# Upper and Lower Bounds on the Quality of the PCA Bounding Boxes

Darko Dimitrov, Christian Knauer, Klaus Kriegel, Günter Rote

Freie Universität Berlin
Institute of Computer Science
Takustrasse 9, D-14195 Berlin, Germany

darko/knauer/kriegel/rote@inf.fu-berlin.de

## ABSTRACT

Principle component analysis (PCA) is commonly used to compute a bounding box of a point set in $\mathbb{R}^d$. The popularity of this heuristic lies in its speed, easy implementation and in the fact that usually, PCA bounding boxes quite well approximate the minimum-volume bounding boxes. In this paper we give a lower bound on the approximation factor of PCA bounding boxes of convex polytopes in arbitrary dimension, and an upper bound on the approximation factor of PCA bounding boxes of convex polygons in $\mathbb{R}^2$.

## Keywords

Bounding Boxes, Principal Component Analysis, Computational Geometry.

## 1. INTRODUCTION

Substituting sets of points or complex geometric shapes with their bounding boxes is motivated by many applications. For example, in computer graphics, it is used to maintain hierarchical data structures for fast rendering of a scene or for collision detection. Additional applications include those in shape analysis and shape simplification, or in statistics, for storing and performing range-search queries on a large database of samples.

Computing a minimum-area bounding box of a set of $n$ points in $\mathbb{R}^2$ can be done in $O(n\log n)$ time, for example with the rotating caliper algorithm [Tou83]. O'Rourke [O'R85] presented a deterministic algorithm, a rotating caliper variant in $\mathbb{R}^3$, for computing the exact minimum-volume bounding box of a set of $n$ points in $\mathbb{R}^3$. His algorithm requires $O(n^3)$ time and $O(n)$ space. Barequet and Har-Peled [BHP99] have contributed two $(1+\varepsilon)$-approximation algorithms for computing the minimum-volume bounding box for point sets in $\mathbb{R}^3$, both with nearly linear com-

plexity. The running times of their algorithms are $O(n+1/\varepsilon^{4.5})$ and $O(n\log n + n/\varepsilon^3)$, respectively.

Numerous heuristics have been proposed for computing a box which encloses a given set of points. The simplest heuristic is naturally to compute the axis-aligned bounding box of the point set. Two-dimensional variants of this heuristic include the well-known *R-tree*, the *packed R-tree* [RL85], the $R^*$-*tree* [BKSS90], the $R^+$-*tree* [SRF87], etc.

A frequently used heuristic for computing a bounding box of a set of points is based on *principal component analysis*. The principal components of the point set define the axes of the bounding box. Once the axis directions are given, the dimension of the bounding box is easily found by the extreme values of the projection of the points on the corresponding axis. Two distinguished applications of this heuristic are the OBB-tree [GLM96] and the BOXTREE [BCG$^+$96], hierarchical bounding box structures, which support efficient collision detection and ray tracing. Computing a bounding box of a set of points in $\mathbb{R}^2$ and $\mathbb{R}^3$ by PCA is quite fast, it requires linear time. To avoid the influence of the distribution of the point set on the directions of the PCs, a possible approach is to consider the convex hull, or the boundary of the convex hull $CH(P)$ of the point set $P$. Thus, the complexity of the algorithm increases to $O(n\log n)$. The popularity of this heuristic, besides its speed, lies in its easy implementation and

in the fact that usually PCA bounding boxes are tight-fitting (see [LKM$^+$00] for some experimental results).

Given a point set $P \subseteq \mathbb{R}^d$ we denote by $BB_{pca}(P)$ the PCA bounding box of $P$ and by $BB_{opt}(P)$ the bounding box of $P$ with smallest possible volume. The ratio of the two volumes $\lambda_d(P) = Vol(BB_{pca}(P))/Vol(BB_{opt}(P))$ defines the approximation factor for $P$, and

$$\lambda_d = \sup\left\{\lambda_d(P) \mid P \subseteq \mathbb{R}^d, Vol(CH(P)) > 0\right\}$$

defines the general PCA approximation factor. We are not aware of any previous published results about this quality feature of PCA. Here, we give lower bounds on $\lambda_d$ for arbitrary dimension $d$, and an upper bound on $\lambda_2$.

The paper is organized as follows. In Section 2. we review the basics of principal component analysis. In particular, we present the continuous version of PCA, which results in the introduction of a series of approximation factors $\lambda_{d,i}$, where $i$ ranges from 0 to $d$ and denotes the dimension of the faces of the convex hull that contribute to the continuous point set for which the principal components are computed. In Section 3. we give lower bounds on $\lambda_{d,i}$ for arbitrary values of $d$ and $1 \leq i \leq d$. An upper bound on $\lambda_{2,1}$ is presented in Section 4. We conclude with future work and open problems in Section 5.

