathbf{x} + O(1). \quad (3.9)$$

We will now show that, for $i > n^{1/3} \lg n$ and any $x \in D_i$,

$$\mathbf{E}(\mathbf{N}_i | \mathbf{X}_{(i+1)} = x) = \Theta\left(\sqrt{\frac{n}{i}}\right).$$

The proof of the theorem will follow by the insertion of this last statement into (3.9).

Henceforth we will assume that $\mathbf{X}_{(i+1)} = x$ for some constant $x \in D_i$. As before, \mathbf{R} is the rectangular region to the left of the line $x = \mathbf{X}_{(i+1)} : \mathbf{R} = [0, x] \times [0, 1]$. Conditioning on $\mathbf{X}_{(i+1)} = x$ we know that there are exactly i points in \mathbf{R} . We can randomly relabel the points so that $q_1 = (x_1, y_1), \dots, q_i = (x_i, y_i)$ are the ones in \mathbf{R} . The crucial fact upon which our analysis will be based is that *after the points are relabeled the q_1, \dots, q_i are I.I.D. uniformly in \mathbf{R} .*

We define new random variables

$$M_j = \begin{cases} 1 & \text{if } x - \delta_i < x_j \leq x \\ 0 & \text{otherwise.} \end{cases}$$

These are the indicator functions for the events q_j is in the δ_i -interval, so $\mathbf{N}_i = \sum_{j \leq i} M_j$. Because the q_j -s are I.I.D., linearity of the expectation operator yields⁴

$$\mathbf{E}(\mathbf{N}_i) = i\mathbf{E}(M_1). \quad (3.10)$$

Since q_1 is (uniformly) distributed in \mathbf{R} we can find $\mathbf{E}(M_1)$ by integrating over the conditional expectation given that q_1 is fixed, i.e.

$$\mathbf{E}(M_1) = \frac{1}{x} \int_{y_1=0}^1 \int_{x_1=0}^x \mathbf{E}(M_1 | q_1 = (x_1, y_1)) dx_1 dy_1. \quad (3.11)$$

At first glance it might seem that all we've done is exchange one complicated integral, (3.8), for another, (3.11). In the next few lines we give a few simple transformations that enable (3.11) to be rewritten in a more pliable form.

$$\begin{aligned} \mathbf{E}(M_1) &= \frac{1}{x} \int_{y_1=0}^1 \int_{x_1=0}^x \Pr(M_1 = 1 | q_1 = (x_1, y_1)) dx_1 dy_1 \\ &= \frac{1}{x} \int_{y_1=0}^1 \int_{x_1=0}^x \Pr(x - \delta_i < x_1 \leq x | q_1 = (x_1, y_1)) dx_1 dy_1 \\ &= \frac{1}{x} \int_{y_1=0}^1 \int_{x_1=0}^x \Pr(\delta_i \geq x_1 | q_1 = (x - x_1, y_1)) dx_1 dy_1. \end{aligned} \quad (3.12)$$

The first equality follows from M_1 being an indicator function, the second from the definition of M_1 , and the third from the transformation $x_1 \rightarrow (x - x_1)$.

Our plan is to bound $\mathbf{E}(M_1)$ by calculating upper and lower bounds on $\Pr(\delta_i \geq x_1 | q_1 = (x - x_1, y_1))$ and integrating over them in the last line of (3.12). Since $\delta_i = Z(q_1, \dots, q_i)$ we really want bounds on

$$\Pr(Z(q_1, \dots, q_i) \geq x_1 | q_1 = (x - x_1, y_1)).$$

If this random variable wasn't conditioned on the event $q_1 = (x - x_1, y_1)$ Lemma 1 would give us good bounds. But, as pointed out in the remark after Lemma 1, the proof of the lemma never used the fact that

⁴ As we said before, all random variables and expectations in this part are conditioned on $\mathbf{X}_{(i+1)} = x$. For the sake of clarity we don't write this explicitly. For example, (3.10) should be read as $\mathbf{E}(\mathbf{N}_i | \mathbf{X}_{(i+1)} = x) = i\mathbf{E}(M_1 | \mathbf{X}_{(i+1)} = x)$.

