# A novel semisupervised SVM classifier based on active learning and context information

Fei Gao · Wenchao Lv · Yaotian Zhang · Jinping Sun · Jun Wang · Erfu Yang

**Abstract** This paper proposes a novel semisupervised support vector machine classifier ($S^3VM$) based on active learning (AL) and context information to solve the problem where the number of labeled samples is insufficient. Firstly, a new semisupervised learning (SSL) method is designed using AL to select unlabeled samples as the semilabled samples, then the context information is exploited to further expand the selected samples and relabel them, along with the labeled samples train $S^3VM$ classifier. Next, a new query function is designed to enhance the reliability of the classification results by using the Euclidean distance between the samples. Finally, in order to enhance the robustness of the proposed algorithm, a fusion method is designed. Several experiments on change detection are performed by considering some real remote sensing images. The results show that the proposed algorithm in comparison with other algorithms can significantly improve the detection accuracy and achieve a fast convergence in addition to verify the effectiveness of the fusion method developed in this paper.

**Keywords** semisupervised support vector machine · active learning · context information · remote sensing

Fei Gao • Wenchao Lv • Yaotian Zhang • Jinping Sun • Jun Wang

School of Electronic and Information Engineering, Beihang University, Beijing, China

Fei Gao

e-mail: feigao2000@163.com

Wenchao Lv

e-mail: 13161111047@163.com

Yaotian Zhang(✉)

e-mail: zhangyaotian@buaa.edu.cn

Jinping Sun

e-mail: sunjinping@buaa.edu.cn

Jun Wang

E-mail: wangj203@buaa.edu.cn

Erfu Yang

Space Mechatronic Systems Technology Laboratory, University of Strathclyde, Glasgow G1 1XJ, UK

e-mail: erfu.yang@strath.ac.uk

**1 Introduction**

With the rapid development of remote sensing (RS) technology, RS data has been expanded constantly and used widely in geologic monitoring, environmental protection and disaster relief, where the distinction of the land-cover change(Hansen and Loveland 2012), classification of the land-use(Gong et al. 2015) and assessment of the earthquake influence(Geiß and Taubenböck 2013)can be attributed to the problem of object detection under certain conditions. In these fields, it is a very difficult task to detect useful objects from a large amount of RS images (RSIs), which presents a high requirement to the information extraction ability of the correlated algorithms and increasingly attracts the research interest.

The traditional technologies for object detection are mainly based on unsupervised learning and supervised learning. Unsupervised method detects the object from its background via grouping the samples into different clusters and does not require prior knowledge, which reduces the time and cost for human-labeled training samples. However, such methods only consider the characteristic information between the different clusters and lack the effective guidance of the labeled samples, which makes it impossible to extract objects of interest(Anand et al. 2014). Therefore, the supervision method based on the labeled samples is adopted to change the detection problem into classification problem and quickly captures the effective RS information with the high accuracy. However, to train the classifier, the supervised method requires enough labeled samples(Ujjwal Maulik and Chakraborty 2012), which are scarce in practice and easily affected by noises. Moreover, it is very difficult to obtain the labeled samples manually because of the limit of the time and cost (Li et al. 2010). In contrast, the unlabeled samples are abundant and contain a wealth of information. By using them, not only the problem of insufficient training samples can be solved, but also the working efficiency is improved(Shahshahani and Landgrebe 1994). Therefore, in recent years the algorithms based on the semisupervised learning (SSL) have been concerned greatly.

The SSL based method iteratively selects the unlabeled samples with the labeled samples together to train the classifier(Kawakita and Kanamori 2013), and can reduce human intervention and classify the unlabeled samples more effectively. In general, the SSL can be grouped into the following five categories, i.e.: self-training, cooperative training, generative probability model, semisupervised support vector machine($S^3VM$) and graph based methods(Chapelle et al. 2006; Zhu 2010), where the key to their effective functioning is the way for screening the unlabeled samples. The progressive $S^3VM$ with diversity (PS$^3$VM-D) takes $\gamma$ samples which are within and closer to the margin band, to define the   candidate semilabeled samples, then incrementally selects $\rho$ diverse samples among the $\gamma$ candidates by considering the kernel cosine-angular similarity(Persello and Bruzzone 2014). The context-sensitive $S^3VM$ (CS$^4$VM) selects the context patterns of training samples as the semilabeled samples, then weights them depending on their consistency with the center sample by SVM(Bruzzone and Persello 2009).Junwei et al. (2015a)built a high-level object feature representation using a deep Boltzmann machine(DBM) and applies the proposed Bayesian principle to characterize the objects information, then classifies the objects by the liner SVM to select the semilabeled samples.

The performance of the SSL algorithms can still be improved in the following aspects. First, the SSL based methods can be deteriorated when exploiting huge amount of the unlabeled samples, which

would seriously hinder the adjustment of the relevant parameters(Munoz-Mari et al. 2012; Gomez-Chova et al. 2008), and the effective strategy for the selection of the semilabeled samples is necessary. The SSL based methods can be greatly affected by initial training set(Didaci et al. 2012), where the validity of samples can't be guaranteed sometimes, leading to random classification result.

For the selection strategy of the semilabeled samples, Pasolli et al. (2014) applied the active learning (AL) and the spatial information, including Parzen window, spatial entropy, etc., to process the unlabeled samples, where the semilabeled samples were selected by estimating the probability density distribution function of the random variables and the entropy of the discrete random variables. Demir et al. (2011) applied the AL based on kernel clustering to deal with the unlabeled samples, then designed a new strategy to select the most representative semilabeled samples. Inspired by the above algorithms, the selection strategy was designed utilizing AL in the proposed SSL based algorithm. The AL iteratively extracts the most informative unlabeled samples to be labeled manually and can obtain a higher classification accuracy than other SSL algorithms, reducing the dependence on the labeled samples (Tuia et al. 2011; Persello and Bruzzone 2014).However, the existing feature descriptors are insufficiently powerful to characterize the information of objects, which inevitably leads samples mislabeling and seriously reduces the classification accuracy. To fully describe the sample characteristics, Junwei et al. (2015b)proposed a novel framework via the deep learning methods for the salient object detection, which extracts four neighborhood image boundaries of the object as a reference, then applied the stacked de-noising auto-encoder(SDAE) with the deep learning architectures to model image background. The robustness of the mislabeled patterns is improved by exploiting spatial information. Inspired by the above method, the structure of the AL is modified by using the context information, which can expand the amount of semilabeled samples and improve the accuracy to a certain extent when labeling them.

The reason why the $S^3VM$ is chosen to process the selected samples is detailed in the following. The support vectors of $S^3VM$ make the decision rather than many redundant samples to improve the operational efficiency and the robustness to the polluted samples. The nonlinear mapping based $S^3VM$ can map the samples into a high dimensional characteristic space to make them linearly separable. The functional margin that is the distance between the sample and the optimal hyper plane decided by the $S^3VM$, can be used as the confidence-based decision marking to make the $S^3VM$ be effectively exploited by the SSL based methods(Hearst et al. 1998; U. Maulik and Chakraborty 2014).

For the uncertainty of the initial training set, the fusion is an effective method. The downsampling is commonly used to obtain different data models by generating the images at different resolutions(Bazi et al. 2010). However, the progressively downsampling is not suitable to deal with the image with large difference in areas, because it can lose them with scanty pixels. The initial sets in the SSL does not require many labeled samples, so they can be directly selected from the image to reduce the complexity of the preprocessing. Their SSL results as different models were fused to decrease the influence of the random distribution of the initial set and obtain stable classification result.

In this paper, a novel $S^3VM$ algorithm is proposed to address the classification problems of the RSIs when the available labeled samples are not sufficient. To guarantee the decision-making ability of the proposed method, the AL and context information are applied to select the informative unlabeled

1 samples (semilabeled samples), which are combined with the labeled samples together to train the

2 S³VM. The novelty of this paper lies in the following aspects. First, a new query function is designed

3 based on the Euclidean distance among samples. Second, the amount of the semilabeled samples are

4 increased by exploiting their contextual patterns. Last, the selected samples are relabeled by exploiting

5 the correlation between the sample and context. Hence, the comprehensive utilization of the context

6 information ensures the reliability of the selected samples and increases the classification credibility of

7 the S³VM.

8 The rest of this paper is organized as follows. Section 2 describes the proposed SSL based method in

9 detail, and the necessity of the fusion is also analyzed to design new fusion rules. Section 3 presents the

10 three data sets and experiments for the RSIs change detection. Section 4 draws the conclusion of this

11 paper.

12 **2 Proposed method**

13 The algorithm consists of two parts: the proposed S³VM based on AL and context information, and the

14 fusion method.