## 2. PRINCIPAL COMPONENT ANALYSIS

The central idea and motivation of PCA [Jol02] (also known as the Karhunen-Loeve transform, or the Hotelling transform) is to reduce the dimensionality of a point set by identifying *the most significant directions (principal components)*. Let $X = \{x_1, x_2, \ldots, x_m\}$, where $x_i$ is a $d$-dimensional vector, and $c = (c_1, c_2, \ldots, c_d) \in \mathbb{R}^d$ be the center of gravity of $X$. For $1 \leq k \leq d$, we use $x_{ik}$ to denote the $k$-th coordinate of the vector $x_i$. Given two vectors $u$ and $v$, we use $\langle u, v \rangle$ to denote their inner product. For any unit vector $v \in \mathbb{R}^d$, the *variance of X in direction $v$* is

$$var(X,v) = \frac{1}{m}\sum_{i=1}^{m}\langle x_i - c, v\rangle^2. \qquad (1)$$

The most significant direction corresponds to the unit vector $v_1$ such that $var(X, v_1)$ is maximum. In general, after identifying the $j$ most significant directions $B_j = \{v_1, v_2, \ldots, v_j\}$, the $(j+1)$-th most significant direction corresponds to the unit vector $v_{j+1}$ such that $var(X, v_{j+1})$ is maximum among all unit vectors perpendicular to $v_1, v_2, \ldots, v_j$.

It can be verified that for any unit vector $v \in \mathbb{R}^d$,

$$var(X,v) = \langle Cv, v\rangle, \qquad (2)$$

where $C$ is the *covariance matrix* of $X$. $C$ is a symmetric $d \times d$ matrix where the $(i, j)$-th component, $c_{ij}, 1 \leq i, j \leq d$, is defined as

$$c_{ij} = \frac{1}{m}\sum_{k=1}^{m}(x_{ik} - c_i)(x_{jk} - c_j). \qquad (3)$$

The procedure of finding the most significant directions, in the sense mentioned above, can be formulated as an eigenvalue problem. If $\lambda_1 > \lambda_2 > \cdots > \lambda_d$ are the eigenvalues of $C$, then the unit eigenvector $v_j$ for $\lambda_j$ is the $j$-th most significant direction. All $\lambda_j$s are non-negative and $\lambda_j = var(X, v_j)$. Since the matrix $C$ is symmetric positive definite, its eigenvectors are orthogonal. If the eigenvalues are not distinct, the eigenvectors are not unique. In this case, an orthogonal basis of eigenvectors is chosen arbitrary. However, we can achieve distinct eigenvalues by a slight perturbation of the point set.

The following result summarizes the above background knowledge on PCA. For any set $S$ of orthogonal unit vectors in $\mathbb{R}^d$, we use $var(X,S)$ to denote $\sum_{v \in S} var(X, v)$.
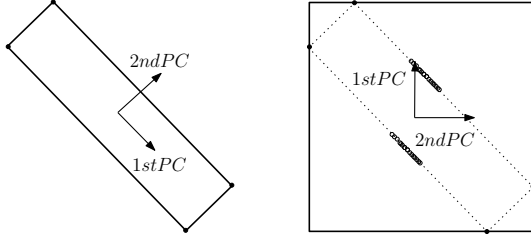
**Lemma 1** *For $1 \leq j \leq d$, let $\lambda_j$ be the $j$-th largest eigenvalue of $C$ and let $v_j$ denote the unit eigenvector for $\lambda_j$. Let $B_j = \{v_1, v_2, \ldots, v_j\}$, $sp(B_j)$ be the linear subspace spanned by $B_j$, and $sp(B_j)^{\perp}$ be the orthogonal complement of $sp(B_j)$. Then $\lambda_1 = \max\{var(X,v) : v \in \mathbb{R}^d, \|v\| = 1\}$, and for any $2 \leq j \leq d$,*

*i) $\lambda_j = \max\{var(X,v) : v \in sp(B_{j-1})^{\perp}, \|v\| = 1\}$.*

*ii) $\lambda_j = \min\{var(X,v) : v \in sp(B_j), \|v\| = 1\}$.*

*iii) $var(X, B_j) \geq var(X, S)$ for any set $S$ of $j$ orthogonal unit vectors.*

Since bounding boxes of a point set $P$ (with respect to any orthogonal coordinate system) depend only on the convex hull of $CH(P)$, the construction of the covariance matrix should be based only on $CH(P)$ and not on the distribution of the points inside. Using the vertices, i.e., the 0-dimensional faces of $CH(P)$ to define the covariance matrix $C$ we obtain a bounding box $BB_{pca(d,0)}(P)$. We denote by $\lambda_{d,0}(P)$ the approximation factor for the given point set $P$ and by

$$\lambda_{d,0} = \sup\left\{\lambda_{d,0}(P) \mid P \subseteq \mathbb{R}^d, Vol(CH(P)) > 0\right\}$$

the approximation factor in general. The example in Figure 1 shows that $\lambda_{2,0}(P)$ can be arbitrarily large if

**Figure 1: Four points and its PCA bounding-box (left). Dense collection of additional points significantly affect the orientation of the PCA bounding-box (right).**

the convex hull is nearly a thin rectangle, but with a lot of additional vertices in the middle of the two long sides. Since this construction can be lifted into higher dimensions we obtain a first general lower bound.

**Proposition 2** $\lambda_{d,0} = \infty$ *for any* $d \geq 2$.

To overcome this problem, one can apply a continuous version of PCA taking into account (the dense set of) all points on the boundary of $CH(P)$, or even all points in $CH(P)$. In this approach $X$ is a continuous set of $d$-dimensional vectors and the coefficients of the covariance matrix are defined by integrals instead of finite sums.

Note that for for $d = 1$ the above problem is trivial, because the PCA bounding box is always optimal, i.e., $\lambda_{1,0}$ and $\lambda_{1,1}$ are 1.