q_1 was chosen from the uniform distribution and it would still be true even if q_1 was deterministically chosen. Therefore we can apply lemma 1 to find that

$$e^{-\pi n^2 \alpha^2 / 2x} [1 + O(n\alpha^2/x + n^3 \alpha^4/x^2)] \leq \Pr(\delta_i \geq x_1 | q_1 = (x - x_1, y_1)) \leq e^{-\pi n^2 \alpha^2 / 8x} [1 + O(n\alpha^2/x + n^3 \alpha^4/x^2)] \quad (3.13)$$

whenever $x_1 \leq x$ (this is the $\alpha \leq x$ condition). Recall that we are assuming that x is fixed, $x \in D_i$. We can therefore write $x = ci/n$ where $1/2 \leq c \leq 3/2$. Recall too the theorem's assumption that $i > n^{1/3} \log n$. Combining these two assumptions we see that when $x_1 \leq \sqrt{\lg n / in}$ then $x_1 \leq x$ and we can use (3.13). We will now show how to use this fact to upper bound $\mathbf{E}(N_i) = i\mathbf{E}(M_1)$. First we substitute $x = ci/n$ and integrate

$$\begin{aligned} \int_0^1 \Pr(\delta_i \geq x_1 | q_1 = (x - x_1, y_1)) dx_1 &= \int_0^{\sqrt{\lg n / in}} \Pr(\delta_i \geq x_1) dx_1 + \int_{\sqrt{\lg n / in}}^1 \Pr(\delta_i \geq x_1) dx_1 \\ &\leq \int_0^{\sqrt{\lg n / in}} e^{-\pi n^2 x_1^2 / 2c} [1 + O(n x_1^2 + i n^2 x_1^4)] dx_1 + n^{-\Omega(\lg n)} \\ &= \sqrt{\frac{8c}{\pi i n}} \left[1 + O\left(\frac{1}{i}\right) \right]. \end{aligned}$$

Inserting this back into (3.12) gives

$$\mathbf{E}(M_1) \leq \frac{n}{ic} \int_0^1 \sqrt{\frac{8c}{\pi i n}} \left[1 + O\left(\frac{1}{i}\right) \right] dy_1 \leq \sqrt{\frac{8}{\pi c}} \sqrt{\frac{n}{i^3}} + O\left(\frac{\sqrt{n}}{i^5}\right)$$

and therefore

$$\mathbf{E}(N_i) \leq i\mathbf{E}(M_1) \leq \sqrt{\frac{8}{\pi c}} \sqrt{\frac{n}{i}} + O\left(\sqrt{\frac{n}{i^3}}\right).$$

All probabilistic statements in this section were conditioned on $\mathbf{X}_{(i+1)} = x \in D_i$. For these x , we know that $1/2 \leq c \leq 3/2$. Therefore we have actually proven that, for $x \in D_i$,

$$\mathbf{E}(N_i | \mathbf{X}_{(i+1)} = x) \leq \sqrt{\frac{16}{\pi}} \sqrt{\frac{n}{i}} + O\left(\sqrt{\frac{n}{i^3}}\right)$$

Reinserting this into (3.9) immediately proves the upper bound of the theorem. The lower bound is proved similarly. The only difference is that we use the left hand side of (3.13) and not the right hand one.

Q.E.D.

Note: Until now we have assumed that the distance function $d(p, q)$, is the Euclidean L_2 metric. Working through the proof of Theorem 3 one finds that it remains valid when $d()$ is replaced by the L_p ($p \geq 1$) metric

$$d_p(p_i, p_j) = (|x_i - x_j|^p + |y_i - y_j|^p)^{1/p}$$

or the L_∞ metric

$$d_p(p_i, p_j) = \max(|x_i - x_j|, |y_i - y_j|).$$

The only thing that changes in the proof is that π is replaced by ω_p , the area of the unit ball under the L_p metric:

$$\omega_p = \text{Area}\{q | d_p((0, 0), q) \leq 1\}.$$

§4 The Projection Algorithm

In [BP] Bentley and Papadimitriou describe an algorithm for solving the “post-office” problem. Given n points, p_1, \dots, p_n , the “post-office” problem is to preprocess them so that, given a new point, q , we can always find the point p_i which is closest to q , i.e.

$$\forall j, \quad d(q, p_i) \leq d(q, p_j).$$

The preprocessing step proposed by [BP] sorts the p_j 's by increasing x -coordinate. To find the closest point to a new point, q , they compare q to the points which are nearest to q by projection on the x axis. This is a two step procedure. They first use a binary search to find p_i such that $p_i.x < q.x \leq p_{i+1}.x$. They then, using their own phrase, “search out,” q 's nearest neighbor by always comparing q to the point whose x -coordinate is closest to q and to which q has not yet been compared. While doing this they keep track of d , the closest distance between q and all of the points it has been compared to so far. They terminate the search when q has been compared to all points p_j whose x -coordinates are within the range $[q.x - d, q.x + d]$. We call this algorithm NNX : NNX(q) returns the distance between q and its nearest neighbor.

