

Recursive SVM Based on TEDA

Dmitry Kangin^{1(✉)} and Plamen Angelov^{1,2}

¹Data Science Group, School of Computing and Communications,
Lancaster University, Lancaster, UK

²Chair of Excellence, University Carlos III, Madrid, Spain
{d.kangin,p.angelov}@lancaster.ac.uk

Abstract. The new method for incremental learning of SVM model incorporating recently proposed TEDA approach is proposed. The method updates the widely renowned incremental SVM approach, as well as introduces new TEDA and RDE kernels which are learnable and capable of adaptation to data. The slack variables are also adaptive and depend on each point's 'importance' combining the outliers detection with SVM slack variables to deal with misclassifications. Some suggestions on the evolving systems based on SVM are also provided. The examples of image recognition are provided to give a 'proof of concept' for the method.

Keywords: SVM · TEDA · Incremental learning · Evolving system

1 Introduction

Nowadays, there are plenty of models for object classification. Some of them are aimed on off-line classification, which takes all the samples at once, other are aimed for the incremental classification, which work sample-by-sample. Those which do not require to hold all the sample set, are referred as 'online'. Finally, 'evolving' models are designed to change the structure of the system taking into account recent changes and forget those patterns which occurred long ago.

The problem stated in the article is raised by different approaches: SVM models, novel TEDA data analysis concept, and evolving systems.

SVM was first proposed as Generalised Portrait Algorithm by Vapnik and Lerner in 1963 [1], with some works together with Chervonenkis [2], but was not widely used until the beginning of the 1990-s, when new contributions by Vapnik and other authors were proposed [3], [4]. Now, the SVM family contains a huge number of different algorithms, capable of most widely known machine learning problems (clustering[5], outlier detection[6], structured learning[7]).

TEDA approach has recently emerged and gives promising results on data outlier detection[8], classification [9] and other machine learning problems [10]. It is based on data and provides attractive concept of learning the parameters by data, not by some pre-defined constraints.

Evolving systems describe the ability of structural adaptation "from scratch" to consider the changes in the data stream[11], [12]. The distinctive feature of such a system is that it takes into consideration more recent patterns in data, forgetting those

that happened many time ago and hence may be not relevant, by changing the structure of the classifier and forgetting the patterns which are not relevant. One of the popular applications of such systems is classification [13], [14], [15]. Such methods are based on fuzzy systems, as well as neural networks [16].

The outline of the article is as follows:

- first, the SVM systems are reviewed;
- second, the TEDA background and TEDA SVM statement is given;
- third, the novel TEDA kernel is proposed;
- fourth, incremental update procedure is formulated;
- fifth, the new samples incremental update procedure is proposed;
- sixth, the experimental results are discussed.

The article is finalised by conclusion, describing the research results.

2 SVM Model Formulation

Here the accurate statement of the supervised learning problem is performed. Let us have some object set, which we designate as Ω , and a finite space of object classes, referred as Y . Then, let us have a subset, $\Omega_L \subset \Omega$, named training set, for each the following function is defined: $F_L: \Omega_L \rightarrow Y$. The problem is to build a function $F: \Omega \rightarrow Y$, approximating the function F_L on the set Ω_L . The assumption is that F will be a good mapping for further objects, Ω_V where the index V denotes validation. $\Omega_V \cap \Omega_L = \emptyset$; $\Omega_V, \Omega_L \subset \Omega$. The method we propose has its roots in the SVM problem statement [4, 17]. Here we introduce the basic formulation of SVM problem with slack variables and notation we further use. Consider a two class classification problem, where

$$\Omega_L = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}, Y = \{-1, 1\}, F(\mathbf{x}_n) = t_n, n \in [1 \dots k]. \quad (1)$$

Here k is the number of the data samples.

Here we state the problem for overlapping classes (C-SVM) [17]

$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^k \Xi_n \rightarrow \min_{\mathbf{w}, b, \Xi_n}, \quad (2)$$

w.r.t. following constraints:

$$t_n y(\mathbf{x}_n) \geq 1 - \Xi_n, \Xi_n \geq 0, n = 1 \dots k. \quad (3)$$

Here Ξ_n are so-called ‘slack variables’, C is the so-called ‘box constraint’ parameter [17]. This problem statement allows some misclassification of the training set for hardly-divisible data, and $y(x) = \mathbf{w}^T \mathbf{x} + b$, where \mathbf{w} , b are the parameters of the hyperplane. This formulation is widely used [4, 17, 18]. If we write down Lagrangian, differentiate it and find an extremum [17], we get the dual problem

$$\check{L}(\alpha) = -\frac{1}{2} \sum_{n=1}^k \sum_{m=1}^k \alpha_n \alpha_m t_n t_m [\phi(\mathbf{x}_n)]^T \phi(\mathbf{x}_m) + \sum_{n=1}^k \alpha_n \rightarrow \min_{\alpha, b}, \quad (4)$$

$$0 \leq \alpha_n \leq C, \sum_{n=1}^k \alpha_n t_n = 0.$$

Then, we can consider another problem (ν -SVM) [18]:

$$\frac{1}{2} \|\mathbf{w}\|^2 - \nu\gamma + \frac{1}{k} \sum_{n=1}^k \Xi_n \rightarrow \min_{\mathbf{w}, b, \gamma, \Xi_n} \quad (5)$$

w.r.t. the following constraints:

$$t_n \gamma(\mathbf{x}_n) \geq \gamma - \Xi_n, \Xi_n \geq 0, \quad n = 1, \dots, k, \quad (6)$$