15 2.1 Proposed SSL-based method

16 Assuming that the training set **L** is composed of $n$ pairs $(x_i, y_i)$ available, where $x_i$ denotes the training

17 samples and $y_i \in \{+1, -1\}$ denotes the label for the binary classifier. We assume that a learning set **U**

18 with $m$ unlabeled samples $(x'_j)$. They are defined as

19
$$\mathbf{L} = \left\{ (x_i, y_i) \mid i = 1 \ldots \ldots n \right\} \tag{1}$$

20
$$\mathbf{U} = \{ x'_j \mid j = 1 \ldots \ldots m \} \tag{2}$$

21 In order to expand **L**, the semilabeled samples are selected from **U** iteratively. Fig. 3 shows the
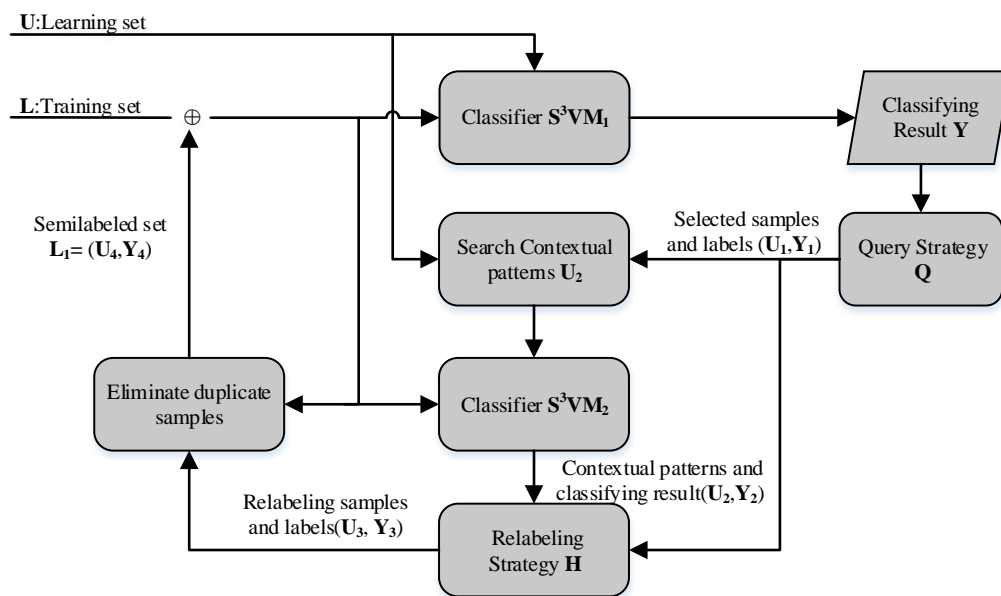
22 flowchart of the SSL method based on the S³VM.



23

24 **Fig. 1** Flowchart of the proposed method

1    First, **L** is processed by the S³VM₁,during which the label **Y** of the **U** is obtained. The subset (**U₁**, **Y₁**) is

2    selected by the query function Q. Then, the context patterns **U₂**of the **U₁** is found from the **U** and

3    classified by the S³VM₂.**U₂**along with the central pattern (**U₁**,**Y₁**) are relabeled by H to get the

4    semilabeled set (**U₃**,**Y₃**). Last, the repeated samples between **L**and (**U₃**,**Y₃**) are eliminated to obtain the

5    ultimate semilabeled set **L₁**= (**U₄**, **Y₄**)to be added into **L**. The entire process is iterated until the

6    predefined convergence condition is satisfied, e.g., the total number of the semilabeled samples or the

7    classification accuracy is reached (Pasolli et al. 2014).

8    Next, the implementation of the S³VM, the query function Q, the search contextual patterns and

9    relabeling strategy H are described.

10    2.2S³VM method

11    The S³VM is the expansion of the SVM. A standard SVM is based on the structural risk minimization

12    to classify the learning set by extracting the support vectors from the training set to find the optimal

13    hyper plane(Scholkopf et al. 1997). In case of the binary SVM, given the training set **L**and the testing

14    set **U**, it is limited to the following constrained optimization problem (Izquierdo-Verdiguier et al.

15    2013):

16

$$\min \Psi(\boldsymbol{w}) = \frac{1}{2}(\boldsymbol{w}^T\boldsymbol{w}) + C\sum_{i=1}^{n}\xi_i$$
$$\text{Subject to}: \quad y_i[\boldsymbol{w}^T\Phi(\boldsymbol{x}_i)+b] \geq 1-\xi_i \quad i=1,\ldots,n$$
$$\xi_i \geq 0$$

(3)

17    Where $x_i$ is the training sample and $y_i$ is the corresponding label,$(x_i,y_i) \in$ **L**;$\Phi(\cdot)$ maps the data into the

18    feature space; $\boldsymbol{w}$ is the orthogonal vector between $x_i$ and the hyper plane; $b$ is the bias to measure

19    the distance between **L** and the hyper plane; $\xi_i$ is the slack variable to represent offset of $x_i$; $C$is the cost

20    factor to measure the weight between the optimal hyper plane and the minimum offset; $n$ is the number

21    of the training samples.

22    After the initialization, the iterative process is operated and the semilabeled samples (selected from **U**

23    in the previous step) are added to **L**. Their confidence is diverse in different iterative steps and they are

24    given different cost factors(Bovolo et al. 2008), which leads to the following cost function for the

25    classifier learning :

26

$$\min \Psi(\boldsymbol{w}) = \frac{1}{2}(\boldsymbol{w}^T\boldsymbol{w}) + C\sum_{i=1}^{n}\xi_i + \sum_{j=1}^{m}C_j \square \varepsilon_j$$
$$\text{Subject to}: \quad y_i[\boldsymbol{w}^T\Phi(x_i)+b] \geq 1-\xi_i \quad i=1,\ldots,n$$
$$\hat{y}_j[\boldsymbol{w}^T\Phi(\hat{x}_j)+b] \geq 1-\varepsilon_j \quad j=1,\ldots,n'$$
$$\xi_i \geq 0, \varepsilon_j \geq 0$$

(4)

27    where$\hat{x}_j$is the semilabled sample selected from **U**, with the slack variable($\varepsilon_j$),cost factor ($C_j$) and

28    semilabel($\hat{y}_j$), and $n'$is the number of the semilabeled samples.

29    Applying the Lagrange Duality, the equation (4) can be transformed into the dual problem, which can

30    be solved by applying the Karush–Kuhn–Tucker(KKT)(Yi et al. 2011)optimality conditions:

$$\max_{\alpha,\beta}\left(\sum_{i=1}^{n}\alpha_i+\sum_{j=1}^{m}\beta_j-\frac{1}{2}\sum_{i,j=1}^{n}\alpha_i\alpha_j y_i y_j K(x_i,x_j)-\frac{1}{2}\sum_{i,j=1}^{m}\beta_i\beta_j\hat{y}_i\hat{y}_j K(\hat{x}_i,\hat{x}_j)-\sum_{i=1}^{n}\sum_{j=1}^{m}\alpha_i\beta_j y_i\hat{y}_j K(x_i,\hat{x}_j)\right)$$

$$\text{Subject to:}\quad \sum_{i=1}^{n}\alpha_i y_i+\sum_{j=1}^{m}\beta_j\hat{y}_j=0 \tag{5}$$

$$0\le\alpha_i\le C \qquad i=1,\ldots,n$$
$$0\le\beta_j\le C_j \qquad j=1,\ldots,m$$

where $K=(x_i,\hat{x}_{jj})$ is a kernel function to calculates the inner product $<\Phi(x_i)\cdot\Phi(\hat{x}_j)>$; $\alpha_i$ and $\beta_j$ are the Lagrange multipliers corresponding to the labeled sample $(x_i)$ and the semilabeled sample $(\hat{x}_j)$ respectively. When $\alpha_i$ and $\beta_j$ are solved by(5), $b$ is solved by

$$f(x)=w^T\bullet\phi(x)+b=\left(\sum_{i=1}^{n}a_i y_i K(x,x_i)+\sum_{j=1}^{m}\beta_j\hat{y}_j K(x,\hat{x}_j)\right)+b \tag{6}$$

Finally, the unlabeled sample $x'(x'\in U)$ can be labeled by the following decision function:

$$y'=\text{sgn}\{(w^T\bullet x')+b\}=\text{sgn}\left\{\left(\sum_{i=1}^{n}a_i y_i K(x',x_i)+\sum_{j=1}^{m}\beta_j\hat{y}_j K(x',\hat{x}_j)\right)+b\right\} \tag{7}$$

It is worth noting that with the increase of iteration, the previous semilabeled samples are more convincing and the corresponding $C_j$ would get bigger. The S³VM can deal with the nonlinear problem using many different kernel functions, such as Gauss kernel function, radial kernel function or exponential kernel function.

