

ε -Nets of Two Sets and Their Application to the Classification Problem

Maria A. Ivanchuk^{1*}, and Igor V. Malyk²

¹ Department of Biological Physics and Medical Informatics, Bukovinian State Medical University, Ukraine

² Department of the System Analysis and Insurance and Financial Mathematics, Yuriy Fedkovych Chernivtsi

National University, Ukraine

Email: mgracia@ukr.net

Abstract. The separation algorithm of linear two sets using their ε -nets in the range space (R^d, H^d) is proposed in the paper. The algorithm is illustrated by two examples for normal and uniform distributions. The set of possible values of ε and its properties are considered in the manuscript.

Keywords: Epsilon-nets, sets' separation, linear classification

1 Introduction

In 1987 D. Haussler and E. Welzl [6] introduced ε -nets. Since that time ε -nets are used in computational and combinatorics geometry [1,4,5,10,11]. Numerous works study ε -nets of one set. In this paper we will build ε -nets of two sets for solving the classification problem. In the first part the algorithm of building two separable ε -nets with respect to halfspaces is proposed. The algorithm is illustrated by examples for normal and uniform distributions. In the second part the separating algorithm of two ε -nets is proposed and illustrated. The result of classification using ε -nets is compared with SVM-method.

2 Building Separable ε -Nets

Definition 1. Sets A and B are called **ε -separable** if there exist sets $A_1 \subset A$, $B_1 \subset B$, such that

$$\text{conv}(A \setminus A_1) \cap \text{conv}(B \setminus B_1) = \emptyset \quad (1)$$

and

$$|A_1| + |B_1| < \varepsilon(n_A + n_B) \quad (2)$$

Definition 2. Hyperplane L is called **separating** for the sets A and B if $\text{conv}_A \subset L^+$, $\text{conv}_B \subset L^-$.

Definition 3. Hyperplane L_ε is called **ε -separating** for the sets A and B if

$$\frac{|A \cap L_\varepsilon^+| + |B \cap L_\varepsilon^-|}{n_A + n_B} \geq 1 - \varepsilon$$

Consider an infinite range space (R^d, H^d) , where H^d is the closed halfspaces in R^d bounded by hyperplanes.

Theorem 1. [7] A necessary and sufficient condition that two sets of points A and B are ε -separable is there exist $\varepsilon_A, \varepsilon_B$ and corresponding ε -nets $N_{\varepsilon_A}^A$, $N_{\varepsilon_B}^B$ in (R^d, H^d) such that

$$\varepsilon_A n_A + \varepsilon_B n_B < \varepsilon(n_A + n_B) \quad (3)$$

and

$$\text{conv}N_{\varepsilon_A}^A \cap \text{conv}N_{\varepsilon_B}^B = \emptyset \quad (4)$$

In the manuscript we will consider two examples. In the first example we will consider two sets which are generated by the normal distribution. In the second one the sets are generated by the uniform distribution. Algorithms proposed in the paper are implemented in the Matlab.

Example 1.

1. Normal distribution

Consider two sets A and B , which are generated from the normal distributions with parameters $n_A = 500$, $\mu_A = (3;5)$, $\sigma^A = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$, $n_B = 500$, $\mu_B = (9;9)$, $\sigma^B = \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}$ (Fig. 1).

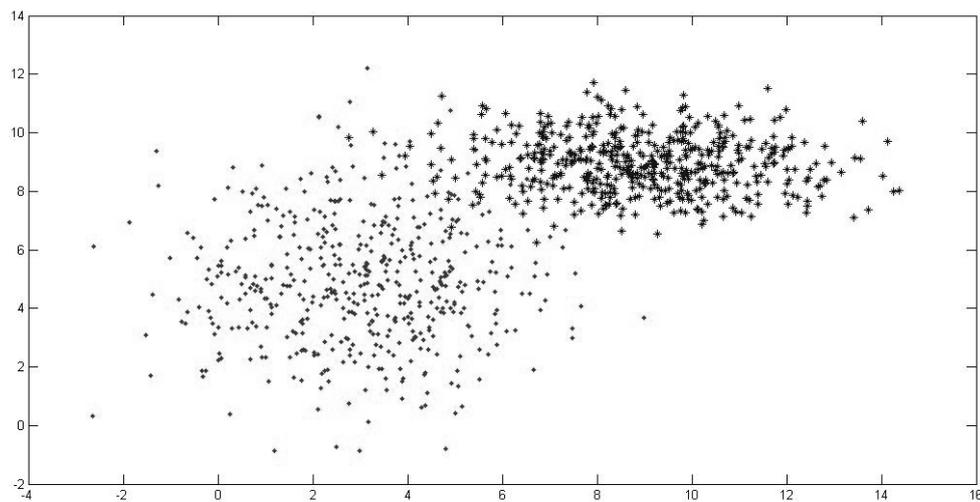


Figure 1. Two sets generated by the normal distribution

2. Uniform distribution

Consider two sets A and B , which are generated from the uniform distributions with parameters $n_A = 1000$, $a_A = (1;1)$, $b_A = (5;5)$, $n_B = 1000$, $a_B = (4;4)$, $b_B = (8;8)$ (Fig. 2)