The parameter ν gives a lower bound for the fraction of support vectors, as well as an upper bound for the fraction of margin errors in case of $\gamma > 0$. This formulation can be proven to be equivalent to C-SVM with $\hat{C} = \frac{1}{k\nu}$, if $\gamma > 0$ [18].

Then we can denote a kernel function $K(\mathbf{x}_n, \mathbf{x}_m) = [\phi(\mathbf{x}_n)]^T \phi(\mathbf{x}_m)$, where $\phi(\cdot)$ is some mapping function, taken from the (possibly infinite) mapping functions space. It can be recognised as a replacement to the ordinary scalar product, introducing non-linearity into the model. The widely known kernels are: linear, Mahalanobis [19], cosine [20], Gaussian [21], histogram intersection [22] or survival [23] kernels). One can also restore it based on the data using metric learning techniques. Hence, the feature transformation (for finite or infinite space) can be replaced by changing the distance metric, given by the kernel.

3 TEDA Approach Summary

TEDA framework provides a novel systematic method of a “per point” online data analysis [8], [9].

Consider that we have object space $\Omega \subseteq \mathbb{R}^p$ containing data samples, where p is the data dimensionality. This space is equipped with some distance (i.e. Euclidean, L_1 , cosine or any others). Further we refer this distance as $d(\mathbf{x}, \mathbf{y})$. We can pick data samples sequence from this object space:

$$\{\mathbf{x}_1, \mathbf{x}_2 \dots \mathbf{x}_k \dots\}, \mathbf{x}_k \in \Omega, k \in \mathbb{N}, \quad (7)$$

where k is a sequence number of the object which can be represented as the instant of time, when the data sample has arrived. For this reason, index k will be referred to as time instant further for simplicity. We can construct sum distance to some particular point $\mathbf{x} \in \Omega$, for each element up to the k -th one [8], [10]:

$$\pi^k(\mathbf{x}) = \sum_{i=1}^k d(\mathbf{x}, \mathbf{x}_i), k \geq 1. \quad (8)$$

Based on this definition, we define the eccentricity at the time instant k :

$$\xi^k(\mathbf{x}) = \frac{2\pi^k(\mathbf{x})}{\sum_{i=1}^k \pi^k(\mathbf{x}_i)} = 2 \frac{\sum_{i=1}^k d(\mathbf{x}, \mathbf{x}_i)}{\sum_{i=1}^k \sum_{j=1}^k d(\mathbf{x}_i, \mathbf{x}_j)}, k \geq 2, \sum_{i=1}^k \pi^k(\mathbf{x}) > 0. \quad (9)$$

The complement of eccentricity, typicality is defined as:

$$\tau^k(\mathbf{x}) = 1 - \xi^k(\mathbf{x}). \quad (10)$$

The eccentricity and typicality are both bounded [8], [10]:

$$0 \leq \xi^k(\mathbf{x}) \leq 1, \sum_{i=1}^k \xi^k(\mathbf{x}_i) = 2, k \geq 2, \sum_{i=1}^k \pi^k(\mathbf{x}_i) > 0. \quad (11)$$

$$0 \leq \tau^k(\mathbf{x}) \leq 1, \sum_{i=1}^k \tau^k(\mathbf{x}_i) = k - 2, k \geq 2, \sum_{i=1}^k \pi^k(\mathbf{x}_i) > 0. \quad (12)$$

Normalised eccentricity and typicality can also be defined as [8], [10]:

$$\zeta^k(\mathbf{x}) = \frac{\xi^k(\mathbf{x})}{2}, \sum_{i=1}^k \zeta^k(\mathbf{x}_i) = 1, k \geq 2, \sum_{i=1}^k \pi^k(\mathbf{x}_i) > 0. \quad (13)$$

$$t^k(\mathbf{x}) = \frac{\tau^k(\mathbf{x})}{k-2}, \sum_{i=1}^k t^k(\mathbf{x}_i) = 1, k > 2, \sum_{i=1}^k \pi^k(\mathbf{x}_i) > 0. \quad (14)$$

The method's capability of online problems resolution follows from existence of the formulae of the incremental update [8], [9]. There exist convenient formulae for incremental calculation for Euclidean and Mahalanobis distance [19], but here we do not discuss it. Generally, for any distance, we can use the formula $\pi^{k+1}(\mathbf{x}) = \pi^k(\mathbf{x}) + d(\mathbf{x}, \mathbf{x}_{k+1})$ and calculate all other quantities based on it.