The basic structure of the S³VM is described above. The learning set $\mathbf{U}$ and the result $\mathbf{Y}$ classified by the S³VM₁ are the inputs to Q that is described next.

2.3 Query Strategy Q

Fig. 2 shows the structure of the query strategy Q consisting of two parts Q₁ and Q₂. Q₁is used to screen $(\mathbf{U'}, \mathbf{Y'})$ from $(\mathbf{U}, \mathbf{Y})$.Then, $(\mathbf{U'}, \mathbf{Y'})$ is refined by Q₂using the proposed rules $f_1$ and $f_2$, and the output is$(\mathbf{U_1}, \mathbf{Y_1})$.
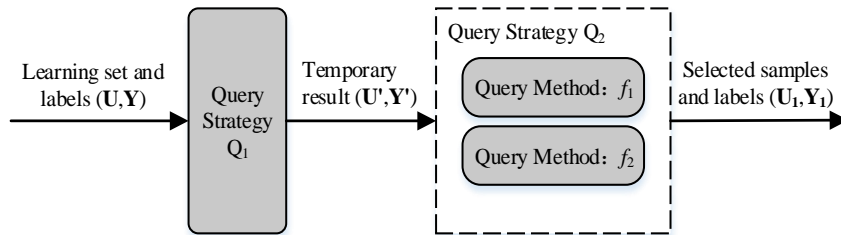


**Fig. 2**Structure of Q

2.3.1 The implementation of Q₁ and Q₂

Q₁is designed on the basis of the MS. The MS is conducive to the convergence of the algorithm because the samples within the margin band are more informative than others. Therefore, $\mathbf{U'}$is selected within the margin band, where the absolute value of the decision function is limited between the $\rho$ and 1, defined as follows:

$$\mathbf{U}' = \{x_i \mid x_i \in \mathbf{U}, \rho < abs(value_i) < 1, i = 1, \ldots, m\}$$
$$\text{Subject to}: \quad value_i = f(x_i) = \mathbf{w}^T \Box x_i + b \tag{8}$$

2    $Q_2$is designed by using the Euclidean distance between the samples and the screening samples from the

3    $\mathbf{U}'$by applying the rules $f_1$and $f_2$. The implementation is shown as follows:

4    I. The distance ($\mathbf{D}$) between the sample of $\mathbf{U}'$and the training set $\mathbf{L}$ is calculated by

$$\mathbf{D} = \{d_{ij} \mid d_{ij} = \text{norm}(x_i, x_j), x_i \in \mathbf{U}', x_j \in \mathbf{L}\} \tag{9}$$

6    II. The minimum, maximum and their ratios between $\mathbf{U}'$and $\mathbf{L}$ are calculated by

$$\mathbf{D}_{\min} = \left\{d_i \mid d_i = \min_j(d_{ij}), d_{ij} \in \mathbf{D}\right\} \tag{10}$$

$$\mathbf{D}_{\max} = \left\{d_i \mid d_i = \max_j(d_{ij}), d_{ij} \in \mathbf{D}\right\} \tag{11}$$

$$\Delta = \left\{\delta_i \mid \delta_i = \mathbf{D}_{\min}(i) \Big/ \mathbf{D}_{\max}(i)\right\} \tag{12}$$

10    III. The rule$f_1$ is used to select $x_1$ from $\mathbf{U}'$:

$$f_1: \quad x_1 = \mathbf{U}'(i)$$
$$\text{Subjected to}: \quad \mathbf{D}_{\max}(i) == \max(\mathbf{D}_{\max}), \quad i = 1, \ldots, k \tag{13}$$

12    IV. The rule$f_2$ is used toselect$x_2$from$\mathbf{U}'$:

$$f_2: \quad x_2 = \mathbf{U}'(i)$$
$$\text{Subjected to}: \quad \Delta(i) = \max(\Delta), \quad i = 1, \ldots, k \tag{14}$$

14    V. Finally, update $\mathbf{U_1}$:

$$\mathbf{U_1} = \mathbf{U_1} + x_1 + x_2 \tag{15}$$

16    Through the loop running I-V, $\mathbf{U_1}$ is updated continuously until the number of samples in $\mathbf{U_1}$ reaches

17    the threshold.

18    The implementation of Q has been introduced. Next, the significance of$f_1$ and $f_2$is described to prove

19    the rationality of Q.

20    2.3.2 Significance of $f_1$ and$f_2$

21    In order to facilitate the description, Fig.3 illustrates the distribution of the sample $x(x \in \mathbf{U}')$ and the
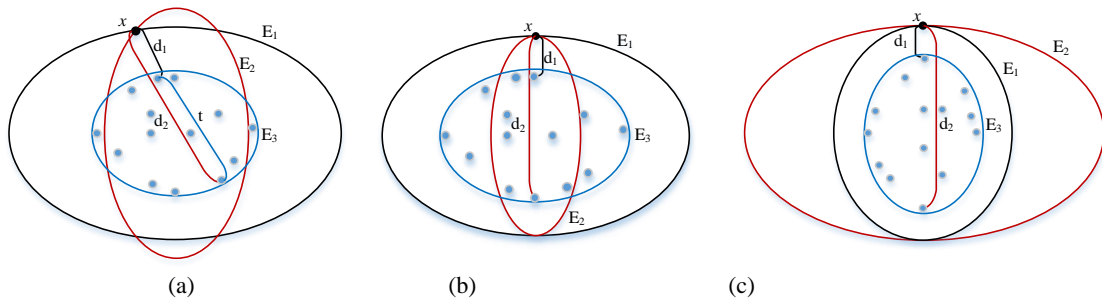
22    training set $\mathbf{L}$.



        (a)            (b)            (c)

25    **Fig. 3**Distribution of the sample $x$ and training set $\mathbf{L}$. (a) $x$ distributed around E$_3$randomly. (b)$x$ distributed in the direction of the

26    short axis of E$_3$. (c)$x$ distributed in the direction of the long axis of E$_3$.

27    For all the figures in Fig.3, $\mathbf{L}$ is shown as the blue dots and E$_3$is the boundary; $x$ is shown as the black

28    dots, $d_1 \in \mathbf{D_{min}}$ and $d_2 \in \mathbf{D_{max}}$(shown in the equation(10) and (11)); the points in E$_1$ have the same

29    minimum distance with $d_1$; the points in E$_2$ have the same maximum distance with $d_2$.

For the $f_1$, Fig. 3(a) shows that the bigger $d_2$ generates a longer distance between $x$ and $\mathbf{L}$, and a smaller probability of $x$ belonging to $\mathbf{L}$, so more information is contained in $x$ which can decrease the amount of the semilabeled samples used by the SSL. Therefore, $f_1$ can select the informative sample and reduce the semilabeled samples to the greatest extent.

For the $f_2$, Fig. 3(a) shows that the $\delta$ ($\delta \in \Lambda$, in (12)) in $x$ can be written as

$$\delta = d_1/d_2 \tag{16}$$

In order to facilitate the description of $f_2$, the expression of $\delta$ needs to be changed.

Let: 
$$d_1 = d_2 - t \tag{17}$$

Where $t$ is the secant of $x$ in $E_3$, then substitute (17) into (16):

$$\delta = (d_2 - t)/d_2 \tag{18}$$

Get: 
$$\delta = 1 - t/d_2 \tag{19}$$

where it can be seen that: when $d_2$ is invariant, the bigger $\delta$ is, the smaller $t$ is, which means that $x$ would be closer to the short axis of $E_3$ (shown in Fig. 3(b)) and can make distribution of $\mathbf{L}$ more uniform (compared with Fig. 3(c)); when t is constant, the bigger $\delta$ is, the smaller $d_2$ is, and the longer distance between $x$ and $\mathbf{L}$ is, which means that the smaller of the probability of $x$ belonging to the $\mathbf{L}$, so the more information that $x$ contains. Therefore, $f_2$ cannot only reduce the semilabeled samples by selecting informative ones, but also make the distribution of $\mathbf{L}$ better.

2.4 Search contextual patterns

This function makes the current pixel $x(x \in \mathbf{U_1})$ as the center sample to find the context patterns $\mathbf{U_2}$ from the neighborhood. There are many choices for $\mathbf{U_2}$. Espinola et al. (2015) introduced the most common types shown in Fig.4, where Von Neumann neighborhood and Moore neighborhood are called the first-order system and the second-order system(Bruzzone and Persello 2009), and Von Neumann neighborhood is also called the four-directly neighbored pixels (Zhao et al. 2013).