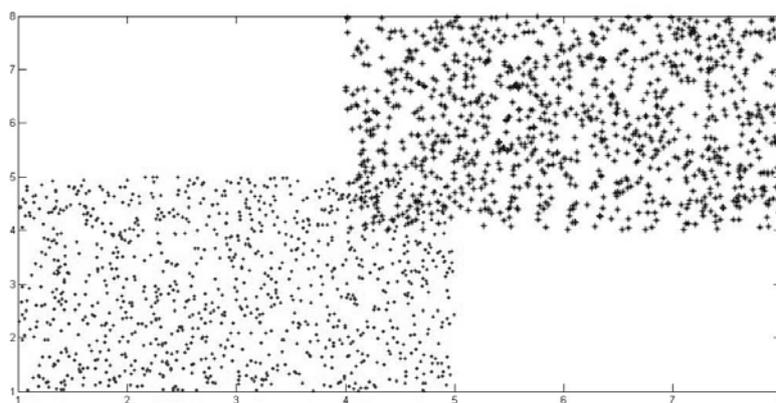


Figure 2. Two sets generated by the uniform distribution

We will build the separable ε -nets $N_{\varepsilon_A}^A$, $N_{\varepsilon_B}^B$ for two ε -separable sets A, B . According to theorem 1, $\varepsilon_A, \varepsilon_B$ have to satisfy the condition (3).

Definition 4. *The set*

$$D_{A,B} = \left\{ (\varepsilon_1, \varepsilon_2) \in (0,1)^2 : \exists N_A^{\varepsilon_1}, N_B^{\varepsilon_2}, \text{conv}N_A^{\varepsilon_1} \cap \text{conv}N_B^{\varepsilon_2} = \emptyset \right\} \tag{5}$$

is called the separation space for A, B.

It is clear that condition (4) holds for $\varepsilon_A, \varepsilon_B \in D_{A,B}$. So, sets A and B are ε -separable if there exists $(\varepsilon_A, \varepsilon_B) \in D_{A,B}$ that satisfies the condition (3). It's enough to show that

$(\varepsilon_A^0, \varepsilon_B^0) = \arg \min_{(\varepsilon_A, \varepsilon_B) \in D_{A,B}} \frac{\varepsilon_A n_A + \varepsilon_B n_B}{n_A + n_B}$ satisfies the condition (3). If condition (3) does not hold for the

$(\varepsilon_A^0, \varepsilon_B^0)$, it does not hold for all $(\varepsilon_A, \varepsilon_B) \in D_{A,B}$, which means sets A and B are not ε -separable.

In the Fig. 3 you can see the bounded line for $D_{A,B}$ and the point (x_0, y_0) , which is found as solution

of the minimization problem $\frac{\varepsilon_A n_A + \varepsilon_B n_B}{n_A + n_B} \rightarrow \min$ with condition $(\varepsilon_A, \varepsilon_B) \in D_{A,B}$.

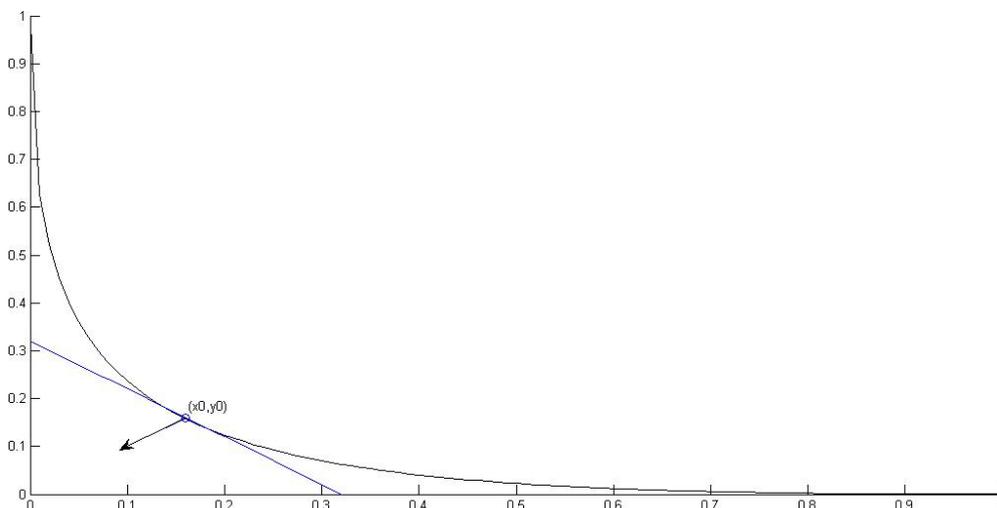


Figure 3. Bounded line for $D_{A,B}$ and the point $(\varepsilon_A^0, \varepsilon_B^0)$

Let $\xi, \eta \in R^1$ be continuous random variables with distribution functions F_ξ, F_η .