## 2.1 Continuous PCA

Variants of the continuous PCA, applied on triangulated surfaces of 3D objects, were presented by Gottschalk et. al. [GLM96], Lahanas et. al. [LKM$^+$00] and Vranić et. al. [VSR01]. In what follows, we briefly review the basics of the continuous PCA in a general setting.

Let $X$ be a continuous set of $d$-dimensional vectors with constant density. Then, the center of gravity of $X$ is

$$c = \frac{\int_{x \in X} x \, dx}{\int_{x \in X} dx}. \tag{4}$$

Here, $\int dx$ denotes either a line integral, an area integral, or a volume integral in higher dimensions. For any unit vector $v \in \mathbb{R}^d$, the *variance of X in direction v* is

$$var(X, v) = \frac{\int_{x \in X} \langle x - c, v \rangle^2 dx}{\int_{x \in X} dx}. \tag{5}$$

The covariance matrix of $X$ has the form

$$C = \frac{\int_{x \in X} (x - c)(x - c)^T dx}{\int_{x \in X} dx}, \tag{6}$$

with its $(i, j)$-th component

$$c_{ij} = \frac{\int_{x \in X} (x_i - c_i)(x_j - c_j) dx}{\int_{x \in X} dx}, \tag{7}$$

where $x_i$ and $x_j$ are the $i$-th and $j$-th component of the vector $x$, and $c_i$ and $c_j$ $i$-th and $j$-th component of the center of gravity. It can be verified that relation (2) is also true when $X$ is a continuous set of vectors. The procedure of finding the most significant directions, can be also reformulated as an eigenvalue problem and consequently Lemma 1 holds.

For point sets $P$ in $\mathbb{R}^2$ we are especially interested in the cases when $X$ represents the boundary of $CH(P)$, or all points in $CH(P)$. Since the first case corresponds to the 1-dimensional faces of $CH(P)$ and the second case to the only 2-dimensional face of $CH(P)$, the generalization to a dimension $d > 2$ leads to a series of $d - 1$ continuous PCA versions. For a point set $P \in \mathbb{R}^d$, $C(P, i)$ denotes the covariance matrix defined by the points on the $i$-dimensional faces of $CH(P)$, and $BB_{pca(d,i)}(P)$, denotes the corresponding bounding box. The approximation factors $\lambda_{d,i}(P)$ and $\lambda_{d,i}$ are defined as

$$\lambda_{d,i}(P) = \frac{Vol(BB_{pca(d,i)}(P))}{Vol(BB_{opt}(P))}, \quad \text{and}$$
$$\lambda_{d,i} = \sup \left\{ \lambda_{d,i}(P) \mid P \subseteq \mathbb{R}^d, Vol(CH(P)) > 0 \right\}.$$

## 3. LOWER BOUNDS

We start with straightforward conclusion from Proposition 2.

**Proposition 3** $\lambda_{d,i} = \infty$ *for any* $d \geq 4$ *and any* $1 \leq i < d - 1$.

**Proof**. We can use a lifting argument to establish $\lambda_{k,i} \leq \lambda_{k+1,i+1}$, and thus $\lambda_{d,i} \geq \lambda_{d-1,i-1} \geq \ldots \geq \lambda_{d-i,0} = \infty$. $\square$

This way, there remain only two interesting cases for a given $d$: the factor $\lambda_{d,d-1}$ corresponding to the boundary of the convex hull, and the factor $\lambda_{d,d}$ corresponding to the full convex hull. The nontrivial lower bounds we are going to derive are based on the following connection between the symmetry of a point set and its principal components.

**Lemma 4** *Let P be a d-dimensional point set symmetric with respect to a hyperplane H and assume that the*

*covariance matrix C has d different eigenvalues. Then, a principal component of P is orthogonal to H.*

**Proof**. Without loss of generality, we can assume that the hyperplane of symmetry is spanned by the last $d-1$ standard base vectors of the $d$-dimensional space and the center of gravity of the point set coincides with the origin of the $d$-dimensional space, i.e., $c = (0,0,\ldots,0)$. Then, the components $c_{1j}$ and $c_{j1}$, for $2 \leq j \leq d$, are 0, and the covariance matrix has the form:

$$C = \begin{bmatrix} c_{11} & 0 & \ldots & 0 \\ 0 & c_{22} & \ldots & c_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & c_{d2} & \ldots & c_{dd} \end{bmatrix} \tag{8}$$

Its characteristic polynomial is

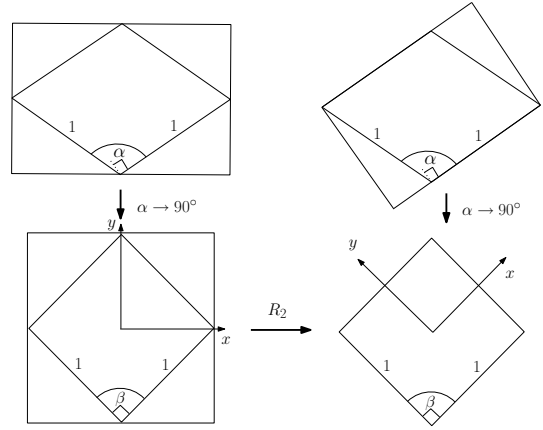$$det(C - \lambda\,I) = (c_{11} - \lambda)f(\lambda), \tag{9}$$

where $f(\lambda)$ is a polynomial of degree $d-1$, with coefficients determined by the elements of the $(d-1) \times (d-1)$ submatrix of $C$. From this it follows that $c_{11}$ is a solution of the characteristic equation, i.e., it is an eigenvalue of $C$ and the vector $(1, 0, \ldots,0)$ is its corresponding eigenvector (principal component), which is orthogonal to the assumed hyperplane of symmetry. □

## 3.1 Lower bounds in $\mathbb{R}^2$

The result obtained in this subsection can be seen as special case of the result obtained in the subsection 3.3. To gain a better understanding of the problem and the obtained results, we consider it separately.