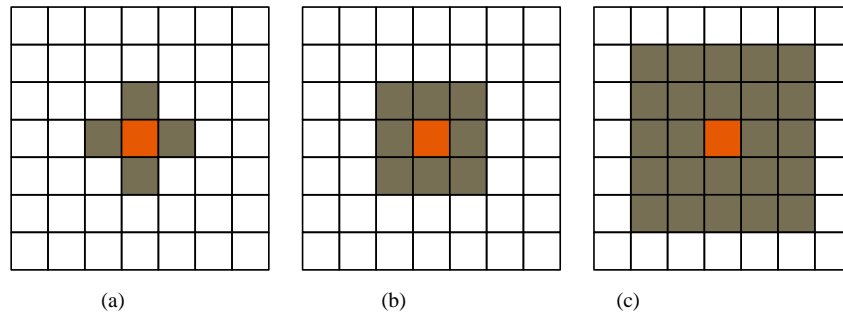(a)                    (b)                    (c)

**Fig. 4** Examples of neighborhood systems. (a) Von Neumann neighborhood, (b) Moore neighborhood and (c) Extended Moore neighborhood

Context information is very significant. Firstly, the center sample and the context patterns are close in the space, which makes their characteristics exist great relevance that can be applied to calculate the image texture feature(D'Elia et al. 2014)or generate the classification map(Bruzzone and Persello 2009). More semilabeled samples are found using context information, which can improve the convergence speed for the SSL. Secondly, their labels can also provide very important information. In(Espinola et al. 2015), the cell automaton distinguishes the different pixels (such as the boundary, noise, or common) at each iteration by utilizing their labels. The robustness of the SSL algorithm can be increased by

1 properly using the context information of labels. The next section describes the way to relabel samples

2 by using the information of labels.

3 2.5 H: Relabeling samples

4 Q selects the samples between the maximum margin bands, where the selected samples are more likely

5 to be contaminated and mislabeled manually. To increase the reliability of the selected samples, the

6 central samples $(\mathbf{U_1},\mathbf{Y_1})$ and the context patterns $(\mathbf{U_2},\mathbf{Y_2})$are relabeled by the function H with the

7 following rules:

8 I. When $\mathbf{Y_2}$are consistent with $\mathbf{Y_1}$, they are left;

9 II. When $\mathbf{Y_2}$are inconsistent with $\mathbf{Y_1}$, $\mathbf{Y_1}$ is incorrect. To make $\mathbf{L}$ representative, the corresponding

10 central sample is discarded and the context patterns are left.

11 III. When a part of $\mathbf{Y_2}$ is consistent with$\mathbf{Y_1}$, for$\mathbf{Y_2}$, the inconsistent context patterns are removed; for $\mathbf{Y_1}$,

12 the proportion of consistent part to context patterns is calculated: if the proportion is greater than

13 the threshold, rules A would be followed, else rules B would be chosen.

14 The semilabeled set $(\mathbf{U_3},\mathbf{Y_3})$ is obtained after being relabeled by H. Then, those samples repeated with

15 $\mathbf{L}$ are eliminated, which generates the semilabeled set $\mathbf{L_1}= (\mathbf{U_4},\mathbf{Y_4})$to add to $\mathbf{L}.$

16 Section 2.1-2.5 describes the proposed SSL based method. In order to overcome the influence of the

17 randomness of the learning set $\mathbf{L}$, the classification results are fused in the next section.

18 2.6 Fusion method

19 Data fusion is a very popular data processing technology to make up for the defects caused by the

20 missing data or noise pollution. It can be used in many methods. The principal component analysis

21 (PCA) extracts the principal components to fuse the training set, and the spectral information and

22 spatial information are combined together to fuse images. Our fusion method is shown in Fig.5,where

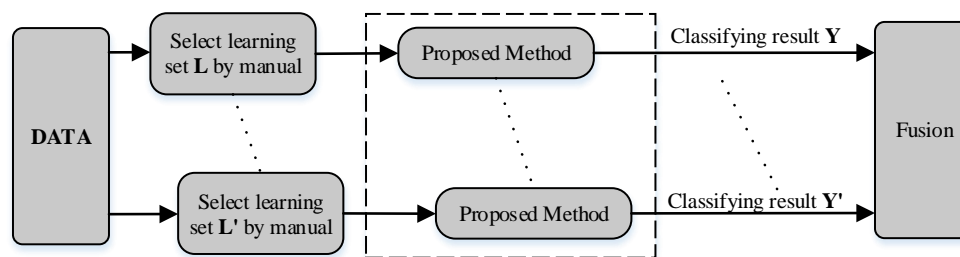23 the part in the dashed rectangular box is the SSL based method proposed in Section 2.1:



24

25 **Fig. 5**Fusion method

26 Firstly, because only a few labeled samples are needed to initialize the training set, the initial $\mathbf{L}$ ($\mathbf{L'}$ and

27 etc.) are selected from the **DATA** manually and would not cost a lot. Then, $\mathbf{L}$ is processed by the

28 proposed method to classify the learning set $\mathbf{U}$ for obtaining the classification result $\mathbf{Y}$. The other initial

29 sets are processed by the same way. Finally, the different classification results ($\mathbf{Y},\mathbf{Y'}$and etc.) are treated

30 as the different data models to be fused by the given different weights based on the detection accuracy.

**3 Experiments**

In this section, the related experiments are designed to measure the performance of the proposed algorithm (Method Proposed, PM) in three aspects. Firstly, to compare PM with the supervised algorithm, the experiments on PM and SVM are carried out. Secondly, to compare PM with the semisupervised algorithm, the experiments on PM, CS[4]VM and PS[3]VM-D are carried out. Finally, to further improve the performance of PM, the different experimental results are fused.

3.1 Dataset description and parameters setting

3.1.1 Dataset Description

Our experiments on change detection are performed on real RSIs, as shown in Fig. 6, where the regions A and regions B indicate the changed regions and unchanged regions respectively. The aim of change detection is to distinguish the changed part and unchanged part of images (Habib et al. 2009; Celik 2010), and can be regarded as the binary classification problem. For the regions A: Fig. 6(a)-(b) are the city images taken before and after the earthquake and possess the rich texture features because of many artificial constructions. The change is due to the accumulation of people. Fig. 6(c)-(d)are the ground images taken at different periods and do not possess obvious texture features. The change is due to the appearance of many white patches. Fig. 6(e)-(f)are the coast images taken at different periods and also do not possess obvious texture features. The change is due to the appearance of a few small-scale gray patches. For the regions B of Fig.6, their geomorphic features are unchanged basically, only the gray changes a little.
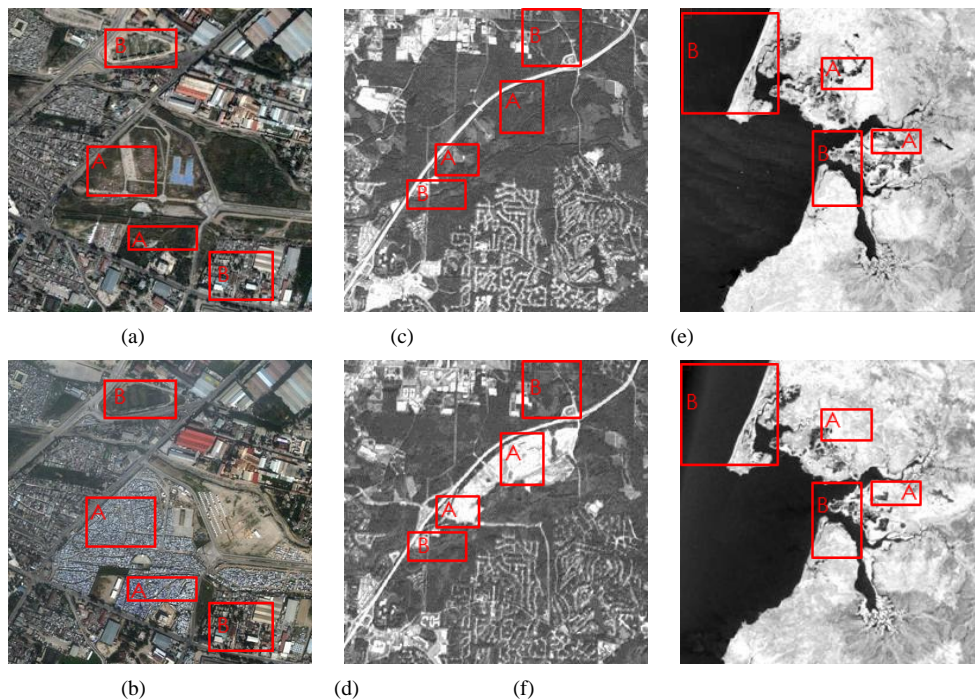


(a)          (c)          (e)

(b)          (d)          (f)

**Fig. 6** Remote sensing images ofthe city (a) before earthquake and (b) after earthquake, the ground (c) before change and (d) after change, and the coast (a) before change and (b) after change. A represents the changed areas, B represents the unchanged areas

1    After all the images are processed by the geometric registration and radiometric calibration, then the

2    ratio map of grayscales is generated. To characterize the images, for Fig. 6(a)-(b) the gray level

3    co-occurrence matrix is calculated as the **DATA₁**; for Fig. 6(c)-(d) and Fig. 6(e)-(f), the sample texture

4    featuresare calculated as the **DATA₂** and **DATA₃** respectively. All the experimental data are normalized

5    and divided into the training set and the learning set respectively.