The problem with this algorithm is that it can take $\Omega(n)$ time, e.g when all of the p_i lie on the same vertical line. In an effort to avoid this problem [BP] notes that there is nothing sacred about projection on the x -axis; we can just as easily write the algorithm so that it projects on the y -axis by “searching-out” the points whose y -coordinates are closest to q 's y -coordinate. This algorithm, which we call NNY, also requires $\Omega(n)$ worst case time. We can also interleave projection on the x -axis with projection on the y -axis. That is, we can run the two versions, NNX and NNY, in lockstep: at each stage comparing q to the point with closest x -coordinate, and then to the point with closest y -coordinate. This last version of the algorithm, which we call NN, has an interesting combinatorial worst case analysis. [BP] prove that if closest is defined by the L_∞ metric then if we calculate ⁵ $NN(p_1), NN(p_2), \dots, NN(p_{n-1}),$ and $NN(p_n)$ the total time taken by the algorithm for these n calls will be $O(n^{3/2})$ as compared to an $\sqrt{n^2}$ worst case bound if we had run just the x -projection version $NNX(p_1), NNX(p_2), \dots, NNX(p_{n-1}),$ and $NNX(p_n)$.

We now modify NN so that it finds the distance between the closest pair. Obviously one way of doing this would be to calculate all of the values $NN(p_1), \dots, NN(p_n)$ and report the minimum value. A more sophisticated technique uses information already calculated. First sort the points so that they are in x sorted order as before relabeling them $p_{(1)}, p_{(2)}, \dots, p_{(n)}$:

$$p_{(1)}.x \leq p_{(2)}.x \leq \dots \leq p_{(n)}.x.$$

We examine the $p_{(i)}$ in this x -sorted order. Suppose we have already calculated ν_i , the minimum distance between the first i points and all of the others:

$$\nu_i = \min(\text{NN}(p_{(1)}), \text{NN}(p_{(2)}), \dots, \text{NN}(p_{(i)})) = \min_{\substack{1 \leq k \leq i \\ 1 \leq j \leq n \\ k \neq j}} d(p_{(k)}, p_{(j)}),$$

We will only be interested in calculating $\text{NN}(p_{(i+1)})$ if $\text{NN}(p_{(i+1)}) < \nu_i$. Thus, if we can convince ourselves that $\text{NN}(p_{(i+1)}) \geq \nu_i$, we can truncate the search performed by $\text{NN}(p_{(i+1)})$. The truncation is implemented by stopping the execution of $\text{NN}(p_{(i+1)})$ if there are no points in $\{(x, y) \mid p_{(i+1)}.x - \nu_i < x < p_{(i+1)}.x + \nu_i\}$ or if there are no points in $\{(x, y) \mid p_{(i+1)}.y - \nu_i < y < p_{(i+1)}.y + \nu_i\}$. We name this interleaved closest pair algorithm CP. [BP]'s result implies that CP takes $O(n^{3/2})$ time in the worst case when *closest* is defined by the L_∞ metric. They also point out that, while, in general, they can't extend this $O(n^{3/2})$ bound to the L_2 metric, in practice, the use of finite precision arithmetic does let them extend it.

⁵ By calculating $\text{NN}(p_j)$ we mean finding the point $p_i \neq p_j$, such that p_i is closest to p_j . This requires us to slightly modify our code so that it p_j is not reported as its own nearest neighbor.

To analyze this algorithm we need to define notation analogous to $L(X, \alpha)$ which, instead of counting the number of points in the α -strip *left* of X counts the number to the *right* of X . We set

$$\begin{aligned} R(X, \alpha) &= \text{number of } p_i \text{ in the } \alpha \text{ strip to the right of } X \\ &= |\{p_j \mid X \leq x_j \leq X + \alpha\}|. \end{aligned}$$

During the truncated execution of $NN(p_{(i)})$ algorithm CP will examine at most

$$2 [L(\mathbf{X}_{(i+1)}, \nu_i) + R(\mathbf{X}_{(i+1)}, \nu_i)] = 2 |\{j : j \neq i, |p_{(j)}.x - p_{(i)}.x| \leq \nu_i\}|$$

points. This is because at most $L(\mathbf{X}_{(i+1)}, \nu_i) + R(\mathbf{X}_{(i+1)}, \nu_i)$ points are examined during the x -projection part of $NN(p_{(i)})$ because their x -coordinates are within ν_i of $p_{(i)}$'s x -coordinate. Furthermore at most one point will be examined in the y -projection part of $NN(p_{(i)})$ for each point examined in the x -projection part.