Definition 5. *The set D_l*

$$D_l := \left\{ (x, y) \in (0,1)^2 : \exists h \in R^1, P\{\xi \in h_+\} \leq x, P\{\eta \in h_-\} \leq y \right\} \tag{6}$$

is called the separation space for ξ, η .

Lemma 1. *Let the inverse function F_ξ^{-1} exist. Then the sets D_l and $\bar{D}_l := (0,1)^2 \setminus D_l$ are separated by the line*

$$y(x) = \min \left(F_\eta \left(F_\xi^{-1} (1-x) \right), 1 - F_\eta \left(F_\xi^{-1} (x) \right) \right) \tag{7}$$

Proof. Consider two possible cases.

1. Let the inequality $F_\xi(h) > F_\eta(h)$ hold for $h \in (-\infty, \infty)$, then the set D_l is described by the system of inequality

$$\begin{cases} x \geq 1 - F_\xi(h) \\ y \geq F_\eta(h) \end{cases}$$

Then the line that separates sets D_l and \bar{D}_l is

$$y(x) = F_\eta \left(F_\xi^{-1} (1 - x) \right)$$

2. Let the inequality $F_\xi(h) \leq F_\eta(h)$ holds for $h \in (-\infty, \infty)$, then the set D_l is described by the system of inequality

$$\begin{cases} x \geq F_\xi(h) \\ y \geq 1 - F_\eta(h) \end{cases}$$

In this case the line that separates sets D_l and \bar{D}_l is

$$y(x) = 1 - F_\eta \left(F_\xi^{-1} (x) \right)$$

So, in general, sets D_l and \bar{D}_l are separated by the line

$$y(x) = \min \left(F_\eta \left(F_\xi^{-1} (1 - x) \right), 1 - F_\eta \left(F_\xi^{-1} (x) \right) \right).$$

Lemma is proved. □

Corollary 1. *Lemma's condition can be changed to the existence of F_η^{-1} . Then separating line is*

$$x(y) = \min \left(F_\xi \left(F_\eta^{-1} (1 - y) \right), 1 - F_\xi \left(F_\eta^{-1} (x) \right) \right)$$

Let's consider the general case, when distribution functions don't have the inverse functions in some points. We will use the generalized inverse [3].

Definition 6. *For an increasing function $T : R \rightarrow R$ with $T(-\infty) = \lim_{x \downarrow -\infty} T(x)$ and $T(\infty) = \lim_{x \uparrow \infty} T(x)$, the generalized inverse*

$$T^-(y) = \inf \{ x \in R : T(x) \geq y \}, y \in R \tag{8}$$

with the convention that $\inf \emptyset = \infty$. If $T : R \rightarrow [0, 1]$ is a distribution function, $T^- : [0, 1] \rightarrow \bar{R}$ is also called the quantile function of T .

Lemma 2. *Sets D_l and $\bar{D}_l := (0, 1)^2 \setminus D$ are separated by the line*

$$y(x) = \min \left(F_\eta \left(\left(F_\xi \right)^- (1 - x) \right), 1 - F_\eta \left(\left(F_\xi \right)^- (x) \right) \right) \tag{9}$$

Proof. Let $x \in (0; 1)$ be the point where inverse function F_ξ^{-1} does not exist. Let's use (6) to find the function $y(x)$. According to (5), for any $\delta > 0$

$$y(x) + \delta \in D_l, \quad y(x) - \delta \in \bar{D}_l.$$

So, sets D_l and \bar{D}_l are separated by the line (6).

Lemma is proved. □

Let's consider the set $\overline{D_{A,B}} = \left\{ (\varepsilon_1, \varepsilon_2) \in (0, 1)^2 : \forall N_A^{\varepsilon_1}, N_B^{\varepsilon_2}, \text{conv}N_A^{\varepsilon_1} \cap \text{conv}N_B^{\varepsilon_2} \neq \emptyset \right\}$.

Theorem 2. *Let the following conditions exist:*

1. *The sets A, B of size n_A, n_B are generated by the independent continuous random variables ξ, η .*

2. *The sets $D_{A,B}$ and $\overline{D_{A,B}}$ are separated by the line $y_{A,B}(x)$.*

Then there exists the following equality

$$\lim_{n_A, n_B \rightarrow \infty} y_{A,B}(x) = y(x),$$

where

$$y(x) = \min\left(F_\eta\left(\left(F_\xi\right)_G^{-1}(1-x)\right), 1 - F_\eta\left(\left(F_\xi\right)_G^{-1}(x)\right)\right)$$

Proof. To prove the theorem it is enough to show that the relations

$$F_{n_B}\left(F_{n_A}^{-1}(y)\right) \rightarrow F_\eta\left(F_\xi^{-1}(y)\right), y \in (0,1) \tag{8}$$

and

$$F_{n_B}\left(1 - F_{n_A}^{-1}(y)\right) \rightarrow F_\eta\left(1 - F_\xi^{-1}(y)\right), y \in (0,1) \tag{9}$$

hold.