6    3.1.2Parameters setting

7    The cost factor $C$ of the PM is nonlinear, defined as follows:

8    $$C = \alpha \cdot (N-1)^2 + C_1 \qquad (20)$$

9    Where $N$ is the number of iterations; $C_1$ is the initial cost factor; $\alpha$is the weight coefficient. $\alpha$ and $C_1$are

10    set to 1, and the change range of $C$ is [1,100]. With the increase of iterations, the reliability of the

11    selected samples would increase, and the corresponding $C$ would become larger than the previous one.

12    For the PS³VM-D, $C$ is shown in (20), where$C_1$ is set to 2 and$\alpha$is calculated by

13    $$\alpha = (C_{max} - C_{min})/(r-1)^2 \qquad (21)$$

14    Where $C_{max}$ and $C_{min}$ are the maximum and minimum of $C$; $r$ is designed artificially, with $r = 10$ herein.

15    For the CS⁴VM, $C$ is defined as:

16    $$C = 2 \bullet (N-1) + C_1 \qquad (22)$$

17    where $N$ is also the number of iterations; The initial value of $C_1$ is 2. $C$ is weighted by the $K$

18    $$K = C/C' = 2 \qquad (23)$$

19    When the label of the contextual pattern is the same as the center sample, C is exploited, otherwise $C'$

20    is explored.

21    All the methods are performed using the Gauss kernel function, with the radial width set 0.6. The Von

22    Neumann neighborhood is applied by the PM and CS⁴VM.

23    3.2 Comparison with SVM

24    Considering the **DATA₁**, the performance of the SVM and PM is compared. Fig. 7 shows the detection

25    grayscales of the two algorithms.



26

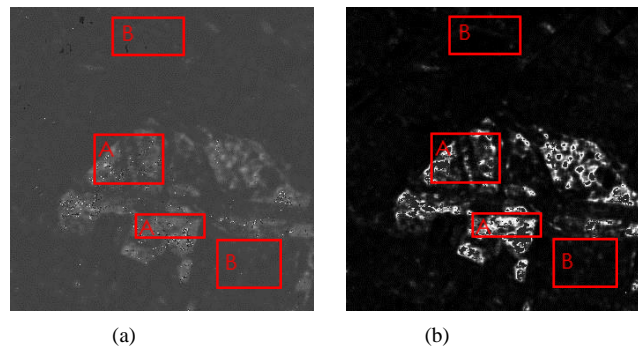27                   (a)                       (b)

28    **Fig. 7** Detection grayscales achieved on **DATA₁**in terms of (a) SVM and (b) PM.A represents the changed areas, B represents the

29    unchanged areas

30    Through the threshold processing, the detection maps are obtained, as shown in Fig. 8 where the

31    white area is detected with the changes and the black area is detected as unchanged; A represents the

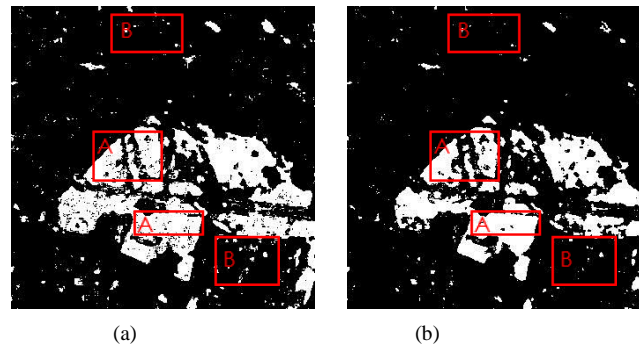1 changed areas and the B represents the unchanged areas.



2
3 (a) (b)

**Fig. 8** Detection maps achieved on the **DATA₁** in terms of (a) SVM and (b) PM, where the white area is detected with the changes and the black area is detected as unchanged; A represents the changed areas and the B represents the unchanged areas

6 There is much omission in A and false alarm in B in Fig. 8(a), while the detection error is much less in
7 Fig. 8(b). In order to compare the performance of the PM in detail, 6 experiments of PM were
8 implemented to obtain the results in terms of the overall accuracy (OA) and Kappa, as shown in Fig. 9.



9
10 (a) (b)

11 **Fig. 9** Results achieved on **DATA₁** in terms of (a) OA and (b) Kappa

12 From Fig.9, except the results of the fifth experiment in which the PM are slightly lower than the SVM,
13 others are significantly higher than the SVM in OA and Kappa. Therefore, in comparison with the
14 SVM, the detection accuracy and performance of the PM have been improved greatly.

15 The training set of the SVM is obtained by the threshold segmentation. The threshold must be large
16 enough to ensure the accuracy of the training set, which means that the samples with a high diversity
17 are probably not within the threshold range. Moreover, the SVM does not exploit the contextual
18 information. Although only a few initial training samples were used, the PM effectively absorbs the
19 abundant and informative unlabeled samples to add to the training set by the AL, which makes it
20 outperform the detection performance compared with the SVM.

21 3.3 Compared with the PS³VM-D and CS⁴VM

22 Considering the **DATA₂**and **DATA₃**, the performance of the PS³VM-D, CS⁴VM and PM is compared.
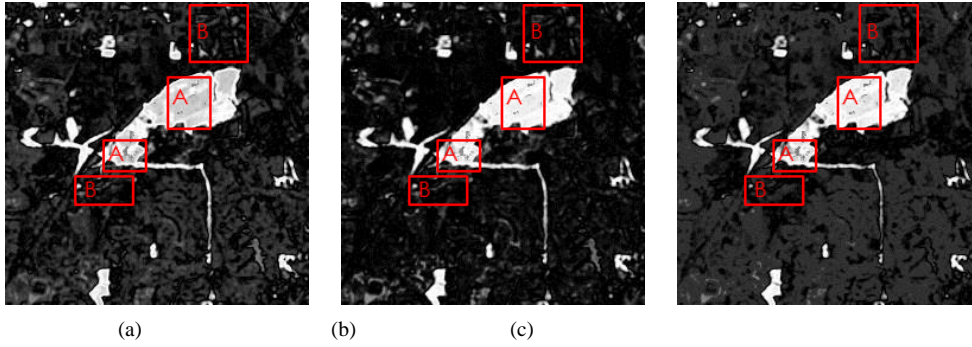23 Fig. 10 and Fig. 11 show the detection grayscales of the three algorithms.

**Fig. 10** Detection grayscales achieved on **DATA₂** in terms of (a) CS⁴VM, (b) PS³VM-D and (c) PM.A represents the changed areas, B represents the unchanged areas
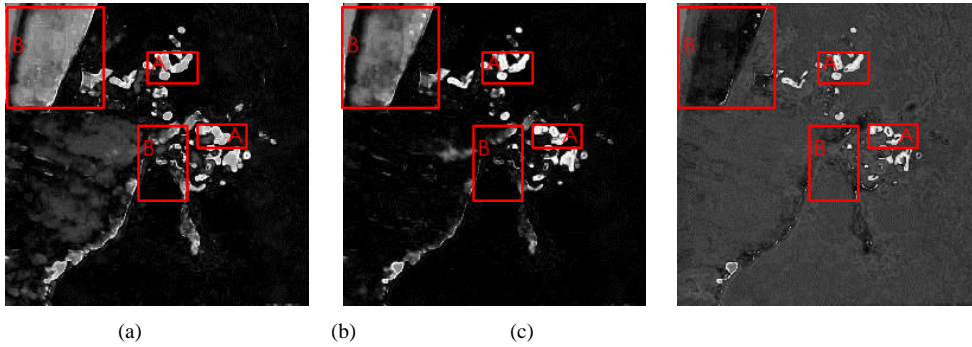


**Fig. 11** Detection grayscales achieved on **DATA₃** in terms of (a) CS⁴VM, (b) PS³VM-D and (c) PM.A represents the changed areas, B represents the unchanged areas