4 The TEDA SVM Statement

In the classical formulation of SVM given above we aim to give an upper bound of the fraction of margin errors. But here we are targeting aiming to make it dependent on the 'importance' of each of the support vectors, i.e. make more 'typical' support vectors penalised more for being out of boundaries, rather than 'anomalous'. In other words, we care less for anomalous data samples and more for typical ones to be covered well by the classifier; we do not try to cover all data sample equally, but proportionally to their typicality.

Here, the "local" typicality in regards to each of the classes is proposed as a weight for each of the support vectors. For each class c_m we define typicality [8]

$$\tau_c^k(\mathbf{x}) = 1 - 2 \frac{\sum_{x_j \in X_c} d(\mathbf{x}, \mathbf{x}_j)}{\sum_{x_i \in X_c} \sum_{x_j \in X_c} d(\mathbf{x}_i, \mathbf{x}_j)}, c \in \{c_m, c_{\bar{m}}\}, c_m \in C, c_{\bar{m}} = C \setminus c_m. \quad (15)$$

Here X_c denotes the objects with labels from the set c , which we build here by the scheme 'one versus the rest'. For each of the classes we build the following model:

$$\frac{1}{2} \|\mathbf{w}\|^2 + \sum_{n=1}^k C \tau_{c_m}^k(\mathbf{x}_n) [y_n > 0] \Xi_n + \sum_{n=1}^k C \tau_{c_{\bar{m}}}^k(\mathbf{x}_n) [y_n < 0] \Xi_n \rightarrow \min_{\mathbf{w}, \Xi, b} \quad (16)$$

w.r.t. following constraints:

$$t_n y(\mathbf{x}_n) \geq 1 - \Xi_n, \Xi_n \geq 0, n = 1 \dots k. \quad (17)$$

We can see that this model preserves the upper boundary property γ for the margin errors [1817], where $\gamma = \frac{1}{k\bar{C}}$, $C > 0$, as the upper boundary is known: $C\tau_{c_{\{m,\bar{m}\}}}^k(\mathbf{x}) \leq C$. But what we gain additionally is possibility to control the ‘importance’ of each data samples. The algorithm takes into account our desire to ‘sacrifice’ the correct recognition only for the most ‘anomalous’ data. The more typical is data in the model, the higher will be the box constraint for it because of higher typicality. Opposite case, when the object is more anomalous, is penalised less due to lower typicality in order to avoid fitting the model to anomalous cases.

5 TEDA Kernel

Further to TEDA formulation, we propose also novel formulation of a kernel within the TEDA framework. We define the following TEDA-like kernel:

$$\begin{aligned} \check{\zeta}^k(\mathbf{x}, \mathbf{y}) &= \langle \mathbf{x}, \mathbf{y} \rangle (\zeta^k(\mathbf{x}) \zeta^k(\mathbf{y}))^\gamma, \\ \check{\zeta}^k(\mathbf{x}, \mathbf{y}) &= \langle \mathbf{x}, \mathbf{y} \rangle \left(\frac{\sum_{i=1}^k d(\mathbf{x}_i, \mathbf{x}) \sum_{i=1}^k d(\mathbf{x}_i, \mathbf{y})}{\left(\sum_{i=1}^k \sum_{j=1}^k d(\mathbf{x}_i, \mathbf{x}_j) \right)^2} \right)^\gamma, \sum_{i=1}^k \sum_{j=1}^k d(\mathbf{x}_i, \mathbf{x}_j) > 0. \end{aligned} \quad (18)$$

Here $\zeta^k(\mathbf{x})$ is a normalised data eccentricity, $\gamma > 0$ is some parbamerter, showing the data eccentricity involvement. It helps us to take into account not only scalar product itself, but also the importance of each of the points. The interpretation of the kernel is to increase the kernel values for the ‘anomalous’ points, and at the same time decrease the kernel values between ‘typical’ points, bringing it closer in the data space.

To be the kernel, $\check{\zeta}^k$ should meet the following requirements:

1. $\check{\zeta}^k(\mathbf{x}, \mathbf{y}) = \check{\zeta}^k(\mathbf{y}, \mathbf{x})$.
2. $\check{\zeta}^k$ is positive semi-definite, i.e. $\forall \mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_m \in \Omega, \forall \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k \in \Omega$, where Ω is a Hilbert space, the matrix $M \in \mathbb{R}^{m \times m}$, $M_{ij} = \check{\zeta}^k(\mathbf{y}_i, \mathbf{y}_j)$ is non-negative definite, that is for any $\alpha \in \mathbb{R}^m$ $\alpha^T M \alpha \geq 0$.