Let's show that relation (8) holds.

$$\begin{aligned} \sup_{y \in [0,1]} \left| F_{n_B}\left(F_{n_A}^{-1}(y)\right) - F_\eta\left(F_\xi^{-1}(y)\right) \right| &= \sup_{y \in [0,1]} \left| F_{n_B}\left(F_{n_A}^{-1}(y)\right) + F_\eta\left(F_{n_A}^{-1}(y)\right) - F_\eta\left(F_{n_A}^{-1}(y)\right) - F_\eta\left(F_\xi^{-1}(y)\right) \right| \leq \\ &\leq \sup_{y \in [0,1]} \left| F_{n_B}\left(F_{n_A}^{-1}(y)\right) - F_\eta\left(F_{n_A}^{-1}(y)\right) \right| + \sup_{y \in [0,1]} \left| F_\eta\left(F_{n_A}^{-1}(y)\right) - F_\eta\left(F_\xi^{-1}(y)\right) \right| \leq \\ &\leq \sup_{x \in R^1} \left| F_{n_B}(x) - F_\eta(x) \right| + \sup_{y \in [0,1]} \left| F_\eta\left(F_{n_A}^{-1}(y)\right) - F_\eta\left(F_\xi^{-1}(y)\right) \right| \end{aligned}$$

According to the Glivenko-Cantelli theorem [8] the first term is

$$\sup_{x \in R^1} \left| F_{n_B}(x) - F_\eta(x) \right| \rightarrow 0, n_B \rightarrow \infty,$$

Let η have the density of distribution $f_\eta < K$. Then we have for a second term

$$\sup_{y \in [0,1]} \left| F_\eta\left(F_{n_A}^{-1}(y)\right) - F_\eta\left(F_\xi^{-1}(y)\right) \right| \leq K \sup_{y \in [0,1]} \left| F_{n_A}^{-1}(y) - F_\xi^{-1}(y) \right|$$

Let's show that $\sup_{y \in [0,1]} \left| F_{n_A}^{-1}(y) - F_\xi^{-1}(y) \right| \rightarrow 0$. Suppose that $F_\xi(x_0) = y_0 \in (0,1)$ is fixed. Assume that

$\sup_{y \in [0,1]} \left| F_{n_A}^{-1}(y) - F_\xi^{-1}(y) \right| \not\rightarrow 0$, then $\left| F_\xi(x_0) - F_{n_A}(x_0) \right| \not\rightarrow 0$. We arrive at contradiction.

So, relation (8) holds. By analogy, relation (9) also holds. □

Example 2.

1. Normal distribution

D_l and \bar{D}_l are separated by the line which is illustrated in the Fig. 4.

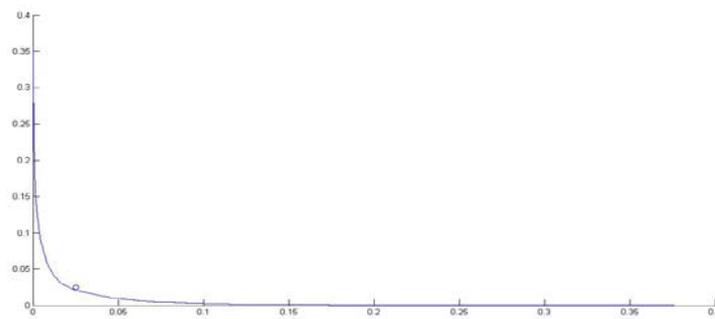


Figure 4. Bound line for the separation space for normal distribution

Since $\frac{|A \cap convB| + |B \cap convA|}{n_A + n_B} = 0.028$, let $\varepsilon_A = 0.025; \varepsilon_B = 0.025$.

According to the Fig. 4, $(0.025; 0.025) \in D_l$

2. Uniform distribution

D_l and \bar{D}_l are separated by the line which is illustrated in the Fig. 5.