Through the threshold processing, the detection maps of **DATA₂** and **DATA₃** are obtained, as shown in Fig. 12 and Fig. 13 respectively, where the white area is detected with the changes and the black area is detected as unchanged; A represents the changed areas and B represents the unchanged areas.
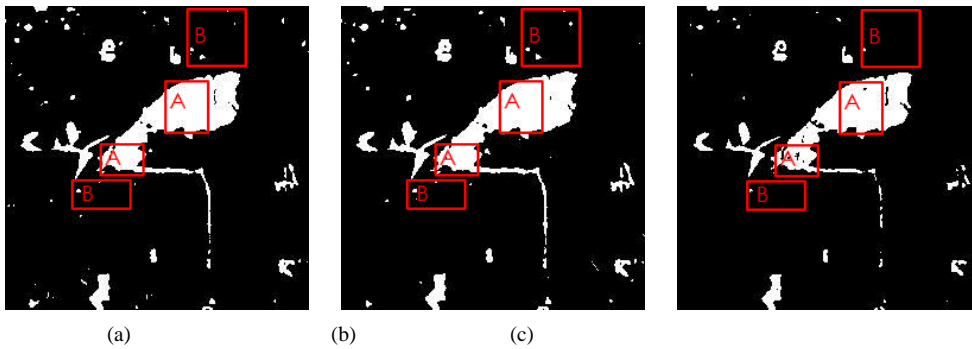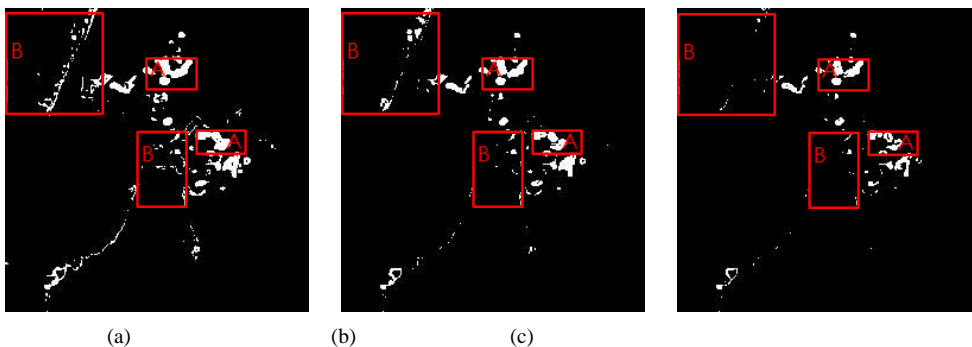


**Fig. 12** Detection maps achieved on **DATA₂** in terms of (a) CS⁴VM, (b) PS³VM-D and (c) PM，where the white area is detected with the changes and the black area is detected as unchanged; A represents the changed areas and the B represents the unchanged areas



**Fig. 13** Detection maps achieved on **DATA₃** in terms of (a) CS⁴VM, (b) PS³VM-D and (c) PM，where the white area is detected

3    In Fig. 12, the omission of PS$^3$VM-D is the least and that of the PM is the greatest in A; the false alarm

4    of PM is the least and that of the CS$^4$VMisthegreatest in B. It can be seen that the performance of the

5    three algorithms is ideal and that of PM is more accurate than others for **DATA$_2$**. In Fig.13, the false

6    alarm of PM is significantly less than that of the other algorithms in B, followed by thePS$^3$VM-D. The

7    omission of three algorithms is almost the same in A. Therefore, the detection result of PM is the best

8    for **DATA$_3$**.

9    In order to know the stability of the three algorithms, 6 different experiments were implemented for the

10    CS$^4$VM, PS$^3$VM-D and PM, respectively with the same initial training sets to obtain the results in

11    terms of OA and Kappa. The line charts of different results are drawn from small to large for

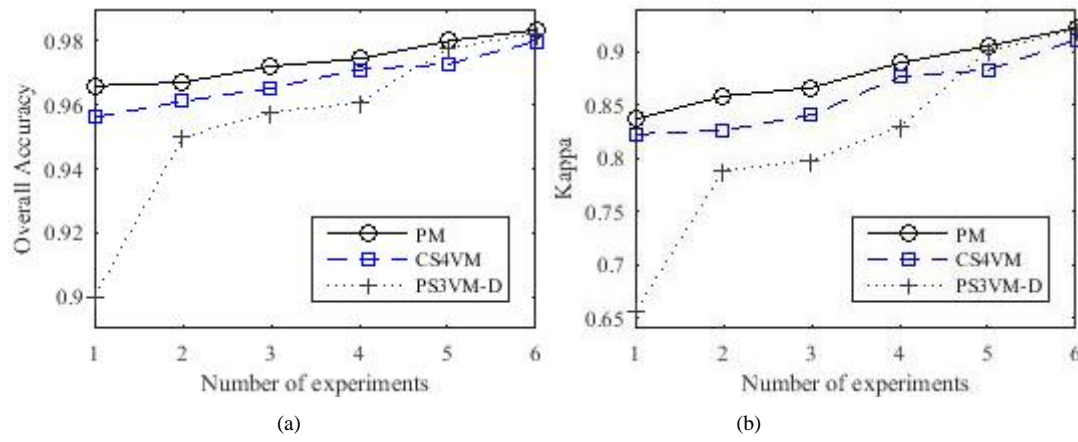12    **DATA$_2$**and **DATA$_3$**, which areshown in Fig.14 and Fig.15.



13

14                  (a)                                   (b)

15    **Fig. 14** Line charts of CS$^4$VM, PS$^3$VM-D and PM achieved on **DATA$_2$** in terms of (a) OA and (b) Kappa
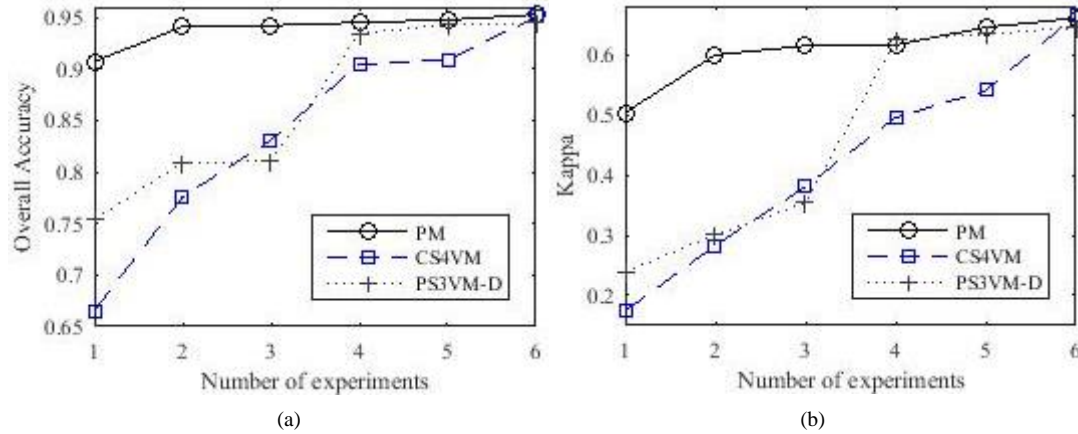


16

17                  (a)                                   (b)

18    **Fig. 15** Line charts of CS$^4$VM, PS$^3$VM-D and PM achieved on **DATA$_3$** in terms of (a) OA and (b) Kappa

19    For **DATA$_2$**and **DATA$_3$**, except that the maximum of the OA and Kappa are basically the same for the

20    three algorithms, the rest of PM are significantly higher than those of thePS$^3$VM-D and CS$^4$VM. For

21    different initial sets, the change range of OA and Kappa for PS$^3$VM-D and CS$^4$VM is great and that for

22    PM is small. Therefore, the detection accuracy and stability of the PMare higher than other algorithms.

23    **DATA$_3$**is selected to illustrate the convergence of the CS$^4$VM, PS$^3$VM-D and PM. The curve of OA is

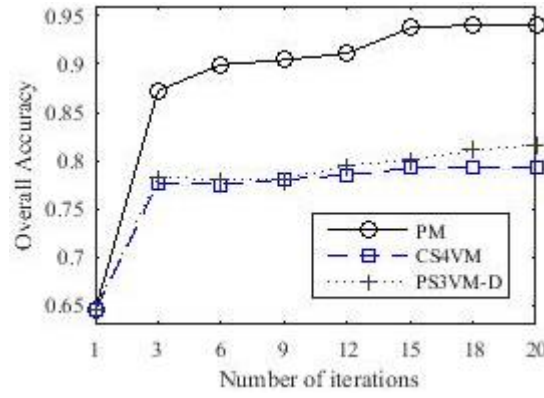24    plotted in Fig. 16, where the X axis represents the number of iterations and the Y axis represents OA.