Using techniques very similar to those used to prove Theorem 3 we can prove

Theorem 7: Choose n points I.I.D. $U[0, 1]^2$. Then

$$\mathbf{E}(L(\mathbf{X}_{(i+1)}, \nu_i) + R(\mathbf{X}_{(i+1)}, \nu_i)) = O\left(\sqrt{\frac{n}{i}}\right), \quad i > n^{1/3} \lg n.$$

Furthermore this result remains valid when the Euclidean distance function is replaced by the L_p ($p \geq 1$) metric

$$d_p(p_i, p_j) = (|x_i - x_j|^p + |y_i - y_j|^p)^{1/p}$$

or the L_∞ metric

$$d_p(p_i, p_j) = \max(|x_i - x_j|, |y_i - y_j|).$$

Algorithm CP starts with a deterministic $O(n \lg n)$ preprocessing stage that sorts the points by x coordinate and also by y -coordinate.

It then performs an expected $\sum_{i \leq n} \mathbf{E} [L(\mathbf{X}_{(i+1)}, \nu_i) + R(\mathbf{X}_{(i+1)}, \nu_i)]$ amount of work in the truncated calls to $NN()$. The theorem tells us that all of the calls to $NN()$ require at most $O(n)$ expected time. Thus this algorithm has the same time bounds as the sweep-line algorithm we analyzed in the previous section. This new algorithm has the advantage that we know – for the L_∞ distance function – that its worst case running time is $O(n^{3/2})$ while the worst case running time for the sweep line algorithm was $O(n^2)$. It has the disadvantage of requiring more complicated code than the sweep-line algorithm and executing two sorts as opposed to the sweep-line's one sort.

In this section we only provided a sketch of the algorithm. For a full description and a complete analysis see [Go].

§5 A “Gridding” algorithm

In this section we describe another algorithm that efficiently finds the closest pair among a set of n points p_1, \dots, p_n , that are chosen I.I.D. $U[0, 1]^2$. Unlike the algorithms in the previous two sections this new algorithm is not robust. It only works well when the input points are actually drawn from the unit square. Counterbalancing this flaw is that fact that it is a truly linear expected time algorithm. Furthermore it does very few costly distance calculations. More specifically the algorithm has a deterministic overhead of $\Theta(n)$ time (floor functions and radix sorting). The remainder of the work it performs, including all distance computations will be shown to take only expected $\Theta(\log^2 n)$ time.

This algorithm is a natural extension of the gridding algorithm found in [BWY]. There they present an expected $O(n)$ time algorithm for finding the Vornoi diagram (and therefore the closest pair) of n points chosen I.I.D. $\mathbf{U}[0, 1]^2$. Their basic idea was to partition the unit square into n smaller squares, or cells, each with dimension $n^{-1/2} \times n^{-1/2}$. They then showed that to construct the Vornoi region associated with each point it is only necessary to examine the distances between the point and an expected constant number of points in the same and “neighboring” cells. This lets them construct the entire Vornoi diagram in expected $O(n)$ time using expected $O(n)$ distance computations.

In this paper we are examining closest pair algorithms so we don't need to calculate the entire Vornoi structure. We will show that it is possible to modify their algorithm in such a way that only an expected $O(\log^2 n)$ number of distance calculations are done. Given n points p_1, \dots, p_n , I.I.D. $\mathbf{U}[0, 1]^2$ we will set $M = \lfloor n/\log n \rfloor$. The algorithm essentially divides the unit square into $\Theta(M^2)$ smaller squares and only calculates the distance between points if they are in the same smaller square.

Now the specifics: First note that Lemma 1 tells us that

$$\Pr(Z(p_1, \dots, p_n) \geq 1/2M) = n^{-\Omega(\log n)}. \quad (5.1)$$

We will demonstrate an efficient procedure that will find all pairs of points (p_i, p_j) such that

$$\max(|p_i.x - p_j.x|, |p_i.y - p_j.y|) < 1/2M. \quad (5.2)$$

Our algorithm will execute this procedure. If it finds that there are no pairs that satisfy (5.2) then it finds the closest pair by examining all $\binom{n}{2}$ pairs of points. (5.1) tells us that this will happen with such a small probability that this second part of the algorithm only takes $n^{-\Omega(\log n)}$ expected time.