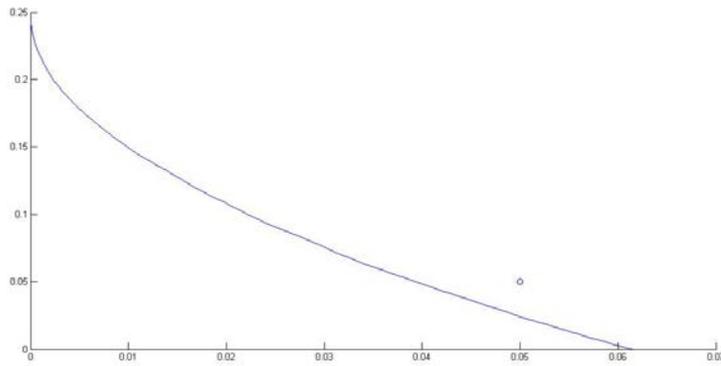


Figure 5. Bound line for the separation space for uniform distribution

Since $\frac{|A \cap convB| + |B \cap convA|}{n_A + n_B} = 0.06$, let $\varepsilon_A = 0.05; \varepsilon_B = 0.05$

According to the Fig. 5, $(0.05; 0.05) \in D_i$

Let's build the ε -net for the set A by the following:

Algorithm of the ε -net building

Let's select $\varepsilon_A, \varepsilon_B \in D_{A,B}$.

Let set A contain the point with minimal y-coordinate. Let's denote the point of set A with minimal x-coordinate by a_{\min} , and the point with maximal x-coordinate by a_{\max} . Let's draw

$k = \left\lceil \frac{1}{\varepsilon_A} \right\rceil + 1$ vertical lines from a_{\min} to a_{\max} in a manner that there are $\varepsilon_A n_A$ points in each of $\left\lceil \frac{1}{\varepsilon_A} \right\rceil$ bands. Vertical lines which separate the bands are described by the equations

$$x = C_i, i = \overline{1, k},$$

where constant C_i can be found from the equation

$$F(C_i) = i\varepsilon_A$$

For each i -th band, $1 \leq i \leq \left\lceil \frac{1}{\varepsilon} \right\rceil$, let's denote:

A^i is the set which contains points from the set A that are contained in the i -th band;

B^i is the set (may be empty) which contains points from the set B that are contained in the i -th band;

ay_{\min}^i, ay_{\max}^i are points from the set A^i with minimal and maximal y-coordinates;

by_{\min}^i, by_{\max}^i are points from the set B^i with minimal and maximal y-coordinates.

$N_A^{\varepsilon_A}$ is the set of points which we will select in the ε -net of the set A .

From the i -th band we will select two points in the set N_A . The first point is the point ay_{\min}^i . According to the assumption, set A is placed below the set B . The second point from the set A^i in the set $N_A^{\varepsilon_A}$ will be selected according to the following rule.

If $B^i = \emptyset$ (it means that i -th band does not contain points from the set B),

add point ay_{\max}^i to the set $N_A^{\varepsilon_A}$

else

if $ay_{\max}^i < by_{\min}^i$ (it means that in i -th band convex hulls of sets A, B are not intersected),

add point ay_{\max}^i to the set $N_A^{\varepsilon_A}$

else (some points of the set B exist in the i -th band and they are placed below some points of set A)

add point $a^i \in A$ to the set $N_A^{\varepsilon_A}$ such that point a^i is the nearest neighbor to the point by_{\min}^i . We will call point a^i the basis point of the set A .

In the same way we build horizontal bands and select two points from each band to the set $N_A^{\varepsilon_A}$.

Example 3.

1. Normal distribution

Vertical and horizontal bands for the set A are illustrated in the Fig. 6

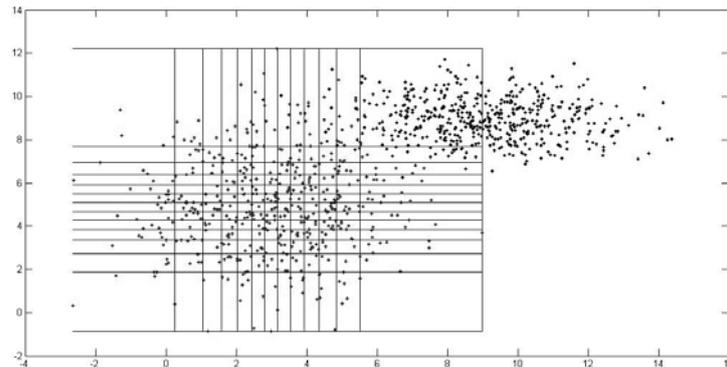


Figure 6. Vertical and horizontal bands for the set A which is generated by the normal distribution

2. Uniform distribution

Vertical and horizontal bands for the set A are illustrated in the fig.7

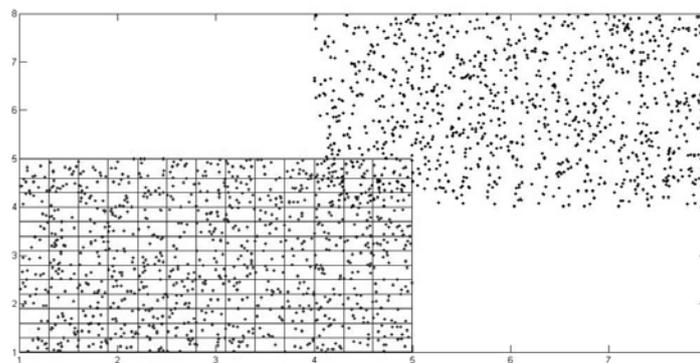


Figure 7. Vertical and horizontal bands for the set A which is generated by the uniform distribution

Lemma 3. The set $N_A^{\varepsilon_A}$ is ε -net for the set A .