With the increase of iterations, the OAs of the three algorithms all change, where the OA of the $CS^4VM$

increases from 0.64 to 0.79, the OA of the $PS^3VM$-D increases from 0.64 to 0.82 and the OA of PM

increases from 0.64 to 0.94. The curve of the PM rises fastest, which suggests that the PM can get the

faster convergence speed than other algorithms.

The $CS^4VM$ can increase the relevance of the training samples by searching the contextual samples,

but it does not consider whether those samples are beneficial to improve the convergence rate. The

$PS^3VM$-D guarantees the information of the selected samples by searching them within the margin

band, but the probability that these samples are polluted is also great. Compared to the $CS^4VM$ and the

$PS^3VM$-D, the PM uses the AL to select the informative samples, then use the context information to

increase the samples correlation, which improves the detection accuracy and ensures the reliability of

the selected samples. We use a new query function Q to make the distribution of the training samples

more uniform for improving the detection efficiency and accuracy. Moreover, the convergence of the

PM is further improved by applying the semilabeled samples.

3.4 Fusion for PM

To obtain the more accurate results and verify the effectiveness of the fusion for the SSL methods, the

experimental results of the PM are fused. The OA and Kappa are listed in Table 1, and the fusion

detection maps are shown in Fig. 17, where the white area is detected with the changes and the black

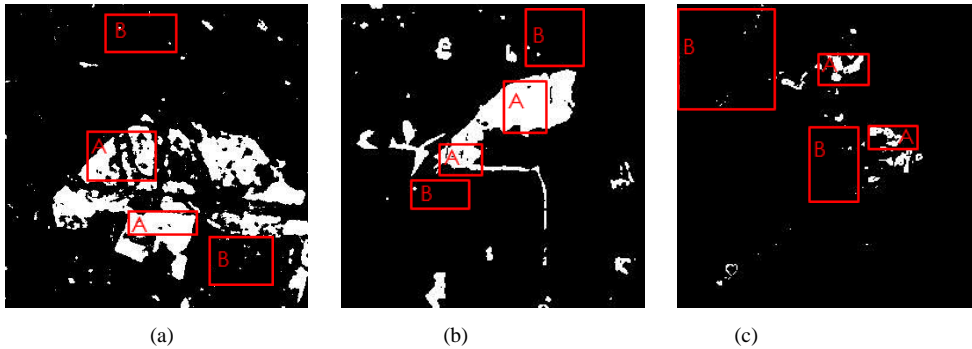area is detected as unchanged. A represents the changed areas and B represents the unchanged areas.



| (a) | (b) | (c) |

| | 1 | | 2 | | 3 | | 4 | | 5 | | 6 | | Fusion | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | OA | Kappa | OA | Kappa | OA | Kappa | OA | Kappa | OA | Kappa | OA | Kappa | OA | Kappa |
| DATA1 | 0.9718 | 0.8955 | 0.9662 | 0.8839 | 0.9668 | 0.8760 | 0.9717 | 0.8997 | 0.9413 | 0.8045 | 0.9724 | 0.8977 | **0.9744** | **0.9060** |
| DATA2 | 0.9659 | 0.8582 | 0.9721 | 0.8661 | 0.9835 | 0.9225 | 0.9800 | 0.9053 | 0.9670 | 0.8365 | 0.9744 | 0.8894 | **0.9837** | **0.9250** |
| DATA3 | 0.9475 | 0.6155 | 0.9420 | 0.6139 | 0.9409 | 0.6447 | 0.9451 | 0.5991 | 0.9072 | 0.5031 | 0.9528 | 0.6599 | **0.9534** | **0.6638** |

In TABEL 1, both the OA and Kappa get higher after fusion. By comparing Fig. 17(a), Fig. 17(b) and Fig. 17(c) with Fig. 8(b), Fig. 12(c) and Fig. 13(c) respectively, it is found that the false alarm in B and the omission in A reduces a lot after the fusion for **DATA1** and **DATA3**. The false alarm in B region and the omission in A decreases slightly for **DATA2**. Obviously, the data fusion could effectively compensate the detection error, reduce the risk caused by the different initial training sets and improve the detection accuracy.

**4. Conclusion**

In this paper, a novel method is proposed for the change detection of RS images. It fully incorporates the methods and concepts in the SSL, and adopts them to fit the situation where the labeled samples are insufficient.

The novelty of this paper lies in: a) By considering the advantages of the AL and the context information, a novel semisupervised method is designed; b) by analyzing the distribution of the samples, a new query function is designed to select the semilabeled samples using the Euclidean distance; c) Based on the idea of data fusion, the discrete results of the PM are effectively fused. In the experiments of change detection for actual RSIs, the PM has made a significant improvement in the detection accuracy, convergence rate, and the stability in comparison with other existing methods. It can be further improved by using other effective fusion methods.

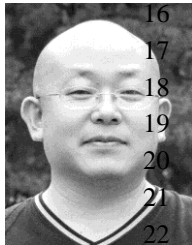**References**

Anand, S., Mittal, S., Tuzel, O., & Meer, P. (2014). Semi-Supervised Kernel Mean Shift Clustering. *Pattern Analysis and Machine Intelligence, IEEE Transactions on, 36*(6), 1201-1215, doi:10.1109/TPAMI.2013.190.

Bazi, Y., Melgani, F., & Al-Sharari, H. D. (2010). Unsupervised change detection in multispectral remotely sensed imagery with level set methods. *Geoscience and Remote Sensing, IEEE Transactions on, 48*(8), 3178-3187.

Bovolo, F., Bruzzone, L., & Marconcini, M. (2008). A Novel Approach to Unsupervised Change Detection Based on a Semisupervised SVM and a Similarity Measure. *Geoscience and Remote Sensing, IEEE Transactions on, 46*(7), 2070-2082, doi:10.1109/TGRS.2008.916643.

Bruzzone, L., & Persello, C. (2009). A Novel Context-Sensitive Semisupervised SVM Classifier Robust to Mislabeled Training Samples. *Geoscience and Remote Sensing, IEEE Transactions on, 47*(7), 2142-2154, doi:10.1109/TGRS.2008.2011983.

Celik, T. (2010). Change Detection in Satellite Images Using a Genetic Algorithm Approach. *Geoscience and Remote Sensing Letters, IEEE, 7*(2), 386-390, doi:10.1109/LGRS.2009.2037024.

Chapelle, O., Schölkopf, B., & Zien, A. (2006). Semi-supervised learning.

D'Elia, C., Ruscino, S., Abbate, M., Aiazzi, B., Baronti, S., & Alparone, L. (2014). SAR Image Classification Through Information-Theoretic Textural Features, MRF Segmentation, and Object-Oriented Learning Vector Quantization. *Selected Topics in Applied Earth Observations and Remote Sensing, IEEE Journal of, 7*(4), 1116-1126, doi:10.1109/JSTARS.2014.2304700.

Demir, B., Persello, C., & Bruzzone, L. (2011). Batch-Mode Active-Learning Methods for the Interactive Classification of Remote Sensing Images. *Geoscience and Remote Sensing, IEEE Transactions on, 49*(3), 1014-1031, doi:10.1109/TGRS.2010.2072929.

Didaci, L., Fumera, G., & Roli, F. (2012). Analysis of Co-training Algorithm with Very Small Training Sets. In G. Gimel'farb, E. Hancock, A. Imiya, A. Kuijper, M. Kudo, S. Omachi, et al. (Eds.), *Structural, Syntactic, and Statistical Pattern Recognition* (Vol. 7626, pp. 719-726, Lecture Notes in Computer Science): Springer Berlin Heidelberg.

Espinola, M., Piedra-Fernandez, J. A., Ayala, R., Iribarne, L., & Wang, J. Z. (2015). Contextual and Hierarchical Classification of Satellite Images Based on Cellular Automata. *Geoscience and Remote Sensing, IEEE Transactions on, 53*(2), 795-809, doi:10.1109/TGRS.2014.2328634.

Geiß, C., & Taubenböck, H. (2013). Remote sensing contributing to assess earthquake risk: from a literature review towards a roadmap. *Natural Hazards, 68*(1), 7-48, doi:10.1007/s11069-012-0322-2.

Gomez-Chova, L., Bruzzone, L., Camps-Valls, G., & Calpe-Maravilla, J. Semi-Supervised Remote Sensing Image Classification based on Clustering and the Mean Map Kernel. In *Geoscience and Remote Sensing Symposium, 2008. IGARSS 2008. IEEE International, 7-11 July 2008 2008* (Vol. 4, pp. IV - 391-IV - 394). doi:10.1109/IGARSS.2008.4779740.