Proof sketch:

1.
$$\check{\zeta}^k(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle \left(\frac{\sum_{i=1}^k d(\mathbf{x}_i, \mathbf{x}) \sum_{i=1}^k d(\mathbf{x}_i, \mathbf{y})}{\left(\sum_{i=1}^k \sum_{j=1}^k d(\mathbf{x}_i, \mathbf{x}_j) \right)^2} \right)^\gamma = \langle \mathbf{y}, \mathbf{x} \rangle \left(\frac{\sum_{i=1}^k d(\mathbf{x}_i, \mathbf{y}) \sum_{i=1}^k d(\mathbf{x}_i, \mathbf{x})}{\left(\sum_{i=1}^k \sum_{j=1}^k d(\mathbf{x}_i, \mathbf{x}_j) \right)^2} \right)^\gamma = \check{\zeta}^k(\mathbf{y}, \mathbf{x}).$$

2. For Euclidean distance, it can be proven to be equivalent to:

$$\check{\zeta}^k(\mathbf{x}, \mathbf{y}) \propto \langle \mathbf{x}, \mathbf{y} \rangle \left(\left(\|\mathbf{x} - \boldsymbol{\mu}^k\|^2 + \sigma^{k^2} \right) \left(\|\mathbf{y} - \boldsymbol{\mu}^k\|^2 + \sigma^{k^2} \right) \right)^\gamma, \quad (19)$$

where $\boldsymbol{\mu}^k$ is the mean over $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$, σ^k is a variance over $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$ which depend on the data set and are the values derived from data. It can be proven to be a particular case of the polynomial kernel multiplied on the linear kernel.

The scalar product can be replaced here by any kernel $K(\mathbf{x}, \mathbf{y})$, to get this expression:

$$\zeta^k(\mathbf{x}, \mathbf{y}) \propto K(\mathbf{x}, \mathbf{y}) \left((\|\mathbf{x} - \boldsymbol{\mu}^k\|^2 + \sigma^{k^2}) (\|\mathbf{y} - \boldsymbol{\mu}^k\|^2 + \sigma^{k^2}) \right)^T, \quad (20)$$

We can also make the similar statement via RDE [12]:

$$D(\mathbf{x}, \mathbf{y}) = 1/(1 + \|\mathbf{x} - \mathbf{y}\|^2 + \Sigma_k - \|\boldsymbol{\mu}_k^2\|). \quad (21)$$

It is equivalent to Cauchy kernel [12] (up to the multiplier), and Cauchy kernel itself can be proven to be a proper kernel, hence it is also a kernel:

$$D(\mathbf{x}, \mathbf{y}) \propto \frac{1}{1 + \|\mathbf{x} - \mathbf{y}\|^2/\alpha}, \alpha > 0. \quad (22)$$

6 TEDA SVM Incremental Update

Here the standard problem of the incremental update is stated.

Let us have a training data sequence $\{\mathbf{x}_1, \dots, \mathbf{x}_k \dots\}$, $\mathbf{x}_k \in \Omega$ which arrives one-by-one, where k is an order number of the data element. For each \mathbf{x}_k there is a label $y_k \in Y$. Up to the k -th element we have the problem (16).

Then, the $(k + 1)$ -th element arrives, and the problem is re-formulated as

$$\frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^{k+1} C \tau_{c_m}^i(\mathbf{x}_i) [y_i > 0] \Xi_n + \sum_{i=1}^{k+1} C \tau_{c_m}^i(\mathbf{x}_i) [y_i < 0] \Xi_i \rightarrow \min_{\mathbf{w}, \Xi, b} \quad (23)$$

w.r.t. following constraints:

$$t_i y(\mathbf{x}_i) \geq 1 - \Xi_i, \Xi_i \geq 0, i = 1 \dots k + 1. \quad (24)$$

Also, we should take into account that

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b, \quad (25)$$

as it was mentioned before, and $\phi(\mathbf{x})$ is a feature mapping. Hence, we should also update the feature mapping $\phi(\mathbf{x})$. Generally, it is not feasible as the mapping can be set on functional spaces including infinite-dimensional ones. It can be shown, that dual problem transition allows us to express $y(\mathbf{x})$ in terms of the kernel $K(\mathbf{x}_i, \mathbf{x}_j) = \phi^T(\mathbf{x}_i) \phi(\mathbf{x}_j)$ (this notation denotes scalar product, although the mapping can be given in infinite-dimensional space). Hence we need to update the kernel matrix $K \in \mathbb{R}^{N \times N}$ only, where $K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$.

Then we can see, that the incremental update consists of several stages:

- kernel update (section 6.2);
- solution update for the updated kernel (section 6.2);
- solution update for the updated box constraints (section 6.3);
- solution update incorporating the new data (section 6.1).

Then we can write down the dual problem as

$$\begin{aligned} \check{L}(\alpha) &= \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k \alpha_i \alpha_j t_i t_j [\phi(\mathbf{x}_i)]^T \phi(\mathbf{x}_j) - \sum_{i=1}^k \alpha_i + b \sum_{n=1}^k t_n \alpha_n = \\ &= \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k \alpha_i \alpha_j t_i t_j K(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^k \alpha_i + b \sum_{i=1}^k t_i \alpha_i \rightarrow \min_{\alpha, b}, \end{aligned} \quad (26)$$

$$0 \leq \alpha_i \leq C_i, \forall i \in [1, k], \sum_{i=1}^k \alpha_i t_i = 0.$$