**Theorem 5** $\lambda_{2,1} \geq 2$ *and* $\lambda_{2,2} \geq 2$.

**Proof**. Both lower bounds can be derived from a rhombus. Let the side length of the rhombus be 1. Since the rhombus is symmetric, its PCs coincide with its diagonals. On the right side in Figure 2 its optimal-area bounding boxes, for 2 different angles, $\alpha > 90°$ and $\beta = 90°$, are shown, and on the left side its corresponding PCA bounding boxes. As the rhombus' angles in limit approach $90°$, the rhombus approaches a square with side length 1, i.e., the vertices of the rhombus in the limit are $(\frac{1}{\sqrt{2}},0), (-\frac{1}{\sqrt{2}},0), (0,\frac{1}{\sqrt{2}})$ and $(0,-\frac{1}{\sqrt{2}})$ (see the left side in Figure 2), and the dimensions of its PCA bounding box are $\sqrt{2} \times \sqrt{2}$. According to Lemma 4, the PCs of the rhombus are unique



**Figure 2: An example which gives us the lower bound of the area of the PCA bounding box of an arbitrary convex polygon in $\mathbb{R}^2$.**

as long its angles are not $90°$. This leads to the conclusion that the ratio between the area of the bounding box on the left side in Figure 3, and the area of its PCA bounding box, on the right side in Figure 3, in limit goes to 2. □

Alternatively, to show that the given squared rhombus fits into a unit cube, one can apply the following rotation matrix

$$R_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}. \tag{10}$$

It can be verified easily that all coordinates of the vertices of the rhombus transformed by $R_2$ are in the interval $[-0.5, 0.5]$. We use similar arguments when we prove the lower bounds in higher dimensions.
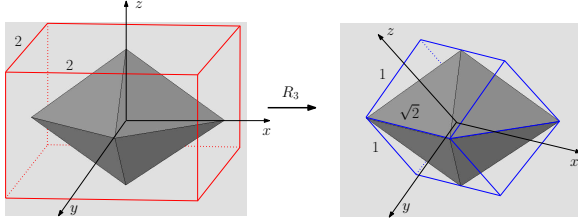
## 3.2 Lower bounds in $\mathbb{R}^3$

**Theorem 6** $\lambda_{3,2} \geq 4$ *and* $\lambda_{3,3} \geq 4$.

**Proof**. Both lower bounds are obtained from a dipyramid, having a rhombus with side length $\sqrt{2}$ as its base. The other sides of the dipyramid have length $\frac{\sqrt{3}}{2}$. Similarly as in $\mathbb{R}^2$, we consider the case when its base, the rhombus, in limit approaches the square, i.e., the vertices of the square dipyramid are $(1,0,0),(-1,0,0),(0,1,0),(0,-1,0),(0,0,\frac{\sqrt{2}}{2})$ and $(0,0,-\frac{\sqrt{2}}{2})$ (see the left side in Figure 3). The dimensions of its PCA bounding box are $2 \times 2 \times \sqrt{2}$. Now, we rotate the coordinate system (or the square

| dimension | $\mathbb{R}$ | $\mathbb{R}^2$ | $\mathbb{R}^3$ | $\mathbb{R}^4$ | $\mathbb{R}^5$ | $\mathbb{R}^6$ | $\mathbb{R}^7$ | $\mathbb{R}^8$ | $\mathbb{R}^9$ | $\mathbb{R}^{10}$ |
|---|---|---|---|---|---|---|---|---|---|---|
| lower bound | 1 | 2 | 4 | 16 | 16 | 32 | 64 | 4096 | 4096 | 8192 |

Table 1: Lower bounds for the approximation factor of PCA bounding boxes for the first 10 dimensions.



Figure 3: An example which gives the lower bound of the volume of the PCA bounding box of an arbitrary convex polygon in $\mathbb{R}^3$.

dipyramid) with the rotation determined by the following orthogonal matrix

$$
R_3 = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{\sqrt{2}} \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{\sqrt{2}} \end{bmatrix}. \qquad (11)
$$

It can be verified easily that the square dipyramid, after rotation with $R_3$ fits into the box $[-0.5, 0.5]^3$ (see the right side in Figure 3). Thus, the ratio of the volume of the bounding box, on the left side in Figure 3, and the volume of its PCA bounding box, on the right side in Figure 3, in limit goes to 4.  $\square$