A pair of points p_i, p_j will satisfy (5.2) only if one of the following conditions is satisfied:

- (a) $\lfloor M * p_i.x \rfloor = \lfloor M * p_j.x \rfloor$ and $\lfloor M * p_i.y \rfloor = \lfloor M * p_j.y \rfloor$:
- (b) $\lfloor M * p_i.x + 1/2 \rfloor = \lfloor M * p_j.x + 1/2 \rfloor$ and $\lfloor M * p_i.y \rfloor = \lfloor M * p_j.y \rfloor$:
- (c) $\lfloor M * p_i.x \rfloor = \lfloor M * p_j.x \rfloor$ and $\lfloor M * p_i.y + 1/2 \rfloor = \lfloor M * p_j.y + 1/2 \rfloor$:
- (d) $\lfloor M * p_i.x + 1/2 \rfloor = \lfloor M * p_j.x + 1/2 \rfloor$ and $\lfloor M * p_i.y + 1/2 \rfloor = \lfloor M * p_j.y + 1/2 \rfloor$.

We will show how to find all pairs of points that satisfy condition (a). Associate with each point p_i the ordered pair of integers $g(p) = (\lfloor M * p_i.x \rfloor, \lfloor M * p_i.y \rfloor)$. Now sort the points p_1, \dots, p_n lexicographically⁶ using $g(p_i)$ as the key for point p_i . It is known [AHO] that n pairs of integers can be lexicographically sorted in $O(n)$ time if the components of the pairs are all less than n . This is certainly the case here because $\lfloor M * p_i.x \rfloor < n$ and $\lfloor M * p_i.y \rfloor < n$. Now notice that $g(p_i) = g(p_j)$ if and only if the pair p_i, p_j fulfills condition (a). Therefore finding all pairs of points that fulfill condition (a) is equivalent to finding all pairs of repeated keys in a sorted list. Let

$$Y_a = |\{(i, j) \mid i < j, \lfloor M * p_i.y \rfloor = \lfloor M * p_j.y \rfloor \text{ and } \lfloor M * p_i.x \rfloor = \lfloor M * p_j.x \rfloor\}|.$$

There will be Y_a such repeated pairs and we can find them in time $O(n) + Y_a$ by walking through the list. A pair of points will satisfy condition (a) with probability M^{-2} and there are $\binom{n}{2}$ such pairs so

$$\mathbb{E}(Y_a) = \binom{n}{2} M^{-2} = \frac{\log^2 n}{2} + O(1).$$

Therefore we can find all pairs of points that fulfill condition (a) in expected time $O(n)$ (without performing any distance calculations) and the total expected number of such pairs will be $\log^2 n/2 + O(1)$.

⁶ The pair (X, Y) is lexicographically less than the pair (X', Y') if $X < X'$, or $X = X'$ and $Y < Y'$.

Similarly we can find all pairs of points fulfilling each of conditions (b), (c) and (d) in expected $O(n)$ time without performing any distance calculations and the expected number of points that fulfill each of these conditions is $E(Y) = \log^2 n/2 + O(1)$.

Thus we have managed to find a superset of all pairs of points (some pairs will be found two or more times depending on how many conditions they fulfill) that satisfy (5.2) and the expected number of pairs reported will be $2 \log^2 n + O(1)$.

The total amount of time that it took to find these pairs was $O(n)$. While finding them no distance calculations were performed because the distance calculations are performed after the points are reported. Thus the entire algorithm takes expected $O(n)$ time and performs only an expected $O(\log^2 n)$ number of distance calculations.

Acknowledgements: The author would like to thank Jon Bentley for the many conversations and suggestions that he provided.