Differentiating the dual problem, we obtain Karush-Kuhn-Tucker conditions (KKT):

$$g_j(\mathbf{x}_j) = \frac{\partial \check{L}(\alpha)}{\partial \alpha_j} = \sum_{i=1}^k \alpha_i t_i t_j K(\mathbf{x}_i, \mathbf{x}_j) - 1 + t_j b = t_j y(\mathbf{x}_j) - 1,$$

$$\frac{\partial \check{L}(\alpha)}{\partial b} = \sum_{i=1}^k t_i \alpha_i = 0, g_j(\mathbf{x}_j) \begin{cases} > 0, & \alpha_j = 0, \\ = 0, & 0 < \alpha_j < C_j. \\ < 0, & \alpha_j = C_j. \end{cases} \quad (27)$$

Let us denote all the training set as Ω_L . This set is divided onto three disjointed sets: margin vectors S , for which $g_j(\mathbf{x}_j) = 0$, error support vectors E , for which $g_j(\mathbf{x}_j) < 0$, and the rest of the vectors vectors R , for which $g_j(\mathbf{x}_j) > 0$, which are not included into the existing solution.

6.1 Adding new Samples

The incremental learning method for SVM was first described in [24].

Apart of this previously proposed method, here we do not choose one global box constraint C for every element of the SVM problem, but use its own box constraint for every object. More, we do not even make it constant, but update it from data.

When transferring from the problem (26) for k elements to the problem for $k + 1$ elements, we should ensure, that the data is in equilibrium, i.e. the conditions of the problem are satisfied.

Let us denote for notation simplicity $M = k + 1$.

We denote as $Q_{ij} = t_i t_j K(\mathbf{x}_i, \mathbf{x}_j)$.

Then, we begin to change the new vector's coefficient α_M until the configuration within the system changes:

$$\begin{aligned} \Delta g_j(\mathbf{x}_j) &= Q_{jM} \Delta \alpha_M + \sum_{n \in S} Q_{jn} \Delta \alpha_n + t_j \Delta b, \forall j \in \Omega_L \cup \{M\}, \\ 0 &= t_M \Delta \alpha_M + \sum_{n \in S} t_n \Delta \alpha_n. \end{aligned} \quad (28)$$

Then, we can define

$$\Theta = \begin{bmatrix} 0 & t_{s_1} & \dots & t_{s_{l(s)}} \\ t_{s_1} & Q_{s_1 s_1} & \dots & Q_{s_1 s_{l(s)}} \\ \vdots & \vdots & \ddots & \vdots \\ t_{s_{l(s)}} & Q_{s_{l(s)} s_1} & \dots & Q_{s_{l(s)} s_{l(s)}} \end{bmatrix} \quad (29)$$

and write the KT conditions changing equations in the vector form as

$$\Theta[\Delta b \quad \Delta\alpha_{s_1} \quad \dots \quad \Delta\alpha_{s_{l(s)}}]^T = -[y_M \quad Q_{s_1 M} \quad \dots \quad Q_{s_{l(s)} M}]^T \Delta\alpha_M \quad (30)$$

for the support vector set.

Then we can continue with

$$\Delta b = \beta \Delta\alpha_M, \quad (31)$$

$$\Delta\alpha_j = \beta_j \Delta\alpha_M, \forall j \in D. \quad (32)$$

$$[\beta \quad \beta_{s_1} \quad \dots \quad \beta_{s_{l(s)}}]^T = -\Theta^{-1}[y_M \quad Q_{s_1 M} \quad \dots \quad Q_{s_{l(s)} M}]^T \quad (33)$$

for all support vectors, and $\beta_n = 0 \forall n \in T \setminus S$.

Then

$$\Delta g_j(x_j) = \Gamma_j \Delta\alpha_M, \forall j \in TU\{M\}; \Gamma_j = Q_{jM} + \sum_{n \in S} Q_{jn} \beta_n + t_j \beta, \forall j \notin S. \quad (34)$$

After that, we find the maximal increment until the following event occur:

- $g_M \leq 0$, with M joining S when $g_M = 0$;
- $\alpha_M \leq 0$, with M joining E when $\alpha_M = 0$;
- $0 \leq \alpha_j \leq C_j, j \in S$ with $\alpha_j = 0$ when j -th vector transfers from S to R , and $\alpha_j = C_j$ when transferring from S to E ;
- $g_j \leq 0, \forall j \in E$, with $g_j = 0$ when j -th vector transfers from E to S ;
- $g_j \geq 0, \forall j \in R$, with $g_j = 0$ when j -th vector transfers from R to S .

After this procedure, the new support vectors should be added to the matrix Θ^{-1} .

It can be proven that

$$\begin{aligned} \Theta^{-1} &\leftarrow \begin{bmatrix} \Theta^{-1} & 0 \\ & \vdots \\ 0 & \dots & 0 \end{bmatrix} + \\ &+ \frac{1}{\Gamma_M} [\beta \quad \beta_{s_1} \quad \dots \quad \beta_{s_{l(s)}} \quad 1]^T [\beta \quad \beta_{s_1} \quad \dots \quad \beta_{s_{l(s)}} \quad 1]. \end{aligned} \quad (35)$$

This formula allows us to add a vector into S .