Proof. Let's make an indirect proof. Assume $N_A^{\varepsilon_A}$ is not an ε -net of the set A . It means that there exists a halfspace $H \subset R^2$ that contains at least $\varepsilon_A n_A$ sets of point A , but each point does not belong to the set $N_A^{\varepsilon_A}$. Let's denote Z as the set of points from the set A that belong to the halfspace H and $|Z| \geq \varepsilon_A n_A$. Consider a point $z \in Z$. This point belongs to one horizontal and one vertical band. Together with point z one of extreme points of horizontal or vertical band or basis point belongs to the halfspace H . According to the building process, set $N_A^{\varepsilon_A}$ consists of extreme and basis points of the set A , so $Z \cap N_A^{\varepsilon_A} \neq \emptyset$. This contradicts the assumption.

Lemma is proved.

According to the algorithm, ε -net $N_A^{\varepsilon_A}$ consists of $\left\lceil \frac{4}{\varepsilon_A} \right\rceil$ points. ε -net $N_B^{\varepsilon_B}$ is built using the same algorithm.

Example 4.

1. Normal distribution

ε -nets of the sets A, B are illustrated in the Fig. 8

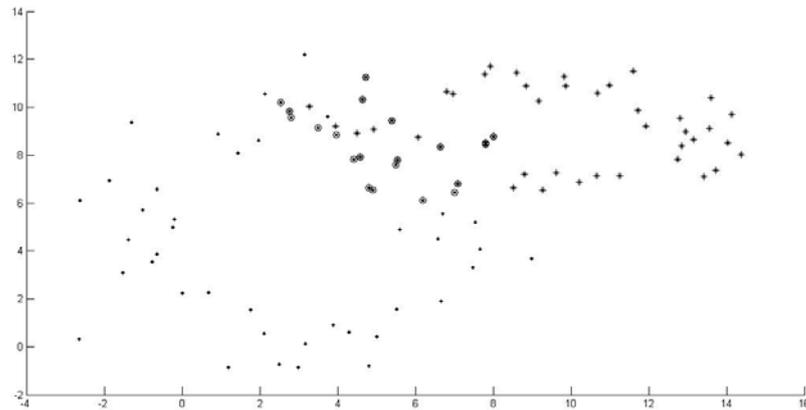


Figure 8. ε -nets of the sets A, B which are generated by the normal distribution.

2. Uniform distribution

ε -nets of the sets A, B are illustrated in the Fig. 9

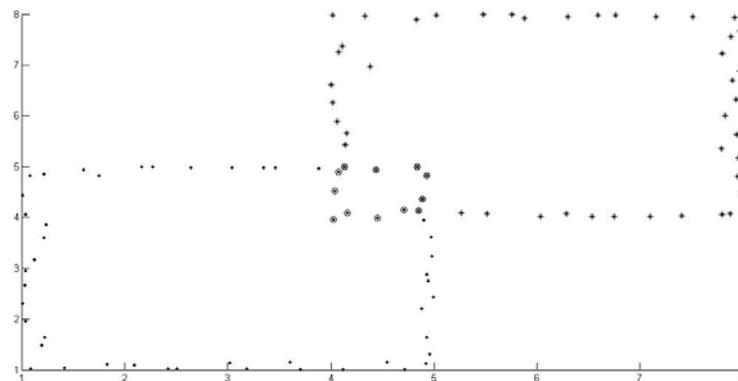


Figure 9. ε -nets of the sets A, B which are generated by the uniform distribution.

3 ε -Nets' Separating

Let's separate sets $N_A^{\varepsilon_A}$ and $N_B^{\varepsilon_B}$ using the separation algorithm of the convex hulls, which is described in [9].

Separation algorithm of the convex hulls

1. Build convex hulls $convN_A^{\varepsilon_A}$ and $convN_B^{\varepsilon_B}$.
2. Find outlier points. In order to minimize the algorithm's time complexity, we find outliers only among the basis points. Point $x \in NB_A$, where $x \in convN_B^{\varepsilon_B}$, is the outlier point of the set $N_A^{\varepsilon_A}$.

The set of outlier points of the set $N_A^{\varepsilon_A}$ is denoted by $P_A^{\varepsilon_A}$.