## 3.3 Lower bounds in $\mathbb{R}^d$

**Theorem 7** *If $d$ is a power of two, then $\lambda_{d,d-1} \geq \sqrt{d}^d$ and $\lambda_{d,d} \geq \sqrt{d}^d$.*

**Proof**. For any $d = 2^k$, let $a_i$ be a $d$-dimensional vector, with $a_{ii} = \frac{\sqrt{d}}{2}$ and $a_{ij} = 0$ for $i \neq j$, and let $b_i = -a_i$. We construct a $d$-dimensional convex polytope $P_d$ with vertices $V = \{a_i, b_i | 1 \leq i \leq d\}$. It is easy to check that the hyperplane normal to $a_i$ is a hyperplane of reflective symmetry, and as consequence of Lemma 4, $a_i$ is an eigenvector of the covariance matrix of $P_d$. To ensure that all eigenvalues are different (which implies that the PCA bounding box is unique), we add $\varepsilon_i > 0$ to the $i$-th coordinate of $a_i$, and $-\varepsilon_i$ to the $i$-th coordinate of $b_i$, for $1 \leq i \leq d$, where $\varepsilon_1 < \varepsilon_2 < \ldots < \varepsilon_d$. When all $\varepsilon_i$, $1 \leq i \leq d$, arbitrary approach 0, the PCA bounding box of the convex polytope $P_d$ converges to

a hypercube with side lengths $\sqrt{d}$, i.e., the volume of the PCA bounding box of $P_d$ converges to $\sqrt{d}^d$. Now, we rotate $P_d$, such that it fits into the cube $[-\frac{1}{2}, \frac{1}{2}]^d$. For $d = 2^k$, we can use a rotation matrix derived from a *Hadamard matrix*[1], recursively defined by

$$
R_d = \frac{1}{\sqrt{2}} \left[ \begin{array}{c|c} R_{\frac{d}{2}} & R_{\frac{d}{2}} \\ \hline R_{\frac{d}{2}} & -R_{\frac{d}{2}} \end{array} \right], \qquad (12)
$$

where we start with the matrix $R_2$ (10) defined above for $d = 2$. A straightforward calculation verifies that $P_d$ rotated with $R_d$ fits into the cube $[-0.5, 0.5]^d$.  $\square$

**Remark:** Theorem 7 holds for all dimensions $d$ for which a $d \times d$ Hadamard matrix exists. Hadamard conjectured that this is the case for all multiples of four. This conjecture is known to be true for $d \leq 664$ [KTR05].
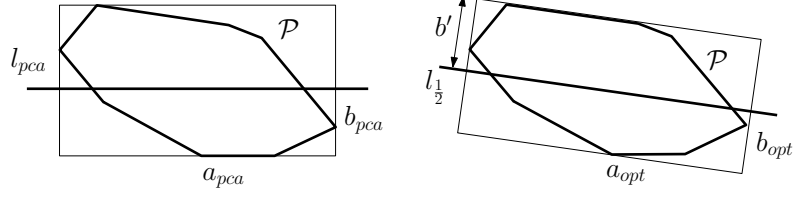
We can combine lower bounds from lower dimensions to get lower bounds in higher dimensions by taking Cartesian products. If $l_{d_1}$ is a lower bound for the ratio between the PCA bounding box and the optimal bounding box of a convex polytope in $\mathbb{R}^{d_1}$, and $l_{d_2}$ is a lower bound in $\mathbb{R}^{d_2}$, then $l_{d_1} \cdot l_{d_2}$ is a lower bound in $\mathbb{R}^{d_1+d_2}$. This observation together with the results from this section enables us to obtain lower bounds in any dimension. For example, for the first 10 dimensions, the lower bounds we obtain are given in Table 1.

## 4. AN UPPER BOUND FOR $\lambda_{2,1}$

Given a point set $P \subseteq \mathbb{R}^2$ and an arbitrary bounding box $BB(P)$ we will denote the two side lengths by $a$ and $b$, where $a \geq b$. We are interested in the side lengths $a_{opt}(P) \geq b_{opt}(P)$ and $a_{pca}(P) \geq b_{pca}(P)$ of $BB_{opt}(P)$ and $BB_{pca(2,1)}(P)$, see Figure 4. The parameters $\alpha = \alpha(P) = a_{pca}(P)/a_{opt}(P)$ and $\beta = \beta(P) = b_{pca}(P)/b_{opt}(P)$ denote the ratios between the corresponding side lengths. Hence, we have $\lambda_{2,1}(P) = \alpha(P) \cdot \beta(P)$. If the relation to $P$ is clear, we will omit the reference to $P$ in the notations introduced above.

Since the side lengths of any bounding box are bounded by the diameter of $P$, we can observe that in

---
[1] A Hadamard matrix is a $\pm 1$ matrix with orthogonal columns.

**Figure 4: A convex polygon $\mathcal{P}$, its PCA bounding box and the line $l_{pca}$, which coincides with the first principal component of $\mathcal{P}$, are given in the left part of the figure. The optimal bounding box and the line $l_{\frac{1}{2}}$, going through the middle of its smaller side, parallel with its longer side, are given in the right part of the figure.**

general $b_{pca}(P) \leq a_{pca}(P) \leq diam(P) \leq \sqrt{2}a_{opt}(P)$, and in the special case when the optimal bounding box is a square $\lambda_{2,1}(P) \leq 2$. This observation can be generalized, introducing an additional parameter $\eta(P) = a_{opt}(P)/b_{opt}(P)$.