Gong, C., Junwei, H., Lei, G., Zhenbao, L., Shuhui, B., & Jinchang, R. (2015). Effective and Efficient Midlevel Visual Elements-Oriented Land-Use Classification Using VHR Remote Sensing Images. *Geoscience and Remote Sensing, IEEE Transactions on, 53*(8), 4238-4249, doi:10.1109/TGRS.2015.2393857.

Habib, T., Inglada, J., Mercier, G., & Chanussot, J. (2009). Support Vector Reduction in SVM Algorithm for Abrupt Change Detection in Remote Sensing. *Geoscience and Remote Sensing Letters, IEEE, 6*(3), 606-610, doi:10.1109/LGRS.2009.2020306.

Hansen, M. C., & Loveland, T. R. (2012). A review of large area monitoring of land cover change using Landsat data. *Remote Sensing of Environment, 122*, 66-74, doi:http://dx.doi.org/10.1016/j.rse.2011.08.024.

Hearst, M. A., Dumais, S. T., Osman, E., Platt, J., & Scholkopf, B. (1998). Support vector machines. *Intelligent Systems and their Applications, IEEE, 13*(4), 18-28, doi:10.1109/5254.708428.

Izquierdo-Verdiguier, E., Laparra, V., Gomez-Chova, L., & Camps-Valls, G. (2013). Encoding Invariances in Remote Sensing Image Classification With SVM. *Geoscience and Remote Sensing Letters, IEEE, 10*(5), 981-985, doi:10.1109/LGRS.2012.2227297.

Junwei, H., Dingwen, Z., Gong, C., Lei, G., & Jinchang, R. (2015a). Object Detection in Optical Remote Sensing Images Based on Weakly Supervised Learning and High-Level Feature Learning. *Geoscience and Remote Sensing, IEEE Transactions on, 53*(6), 3325-3337, doi:10.1109/TGRS.2014.2374218.

Junwei, H., Dingwen, Z., Xintao, H., Lei, G., Jinchang, R., & Feng, W. (2015b). Background Prior-Based Salient Object Detection via Deep Reconstruction Residual. *Circuits and Systems for Video Technology, IEEE Transactions on, 25*(8), 1309-1321, doi:10.1109/TCSVT.2014.2381471.

Kawakita, M., & Kanamori, T. (2013). Semi-supervised learning with density-ratio estimation. *Machine Learning, 91*(2), 189-209, doi:10.1007/s10994-013-5329-8.

Li, J., Bioucas-Dias, J. M., & Plaza, A. (2010). Semisupervised Hyperspectral Image Segmentation Using Multinomial Logistic Regression With Active Learning. *Geoscience and Remote Sensing, IEEE Transactions on, 48*(11), 4085-4098, doi:10.1109/TGRS.2010.2060550.

Maulik, U., & Chakraborty, D. (2012). A novel semisupervised SVM for pixel classification of remote sensing imagery. *International Journal of Machine Learning and Cybernetics, 3*(3), 247-258, doi:10.1007/s13042-011-0059-3.

Maulik, U., & Chakraborty, D. (2014). Fuzzy Preference Based Feature Selection and Semisupervised SVM for Cancer Classification. *NanoBioscience, IEEE Transactions on, 13*(2), 152-160, doi:10.1109/TNB.2014.2312132.

Munoz-Mari, J., Tuia, D., & Camps-Valls, G. (2012). Semisupervised Classification of Remote Sensing Images With Active Queries. *Geoscience and Remote Sensing, IEEE Transactions on, 50*(10), 3751-3763, doi:10.1109/TGRS.2012.2185504.

Pasolli, E., Melgani, F., Tuia, D., Pacifici, F., & Emery, W. J. (2014). SVM Active Learning Approach for Image Classification Using Spatial Information. *Geoscience and Remote Sensing, IEEE Transactions on, 52*(4), 2217-2233, doi:10.1109/TGRS.2013.2258676.

Persello, C., & Bruzzone, L. (2014). Active and Semisupervised Learning for the Classification of Remote Sensing Images. *Geoscience and Remote Sensing, IEEE Transactions on, 52*(11), 6937-6956, doi:10.1109/TGRS.2014.2305805.

Schohn, G., & Cohn, D. Less is more: Active learning with support vector machines. In *ICML, 2000* (pp. 839-846): Citeseer

Scholkopf, B., Kah-Kay, S., Burges, C. J. C., Girosi, F., Niyogi, P., Poggio, T., et al. (1997). Comparing support vector machines with Gaussian kernels to radial basis function classifiers. *Signal Processing, IEEE Transactions on, 45*(11),

2758-2765, doi:10.1109/78.650102.

Shahshahani, B. M., & Landgrebe, D. A. (1994). The effect of unlabeled samples in reducing the small sample size problem and mitigating the Hughes phenomenon. *Geoscience and Remote Sensing, IEEE Transactions on, 32*(5), 1087-1095, doi:10.1109/36.312897.

Tuia, D., Volpi, M., Copa, L., Kanevski, M., & Munoz-Mari, J. (2011). A Survey of Active Learning Algorithms for Supervised Remote Sensing Image Classification. *Selected Topics in Signal Processing, IEEE Journal of, 5*(3), 606-617, doi:10.1109/JSTSP.2011.2139193.

Yi, Y., Wu, J., & Xu, W. (2011). Incremental SVM based on reserved set for network intrusion detection. *Expert Systems with Applications, 38*(6), 7698-7707.

Zhao, C., Li, X., Ren, J., & Marshall, S. (2013). Improved sparse representation using adaptive spatial support for effective target detection in hyperspectral imagery. *International Journal of Remote Sensing, 34*(24), 8669-8684, doi:10.1080/01431161.2013.845924.

Zhu, X. (2010). Semi-Supervised Learning. In C. Sammut, & G. Webb (Eds.), *Encyclopedia of Machine Learning* (pp. 892-897): Springer US.
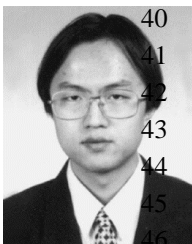
## Author Biographies

**FeiGao** received the B.S. and M.S. degree from the Xi'an Petroleum Institute, Xi'an, China, in 1996 and 1999, respectively, and the Ph.D. degrees from Beijing University of Aeronautics and Astronautics (BUAA), Beijing, China, in 2005. He is currently an Associate Professor with the School of Electronic and Information Engineering, BUAA. He is interested in radar signal processing, moving target detection and image processing.
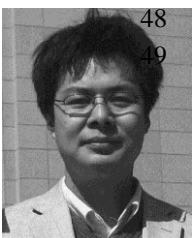
**Wenchao Lv** received the B.S. degree in electronic and information engineering fromBeijing University of Posts and Telecommunications, Beijing, China, in 2014. He is currently pursuing the M.E. degree in electronic and communication engineering at Beijing University of Aeronautics and Astronautics, Beijing, China. His current research activity is in machine learning and synthetic aperture radar image classification.

**Yaotian Zhang** received the B.S. degree in Electronic Engineering from BeihangUniversityin 2003. He was awarded Ph.D degree in Signal Processing from Beihang University in 2009. Since2010, he has worked in School of Electronic and Information Engineering of Beihang University as an assistant professor. His research interest covers image understanding, target detection and Micro-Doppler signal analysis.

**Jinping Sun** received the M.Sc. and Ph.D. degrees from the Beijing University of Aeronautics and Astronautics (BUAA), Beijing, China, in 1998 and 2001, respectively. He is currently a Professor with the School of Electronic and Information Engineering, BUAA. His research interests include high-resolution radar signal processing, image understanding, and robust beamforming.

**Jun Wang** received the B.S. degree from the Northwestern Polytechnical University, Xi'an, China, in 1995 and the M.S. and Ph.D. degrees from the Beijing University of Aeronautics and Astronautics

(BUAA), Beijing, China, in 1998 and 2001, respectively. He is currently a Professor with the School of Electronic and Information Engineering, BUAA. He is interested in signal processing, DSP/FPGA real-time architecture, target recognition and tracking, and so on. His research has resulted in 38 papers in journals, books, and conference proceedings.

**Erfu Yang** is a Lecturer in the Department of Design Manufacture and Engineering Management (DMEM) at the University of Strathclyde, Glasgow, UK. His main research interests include robotics, autonomous systems, mechatronics, manufacturing automation, computer vision, image/signal processing, nonlinear control, process modelling and simulation, condition monitoring, fault diagnosis, multi-objective optimizations, and applications of machine learning and artificial intelligence including multi-agent reinforcement learning, fuzzy logic, neural networks, bio-inspired algorithms, and cognitive computation, etc. He has over 60 publications in these areas, including more than 30 journal papers and 5 book chapters.