3. Reject the outliers from the set $N_A^{\varepsilon_A}$ and build the convex hull of the set $N'_A = N_A^{\varepsilon_A} \setminus P_A^{\varepsilon_A}$.
4. Among the edges of the polygons $\text{conv}N'_A$ and $\text{conv}N_B^{\varepsilon_B}$ find the edge so that points of the set N'_A and points of the set $N_B^{\varepsilon_B}$ are placed in different halfspaces which are generated by the line l containing this edge.
5. The line l is the separating line for the sets $N_A^{\varepsilon_A}$ and $N_B^{\varepsilon_B}$. In order to minimize the algorithm's time complexity, we find separating line among the edges containing basis points.

If ε -nets are not linear separable, we propose to use Voronoi diagram [8].

Example 5.

1. Normal distribution

The separating line for ε -nets is illustrated in the Fig. 10

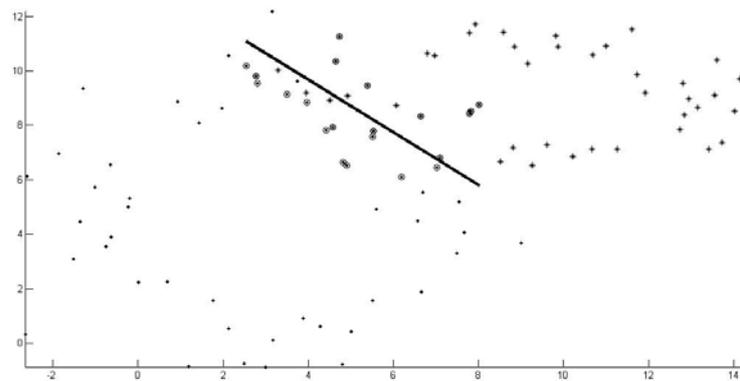


Figure 10. The separating line for ε -nets for normal distribution.

2. Uniform distribution

The separating line for ε -nets is illustrated in the Fig. 11

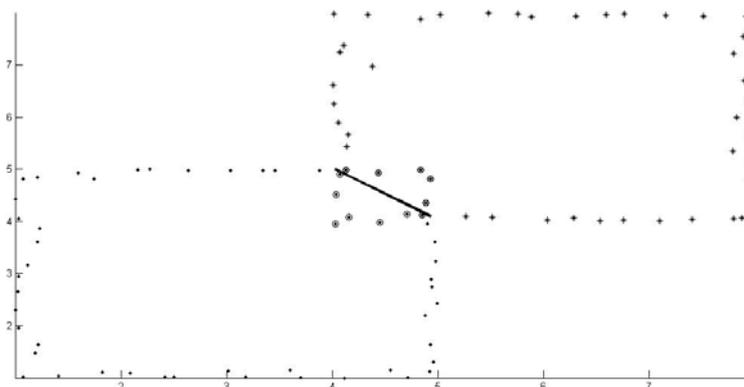


Figure 11. The separating line for ε -nets for uniform distribution

According to the theorem 1, separating line for the ε -nets $N_A^{\varepsilon_A}$ and $N_B^{\varepsilon_B}$ is ε -separating for the sets A and B .

Example 6.

1. Normal distribution

The ε -separating line for sets A and B is illustrated in the Fig. 12

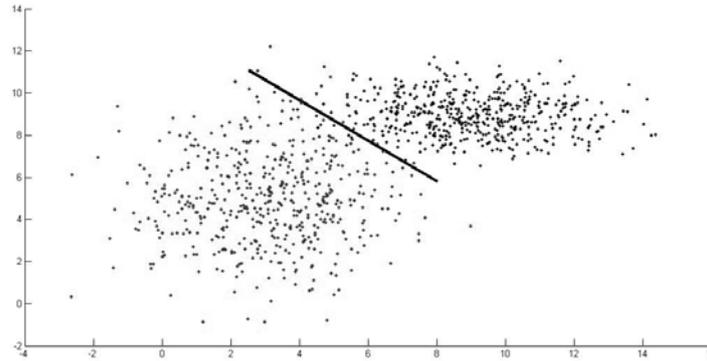


Figure 12. The ε -separating line for the sets A, B which are generated by the normal distribution.

2. Uniform distribution

The ε -separating line for sets A and B is illustrated in the Fig. 13

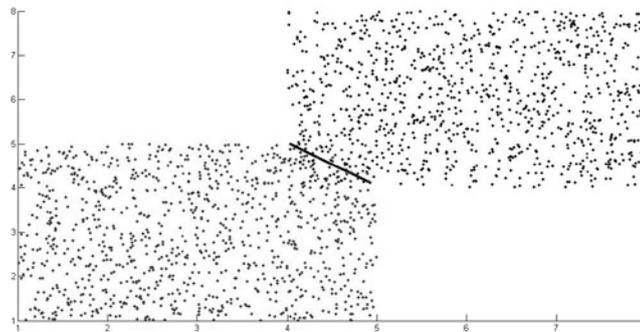


Figure 13. The ε -separating line for the sets A, B which are generated by the uniform distribution.