**Lemma 8** $\lambda_{2,1}(P) \leq \eta + \frac{1}{\eta}$ and $\lambda_{2,2}(P) \leq \eta + \frac{1}{\eta}$ *for any point set P with fixed aspect ratio $\eta(P) = \eta$.*
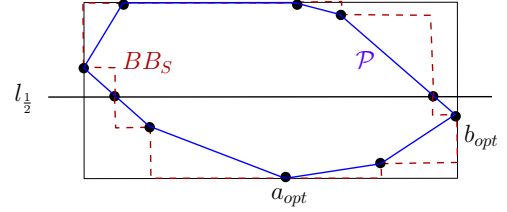
**Proof.** We have for both $a_{pca}$ and $b_{pca}$ the upper bound $diam(P) \leq \sqrt{a_{opt}^2 + b_{opt}^2} = a_{opt}\sqrt{1 + \frac{1}{\eta^2}}$. Replacing $a_{opt}$ by $\eta \cdot b_{opt}$ in the bound for $b_{pca}$ we obtain $\alpha\beta \leq \eta\left(\sqrt{1 + \frac{1}{\eta^2}}\right)^2 = \eta + \frac{1}{\eta}$. $\qquad\square$

Unfortunately, this parametrized upper bound tends to infinity for $\eta \rightarrow \infty$. Therefore we are going to derive another upper bound that is better for large values of $\eta$. In this process we will make essential use of the properties of $BB_{pca(2,1)}(P)$. In order to distinguish clearly between a convex set and its boundary, we will use calligraphic letters for the boundaries, especially $\mathcal{P}$ for the boundary of $CH(P)$ and $\mathcal{BB}_{opt}$ for the boundary of the rectangle $BB_{opt}(P)$. Furthermore, we denote by $d^2(\mathcal{P}, l)$ the integral of the squared distances of the points on $\mathcal{P}$ to a line $l$, i.e., $d^2(\mathcal{P}, l) = \int_{x \in \mathcal{P}} d^2(x, l)ds$. Let $l_{pca}$ be the line going through the center of gravity and parallel to the longer side of $BB_{pca(2,1)}(P)$ and $l_{\frac{1}{2}}$ be the bisector of $BB_{opt(P)}$ parallel to the longer side. By Lemma 1, part ii) $l_{pca}$ is the best fitting line of $P$ and therefore

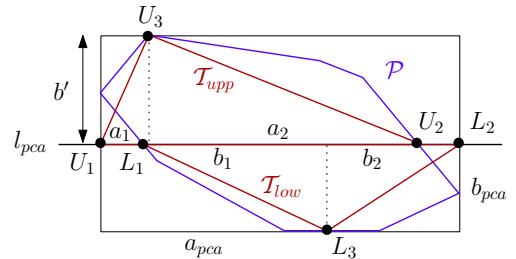$$d^2(\mathcal{P}, l_{pca}) \leq d^2(\mathcal{P}, l_{\frac{1}{2}}). \qquad (13)$$

**Lemma 9** $d^2(\mathcal{P}, l_{\frac{1}{2}}) \leq \frac{b_{opt}^2 a_{opt}}{2} + \frac{b_{opt}^3}{6}$.

**Proof.** If a segment of $\mathcal{P}$ intersects the line $l_{\frac{1}{2}}$, we split this segment into two segments, with the intersection point as a split point. Then, to each segment $f$ of $\mathcal{P}$ flush with the side of the PCA bounding
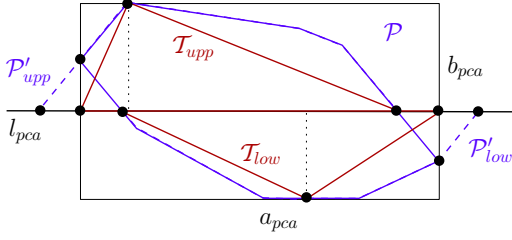


**Figure 5: The convex polygon $\mathcal{P}$, its optimal bounding box, and the staircase polygon $BB_S$ (depicted dashed).**

box, we assign a segment identical to $f$. To each remaining segment $s$ of $\mathcal{P}$, with endpoints $(x_1, y_1)$ and $(x_2, y_2)$, with $|y_1| \leq |y_2|$, we assign two segments: a segment $s_1$, with endpoints $(x_1, y_1)$ and $(x_1, y_2)$, and a segment $s_2$, with endpoints $(x_1, y_2)$ and $(x_2, y_2)$. All these segments form the boundary $\mathcal{BB}_S$ of a staircase polygon (see Figure 5 for illustration). Two straightforward consequences are that $d^2(\mathcal{BB}_S, l_{\frac{1}{2}}) \leq d^2(\mathcal{BB}_{opt}, l_{\frac{1}{2}})$, and $d^2(s, l_{\frac{1}{2}}) \leq d^2(s_1, l_{\frac{1}{2}}) + d^2(s_2, l_{\frac{1}{2}})$, for each segment $s$ of $\mathcal{P}$. Therefore, $d^2(\mathcal{P}, l_{\frac{1}{2}})$ is at most $d^2(\mathcal{BB}_S, l_{\frac{1}{2}})$, which is bounded from above by $d^2(\mathcal{BB}_{opt}, l_{\frac{1}{2}}) = 4\int_0^{\frac{b_{opt}}{2}} x^2\,dx + 2\int_0^{a_{opt}} (\frac{b_{opt}}{2})^2\,dx = \frac{b_{opt}^2 a_{opt}}{2} + \frac{b_{opt}^3}{6}$. $\qquad\square$



**Figure 6: The convex polygon $\mathcal{P}$, its PCA bounding box, and a construction for a lower bound for $d^2(\mathcal{P}, l_{pca})$**