The procedure is repeated until no transition between subsets R , E , and S occurs. Also, it should be noticed, that the procedure is proven to be reversible [24]. The process of deletion data during the learning process is referred as decremental learning [24].

6.2 Updating the Kernel

Here we address the learnable kernel update problem in SVM:

$$\mathcal{L}_k(\alpha) = \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k \alpha_i \alpha_j t_i t_j K(x_i, x_j) - \sum_{i=1}^k \alpha_i + b \sum_{i=1}^k t_i \alpha_i \rightarrow \min_{\alpha, b}. \quad (36)$$

The problem is replaced by one with the new kernel \hat{K}

$$\check{L}_{k+1}(\alpha) = \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k \alpha_i \alpha_j t_i t_j \widehat{K}(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^k \alpha_i + b \sum_{i=1}^k t_i \alpha_i \rightarrow \min_{\alpha, b} \quad (37)$$

Again, the problems are constrained as in (26).

Therefore, the problem differs in kernel we use. As before, we denote

$$Q_{ij} = t_i t_j K(\mathbf{x}_i, \mathbf{x}_j), \widehat{Q}_{ij} = t_i t_j \widehat{K}(\mathbf{x}_i, \mathbf{x}_j). \quad (38)$$

Let us denote also

$$\widehat{Q}_{ij} - Q_{ij} = \Delta Q_{ij}. \quad (39)$$

Then we consider

$$\Delta g_j(x_j) = \sum_{n \in S} (Q_{jn} + \beta \Delta Q_{jn}) \Delta \alpha_n + \beta \sum_{n \in S} \Delta Q_{jn} \alpha_n + t_j \Delta b, \sum_{n \in S} t_n \Delta \alpha_n = 0. \quad (40)$$

Here $0 \leq \beta \leq 1$. The problem is to express the corrections of $\Delta \alpha_n$ as a function of some coefficient $\beta \in [0, 1]$.

Here we denote Θ as in (29) and

$$\check{\Theta} = \begin{bmatrix} 0 & 0 & \dots & 0 \\ 0 & \Delta Q_{s_1 s_1} & \dots & \Delta Q_{s_1 s_{l_S}} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \Delta Q_{s_{l_S} s_1} & \dots & \Delta Q_{s_{l_S} s_{l_S}} \end{bmatrix}. \quad (41)$$

Then

$$(\Theta + \beta \check{\Theta}) [\Delta b \quad \Delta \alpha_{s_1} \quad \dots \quad \Delta \alpha_{s_{l_S}}]^T = \beta \check{\Theta} [b \quad \alpha_{s_1} \quad \dots \quad \alpha_{s_{l_S}}]^T. \quad (42)$$

In this case

$$\begin{aligned} [\Delta b \quad \Delta \alpha_{s_1} \quad \dots \quad \Delta \alpha_{s_{l_S}}]^T &= \beta (\Theta + \beta \check{\Theta})^{-1} \check{\Theta} [b \quad \alpha_{s_1} \quad \dots \quad \alpha_{s_{l_S}}]^T = \\ &= [b \quad \alpha_{s_1} \quad \dots \quad \alpha_{s_{l_S}}]^T - (\Theta + \beta \check{\Theta})^{-1} \Theta [b \quad \alpha_{s_1} \quad \dots \quad \alpha_{s_{l_S}}]^T. \end{aligned} \quad (43)$$

The maximal increment condition is the same as for the new data update, but additionally we should mind $\beta \in [0, 1]$. The update of the matrix is performed the same way like Θ .

6.3 Updating Box Constraints

In this case, we use the same update equations, as the standard update for the new data, but here we check if the KKT constraints are violated on the first stage due to the change of the box constraints.

We remember, that for each of the objects we have (27) with the old constraints. For the stage $k + 1$ we should have

$$g_j(\mathbf{x}_j) \begin{cases} > 0, & \alpha_j = 0, \\ = 0, & 0 < \alpha_j < C_j^{k+1}, \\ < 0, & \alpha_j = C_j^{k+1}. \end{cases} \quad (44)$$

but for some of the vectors $D \subset \Omega_L$ the conditions can be broken: $\alpha_j > C_j^{k+1}$. Also, because generally $C_j^{k+1} \neq C_j^k$, some transitions between sets E and R may occur. Hence, we need to correct the solution w.r.t. all the vectors violating the conditions.

The procedure is proposed in a following way. For each vector \mathbf{x}_m , violating the KKT conditions, we perform the procedure exactly the same way as for the new vector, considering that it is a newly-added vector to k -vector SVM problem with training set $D \setminus \mathbf{x}_m$ and new C_i . For the sets E and S , the vector should be preliminarily deleted by decremental learning procedure should be as [24].

7 Suggestions on the Method Implementation

Additionally, even despite the given method is not ‘*evolving*’ but ‘*online*’, as it does not discard the previous data samples, the decremental learning procedure can be proposed, as the incremental procedure for SVM is reversible. The following simplest method can be considered:

- for each sample \mathbf{x} , update its within-class typicality $\tau_c^k(\mathbf{x})$;
- inspect all the samples which appeared in the train sequence up to the newest element \mathbf{x}_k : $\{\mathbf{x}_1, \dots, \mathbf{x}_n \dots \mathbf{x}_k\}$, $\mathbf{x}_n \in \Omega$. If their typicality is lower than some given threshold T , we can state, that the sample should be rejected.
- For the sample set elements to be rejected, perform decremental learning, analogous to that described in [24].