Let's compare the classification using the algorithm described above and the classification using the Support Vector Machine (SVM) [2].

Example 7.

1. Normal distribution

Classification using ε -nets gives 3.0% errors; classification by SVM 2.9% errors (Fig. 14)

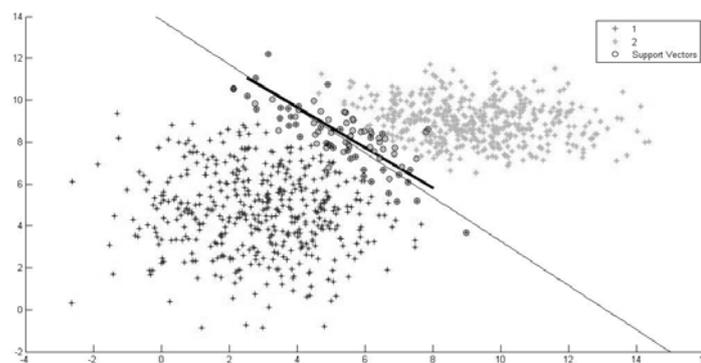


Figure 14. Comparing classification for the normal distribution

2. Uniform distribution

Classification using ε -nets gives 5.9% errors; classification by SVM 5.9% errors (Fig. 15). The separating lines coincide.

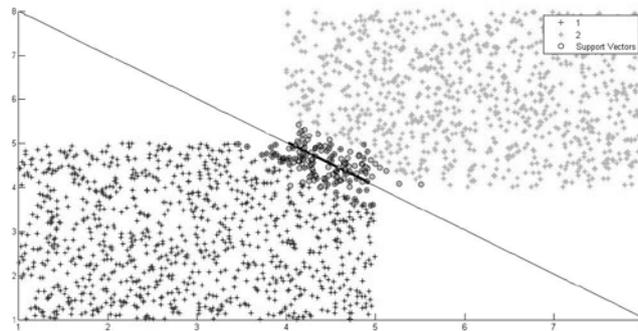


Figure 15. Comparing classification for the uniform distribution

4 Conclusions

The algorithm of building ε -nets for two sets is described in the paper. The ε -nets, constructed according to this algorithm have size $\left\lceil \frac{4}{\varepsilon_A} \right\rceil$, which does not depend on the size of the set. It is shown in the paper that for separating two sets one can use their ε -nets, which considerably reduce the complexity of the separating algorithm for large sets. Two examples in the paper illustrate the algorithm's effectiveness.

References

1. Aronov B., Ezra E., Sharir M. "Small-size epsilon-nets for axis-parallel rectangles and boxes", Symposium on Theory of Computing, 2009, P.P. 639–648
2. Christopher J.C. Burges "A Tutorial on Support Vector Machines for Pattern Recognition", Data Mining and Knowledge Discovery, 2(2), 1998, P.121–167.
3. Embrechts P., Hofert M. A note on generalized inverses Mathematical Methods of Operations Research , 2013, 77(3), 423-432
4. Gärtner B., Hoffmann M. Computational Geometry, <http://www.ti.inf.ethz.ch/ew/lehre/CG12/lecture/CG%20lecture%20notes.pdf>
5. Hausler S. VC Dimension. A Tutorial for the Course Computational Intelligence, <http://www.igi.tugraz.at/lehre/CI>
6. Haussler D. and Welzl E. "Epsilon-nets and simplex range queries", Discrete Comput. Geom., 1987, №2, P.P. 127–151
7. Ivanchuk M. A., Malyk I. V. "Using ε -Nets for Linear Separation of Two Sets in a Euclidean Space R^d ", Cybernetics and Systems Analysis, Vol.51, Issue 6 (2015), P. 965-968.
8. Ivanchuk Maria A., Malyk Igor V. "Building expert medical prognostic systems using Voronoi diagram", Hindawi Publishing Corporation. - International Journal of Computational Mathematics, Volume 2015, Article ID 415146, 4 pages. – DOI 10.1155/2015/415146
9. Ivanchuk Mariya A., Malyk Igor V. "Separation of convex hulls as a way for modeling of systems of prediction of complications in patients", Journal of Automation and Information Sciences, 2015, Vol.47, Issue 4, P.78-84, DOI: 10.1615/JAutomatInfScien.v47.i4.80
10. Kulkarni J., Govindarajan S. "New ε -Net Constructions", Canadian Conference on Computational Geometry (CCCG), 2010, P.P.159-162

11. Matousek J., Seidel R., Welzl E. "How to Net a Lot with Little: Small ε -Nets for Disks and Halfspaces" In Proc. sixth annual symposium on Computational geometry, P.P. 16–22, 1990
12. Tucker H.G. "A Generalization of the Glivenko-Cantelli Theorem", The Annals of Mathematical Statistics, Vol. 30, No. 3, Sep., 1959, pp. 828-830