**Figure 7: Two polylines $\mathcal{P}'_{upp}$ and $\mathcal{P}'_{low}$ (depicted dashed) formed from $\mathcal{P}$.**



**Figure 8: Two types of chains of segments (depicted dashed and denoted by $R$), and their corresponding triangles' edges (depicted solid and denoted by $t$).**
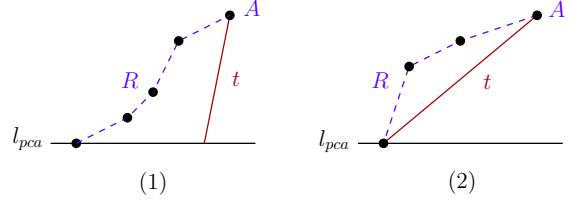
Now we look at $\mathcal{P}$ and its PCA bounding box (Figure 6). The line $l_{pca}$ divides $\mathcal{P}$ into an upper and a lower part, $\mathcal{P}_{upp}$ and $\mathcal{P}_{low}$. $l_{upp}$ denotes the orthogonal projection of $\mathcal{P}_{upp}$ onto $l_{pca}$, with $U_1$ and $U_2$ as its extreme points, and $l_{low}$ denotes the orthogonal projection of $\mathcal{P}_{low}$ onto $l_{pca}$, with $L_1$ and $L_2$ as its extreme points. $\mathcal{T}_{upp} = \triangle(U_1 U_2 U_3)$ is a triangle inscribed in $\mathcal{P}_{upp}$, where point $U_3$ lies on the intersection of $\mathcal{P}_{upp}$ with the upper side of the PCA bounding box. Analogously, $\mathcal{T}_{low} = \triangle(L_1 L_2 L_3)$ is a triangle inscribed in $\mathcal{P}_{low}$.

**Lemma 10**

$$d^2(\mathcal{P}, l_{pca}) \geq d^2(\mathcal{T}_{upp}, l_{pca}) + d^2(\mathcal{T}_{low}, l_{pca}).$$

**Proof.** Let $Q$ denote a chain of segments of $\mathcal{P}$, which does not touch the longer side of the PCA bounding box, and whose one endpoint lies on the smaller side of the PCA bounding box, and the other endpoint on the line $l_{pca}$. We reflect $Q$ at the line supporting the side of the PCA bounding box touched by $Q$. All such reflected chains of segments, together with the rest of $\mathcal{P}$, form two polylines: $\mathcal{P}'_{upp}$ and $\mathcal{P}'_{low}$ (see Figure 7 for illustration). As a consequence, to each of the sides of the triangles $\mathcal{T}_{low}$ and $\mathcal{T}_{upp}$, $\overline{L_1 L_3}$, $\overline{L_2 L_3}$, $\overline{U_1 U_3}$, $\overline{U_2 U_3}$, we have a corresponding chain of segments $R$ as shown in the two cases in Figure 8. In both cases $d^2(t, l_{pca}) \leq d^2(R, l_{pca})$. Namely, we can parametrize both curves, $R$ and $t$, starting at the common endpoint $A$ that is furthest from $l_{pca}$. By comparing two points with the same parameter (distance from $A$ along the curve) we see that the point on $t$ always has a smaller distance to $l_{pca}$ than the corresponding point on $R$. In addition $t$ is shorter, and some parts of $R$ have no match on $t$.

Consequently, $d^2(\mathcal{P}', l_{pca}) \geq d^2(\mathcal{T}_{upp} \bigcup \mathcal{T}_{low}, l_{pca}) = d^2(\mathcal{T}_{upp}, l_{pca}) + d^2(\mathcal{T}_{low}, l_{pca})$, and since, $d^2(\mathcal{P}', l_{pca}) = d^2(\mathcal{P}, l_{pca}) = d^2(\mathcal{P}_{upp} \bigcup \mathcal{P}_{low}, l_{pca})$, the proof is completed. □

Since $\mathcal{P}$ is convex, the following relations hold:

$$|l_{upp}| \geq \frac{b'}{b_{pca}} a_{pca}, \text{ and } |l_{low}| \geq \frac{b_{pca} - b'}{b_{pca}} a_{pca}. \quad (14)$$

The value

$$
\begin{aligned}
d^2(\mathcal{T}_{upp}, l_{pca}) &= \int_0^{\sqrt{a_1^2 + b'^2}} \left(\frac{\alpha}{\sqrt{a_1^2 + b'^2}} b'\right)^2 d\alpha \\
&+ \int_0^{\sqrt{a_2^2 + b'^2}} \left(\frac{\alpha}{\sqrt{a_2^2 + b'^2}} b'\right)^2 d\alpha \\
&= \frac{b'^2}{3} \left(\sqrt{a_1^2 + b'^2} + \sqrt{a_2^2 + b'^2}\right)
\end{aligned}
$$

is minimal when $a_1 = a_2 = \frac{|l_{upp}|}{2}$. With (14) we get

$$d^2(\mathcal{T}_{upp}, l_{pca}) \geq \frac{b'^3}{3 b_{pca}} \sqrt{a_{pca}^2 + 4 b_{pca}^2}.$$

Analogously, we have for the lower part:

$$d^2(\mathcal{T}_{low}, l_{pca}) \geq \frac{(b_{pca} - b')^3}{3 b_{pca}} \sqrt{a_{pca}^2 + 4 b_{pca}^2}.$$