Using this modification, we can employ all the benefits of evolving systems, such as possibility to adopt to the changes on the data neglecting those patterns that appeared long ago, but base it on the SVM framework, rather than fuzzy systems or neural networks. Usage of the SVM framework, instead of fuzzy rule-based evolving systems, gives us an opportunity of using strict optimisation problem statements, which gives benefits of better generalisation of the method.

8 Demonstrations for the Experimental Data

For the proof of concept, an example of human activities images classification was provided (**Fig. 1**, data set [25]). The feature transformation method was exactly borrowed from the previous research described in [13] and is composed from Haar and gist [26] features.



Fig. 1. Human activities data samples

Riding bike	60%	0%	15%	0%	25%	Riding bike	65%	5%	15%	0%	15%	
Playing guitar	10%	75%	10%	5%	0%	Playing guitar	0%	85%	5%	5%	5%	
Riding horse	5%	5%	85%	0%	5%	Riding horse	5%	10%	80%	0%	5%	
Phoning	5%	5%	0%	90%	0%	Phoning	0%	5%	5%	90%	0%	
Running	10%	0%	5%	0%	85%	Running	5%	0%	5%	5%	85%	
		Riding bike	Playing guitar	Riding horse	Phoning	Running		Riding bike	Playing guitar	Riding horse	Phoning	Running

Fig. 2. Results of the recognition for different methods (left picture is SVM with histogram intersection kernel, accuracy rate is 79%, right picture is SVM with TEDA kernel, combined with histogram intersection kernel, and TEDA box constraints, accuracy rate is 81%).

For the experiments, the following methods were compared:

- SVM with Gaussian kernel with $\sigma = 34.4725$, $C = 30$.
- SVM with TEDA kernel, combined with Gaussian kernel, with $\sigma = 34.4725$, $\gamma = 2$, $C = 30$, and TEDA weights.

Here the Gaussian kernel is

$$K_G(x, y) = \exp(-(\mathbf{x} - \mathbf{y})^T(\mathbf{x} - \mathbf{y})/(2\sigma^2)). \tag{45}$$

In contrast to Gaussian distribution probability density, the multiplier is neglected, just because the solution will be the same for any constant.

By TEDA kernel combined with Gaussian kernel we denote

$$K_{TEDA}(x, y) = K_G(x, y) \times \left((\|\mathbf{x} - \boldsymbol{\mu}^k\|^2 + \sigma^{k^2}) (\|\mathbf{y} - \boldsymbol{\mu}^k\|^2 + \sigma^{k^2}) \right)^\gamma. \tag{46}$$

No other modifications are introduced into the model except the TEDA kernel and TEDA weights. Therefore, we suppose that the improvement in the kernel and box variables leads to the increase of the quality of the recognition.

The training data set contains 200 items, 40 images for each of the classes. The testing data set consists of 100 items, 20 images per class. Although the method was initially designed for handwritten symbol images, it can be also applied to describe any other images as well.

The results are given as it is depicted in **Fig. 2**. On the left part of the image, the results for Gaussian kernel are proposed. At the right part of the image, SVM was augmented with TEDA ‘box constraints’ and uses TEDA kernel, combined with Gaussian kernel. One can see, that the results were (slightly) improved. However, the method has more perspectives, as it makes the model more flexible giving different values to the ‘box constraints’ and enabling kernel learning ‘from scratch’.

9 Conclusion

In this paper, the new modification of the SVM with slack variables for classification was proposed which changes slack variables independently for each of the data samples taking into account the ‘typicality’ of each data sample and forcing more ‘anomalous’ data to be misclassified sooner than typical ones. For this purpose, the novel ‘box constraints’ update model was developed based on the recently proposed TEDA framework, and the novel learnable kernels based on TEDA and RDE were proposed. The incremental SVM was modified to take into account changes of the box variables during the training process, as well as learnable kernels. Though the model, proposed here, is not evolving yet, the ideas were described how to make the model ‘evolving’. In the experimental section, a few examples were given to approve the ideas of SVM with new learnable box constraints.

Acknowledgements. Plamen Angelov would like to acknowledge the support of Chair of Excellence award which he was given by Santander Bank, Spain for the period March-September 2015 which is being hosted by Carlos III University, Madrid, Spain.