REFERENCES:

- [AHU] A. Aho, J. Hopcroft and J. D. Ullman, *The Design and Analysis of Computer Algorithms*, Addison-Wesley, Reading, Mass. (1974)
- [BP] J. L. Bentley and C. H. Papadimitriou, "A Worst-Case Analysis of Nearest-Neighbor Searching by Projection," *7th Int. Conf. on Automata, Languages and Programming*, (August 1980) 470-482.
- [BWY] J.L. Bentley, B.W. Weide, and A.C. Yao, "Optimal Expected-Time Algorithms for Closest Point Problems," *ACM Trans. on Mathematical Software*, 6(4) (Dec. 1980) 563-580.
- [De] Luc Devroye, *Non-Uniform Random Variate Generation*, Springer-Verlag, New York. (1986).
- [FH] S. Fortune and J.E. Hopcroft, "A Note on Rabin's Nearest Neighbor Algorithm," *Information Processing Letters*, 8(1) (1979) 20-23.
- [Go] M. J. Golin *Probabilistic Analysis of Geometric Algorithms (Thesis)*, Princeton University Technical Report CS-TR-266-90. June 1990.
- [HNS] Klaus Hinrichs, Jurg Nievergelt, and Peter Schorn, "Plane-Sweep Solves the Closest Pair Problem Elegantly," *Information Processing Letters*, 26 (January 11, 1988) 255-261.
- [MT] U. Manber and M. Tompa "The Complexity of Problems on Probabilistic, Nondeterministic, and Alternating Decision Trees," *JACM*, 32(3) (July 1985) 720-732.
- [Ra] M.O. Rabin, "Probabilistic Algorithms," *Algorithms and Complexity: New Directions and Recent Results (J.F. Traub ed.)*, (1976) 21-39.
- [Sh] Michael Ian Shamos, "Computational Geometry," *Thesis (Yale)*, (1978).
- [SH] Michael Ian Shamos and Dan Hoey, "Closest Point Problems," *16th Annual Symposium on Foundations of Computer Science (FOCS)*, (1975) 151-161.
- [We] Bruce W. Weide, "Statistical Methods in Algorithm Design and Analysis," *Thesis (Carnegie-Mellon University)*. CMU-CS-78-142, (August 1978).

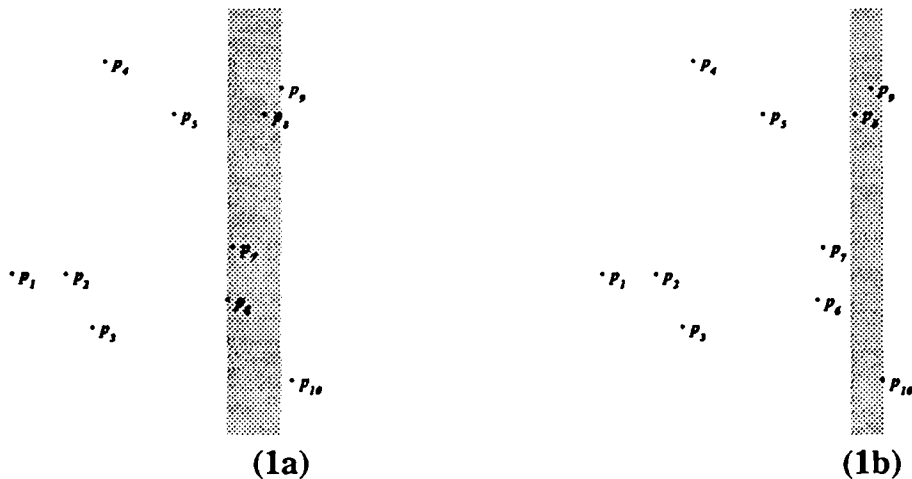


Figure 1. The minimum pair among all points preceding p_9 is (p_1, p_2) so $\delta_8 = \|p_1 - p_2\|$. In figure 1(a) the shaded region has width δ_8 and is what we called the δ_8 -interval. While checking all the points in the δ_8 -interval we find that $\|p_9 - p_8\| < \delta_8$ and thus $\delta_9 = \|p_9 - p_8\|$. We now find the δ_9 interval by scanning rightward from p_6 until we hit the first point (p_8) whose x -coordinate is within δ_9 of $p_{10}.x$. This gives us figure 1(b).

Figure 2. A Pascal listing of code that implements the linear scan for the closest pair. It is assumed that the points are sorted by x -coordinate in the array $p[]$ and that there is some distance function $dist$. $left$ and i point to the leftmost and rightmost points in the δ_i -interval. The **while** loop updates the points in the δ_i -interval. The inner **for** loop updates the current value of δ_i by comparing $p[i+1]$ to all of the points in the δ_i -interval.

```

left := 1; delta := dist(p[1], p[2]);
for i := 1 to (n-1) do
  begin
    while (p[i+1].x - p[left].x) < delta
      left := left + 1;
    for j := left to i do
      if dist(p[j], p[i+1]) < delta
        delta := dist(p[j], p[i+1]);
  end;

```

ISSN 0249 - 6399