The sum $d^2(\mathcal{T}_{upp}, l_{pca}) + d^2(\mathcal{T}_{low}, l_{pca})$ is minimal when $b' = \frac{b_{pca}}{2}$. This, together with Lemma 10, gives:

$$d^2(\mathcal{P}, l_{pca}) \geq \frac{b_{pca}^2}{12} \sqrt{a_{pca}^2 + 4 b_{pca}^2}. \quad (15)$$

Combining (13), (15) and Lemma 9 we have:

$$\frac{1}{2} a_{opt} b_{opt}^2 + \frac{1}{6} b_{opt}^3 \geq \frac{b_{pca}^2}{12} \sqrt{a_{pca}^2 + 4 b_{pca}^2} \geq \frac{b_{pca}^2}{12} a_{pca}. \quad (16)$$

Replacing $a_{opt}$ with $\eta b_{opt}$ on the left side, $b_{pca}^2$ with $\beta^2 b_{opt}^2$ and $a_{pca}$ with $\alpha a_{opt} = \alpha \eta b_{opt}$ on the right side of (16), we obtain:

$$\left(\frac{\eta}{2} + \frac{1}{6}\right) b_{opt}^3 \geq \frac{\beta^2 \alpha \eta}{12} b_{opt}^3$$

which implies

$$\beta \leq \sqrt{\frac{6\eta + 2}{\alpha \eta}}.$$

This gives the second upper bound on $\lambda_{2,1}(P)$ for point sets with parameter $\eta$:

$$\alpha\beta \leq \sqrt{\frac{(6\eta+2)\alpha}{\eta}} \leq \sqrt{\frac{6\eta+2}{\eta}\sqrt{1+\frac{1}{\eta^2}}} \quad (17)$$

**Theorem 11** *The PCA bounding box of a point set $P$ in $\mathbb{R}^2$ computed over the boundary of $CH(P)$ has a guaranteed approximation factor $\lambda_{2,1} \leq 2.737$.*

**Proof**. The theorem follows from the combination of the two parametrised bounds from Lemma 8 and (17) proved above:

$$\lambda_{2,1} \leq \sup_{\eta \geq 1}\left\{\min\left(\eta+\frac{1}{\eta}, \sqrt{\frac{6\eta+2}{\eta}\sqrt{1+\frac{1}{\eta^2}}}\right)\right\}.$$

It is easy to check that the supremum $s \approx 2.736$ is obtained for $\eta \approx 2.302$. $\qquad\square$

## 5. FUTURE WORK AND OPEN PROBLEMS

It should be possible to prove an upper bound on $\lambda_{2,2}$ along the same line as for $\lambda_{2,1}$, but the analogon of Lemma 9 seems to require some new analytical tools, since, e.g., the reflection tricks do not apply in that setting. However, there is some evidence that an upper bound proof for $\lambda_{2,2}$ would give some ideas to attack the 3-dimensional problem for $\lambda_{3,3}$, and, maybe also a generalization to $\lambda_{d,d}$ in higher dimensions.

## REFERENCES

[BCG⁺96]   G. Barequet, B. Chazelle, L. J. Guibas, J. S. B. Mitchell, and A. Tal. Boxtree: A hierarchical representation for surfaces in 3D. *Computer Graphics Forum*, 15:387–396, 1996.

[BHP99]   G. Barequet and S. Har-Peled. Efficiently approximating the minimum-volume bounding box of a point set in 3d. In *10th ACM-SIAM Sympos. Discrete Algorithms*, pages 82–91, 1999.

[BKSS90]   N. Beckmann, H.-P. Kriegel, R. Schneider, and B. Seeger. The R*-tree: An efficient and robust access method for points and rectangles. *ACM SIGMOD Int. Conf. on Manag. of Data*, pages 322–331, 1990.

[GLM96]   S. Gottschalk, M. C. Lin, and D. Manocha. OBBTree: A hierarchical structure for rapid interference detection. In *SIGGRAPH 1996*, pages 171–180, 1996.

[Jol02]   I. Jolliffe. *Principal Component Analysis*. Springer-Verlag, New York, 2nd ed., 2002.

[KTR05]   H. Kharaghani and B. Tayfeh-Rezaie. A hadamard matrix of order 428. In *J. Combin. Designs 13*, pages 435–440, 2005.

[LKM⁺00]   M. Lahanas, T. Kemmerer, N. Milickovic, D. Baltas K. Karouzakis, and N. Zamboglou. Optimized bounding boxes for three-dimensional treatment planning in brachytherapy. In *Med. Phys. 27*, pages 2333–2342, 2000.

[O'R85]   J. O'Rourke. Finding minimal enclosing boxes. In *Int. J. Comp. Info. Sci. 14*, pages 183–199, 1985.

[RL85]   N. Roussopoulos and D. Leifker. Direct spatial search on pictorial databases using packed R-trees. In *ACM SIGMOD*, pages 17–31, 1985.

[SRF87]   T. Sellis, N. Roussopoulos, and C. Faloutsos. The R+-tree: A dynamic index for multidimensional objects. In *13th VLDB Conference*, pages 507–518, 1987.

[Tou83]   G. Toussaint. Solving geometric problems with the rotating calipers. In *IEEE MELECON*, May 1983.

[VSR01]   D. V. Vranić, D. Saupe, and J. Richter. Tools for 3d-object retrieval: Karhunen-Loeve transform and spherical harmonics. In *IEEE 2001 Workshop Multimedia Signal Processing*, pages 293–298, 2001.