References

1. Vapnik, V., Lerner, A.: Pattern recognition using generalised portrait method. *Automation and Remote Control* **24**, 774–780 (1963)
2. Vapnik, V.N., Ya, A., Chervonenkis.: *Teoriya raspoznavaniya obrazov: Statisticheskie problemy obuchenija*. (Russian) = Theory of pattern recognition: Statistical problems of learning. Nauka, Moscow (1974)
3. Boser, B.E., Guyon, I.M., Vapnik, V.N.: A training algorithm for optimal margin classifiers. In: *COLT 1992: Proceedings of the Fifth Annual Workshop on Computational Learning Theory*, pp. 144–152. ACM Press, New York (1992)
4. Vapnik, V.N.: *The Nature of Statistical Learning Theory*. Springer-Verlag New York, Inc. (1995)
5. Yang, J., Estivill-Castro, V., Chalup, S.K.: Support vector clustering through proximity graph modelling. In: *Proceedings of the 9th International Conference on Neural Information Processing, ICONIP 2002*, vol. 2. IEEE (2002)
6. Chen, Y., Zhou, X., Huang, T.S.: One-class SVM for learning in image retrieval. In: *Proceedings of the International Conference on Image Processing*, vol. 1. IEEE (2001)
7. Tsochantaridis, I., et al.: Support vector machine learning for interdependent and structured output spaces. In: *Proceedings of the Twenty-First International Conference on Machine Learning*. ACM (2004)

8. Angelov, P.: Outside the Box: An Alternative Data Analytics Framework. *Journal of Automation, Mobile Robotics and Intelligent Systems* **8**(2), 53–59 (2014)
9. Kangin, D., Angelov, P.: New Autonomously Evolving Classifier TEDAClass. *IEEE Transactions on Cybernetics* (2014, submitted)
10. Angelov, P.: Anomaly detection based on eccentricity analysis. In: 2014 IEEE Symposium Series on Computational Intelligence, 9-12 December, Orlando, Florida (2014, to appear, accepted to publication)
11. Angelov, P., Filev, D., Kasabov, N. (eds.) *Evolving Intelligent Systems: Methodology and Applications*. IEEE Press Series on Computational Intelligence, p. 444. John Willey and Sons (April 2010) ISBN: 978-0-470-28719-4
12. Angelov, P.: *Autonomous Learning Systems: From Data Streams to Knowledge in Real time*. John Willey and Sons (December 2012) ISBN: 978-1-1199-5152-0
13. Angelov, P., Kangin, D., Zhou, X., Kolev, D.: Symbol recognition with a new autonomously evolving classifier autotool. In 2014 IEEE Conference on Evolving and Adaptive Intelligent Systems, pp. 1–6. IEEE (2014); DOI BibTeX
14. Angelov, P., Yager, R.: A new type of simplified fuzzy rule-based systems. *International Journal of General Systems* **41**(2), 163–185 (2012)
15. Angelov, P., Zhou, X., Klawonn, F.: Evolving fuzzy rule-based classifiers. In: First 2007 IEEE International Conference on Computational Intelligence Applications for Signal and Image Processing, April 1-5, Honolulu, Hawaii, USA, pp. 220–225 (2007)
16. Kasabov, N., Liang, L., Krishnamurthi, R., Feigin, V., Othman, M., Hou, Z., Parmar, P.: Evolving Spiking Neural Networks for Personalised Modelling of Spatio-Temporal Data and Early Prediction of Events: A Case Study on Stroke. *Neurocomputing* **134**, 269–279 (2014)
17. Bishop, C.M.: *Pattern recognition and machine learning*, pp. 325–358. Springer, New York (2006)
18. Schölkopf, B., Smola, A., Williamson, R.C., Bartlett, P.L.: New support vector algorithms. *Neural Computation* **12**, 1207–1245 (2000)
19. Mahalanobis, P.C.: On the generalised distance in statistics. *Proceedings of the National Institute of Sciences of India* **2**(1), 49–55 (1936)
20. Li, B., Han, L.: Distance weighted cosine similarity measure for text classification. In: Yin, H., Tang, K., Gao, Y., Klawonn, F., Lee, M., Weise, T., Li, B., Yao, X. (eds.) *IDEAL 2013*. LNCS, vol. 8206, pp. 611–618. Springer, Heidelberg (2013)
21. Keerthi, S.S., Lin, C.-J.: Asymptotic behaviors of support vector machines with Gaussian kernel. *Neural Computation* **15**(7), 1667–1689 (2003)
22. Maji, S., Berg, A.C., Malik, J.: Classification using intersection kernel support vector machines is efficient. In: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2008*. IEEE (2008)
23. Chen, B., Zheng, N., Principe, J.C.: Survival kernel with application to kernel adaptive filtering. In: *The 2013 International Joint Conference on IEEE Neural Networks (IJCNN)* (2013)
24. Cauwenberghs, G., Poggio, T.: Incremental and Decremental Support Vector Machine Learning. *NIPS*, pp. 409–415 (2000)
25. Li, L.-J., Fei-Fei, L.: What, where and who? Classifying events by scene and object recognition. In: *IEEE 11th International Conference on Computer Vision, ICCV 2007*, pp. 1–8 (2007)
26. Oliva, A., Torralba, A.: Modeling the shape of the scene: a holistic representation of the spatial envelope. *International Journal of Computer Vision* **42**(3), 145–175 (2001). Gist Descriptor (Matlab code). <http://people.csail.mit.edu/torralba/code/spatialenvelope/b>