

Unsupervised Deep Slow Feature Analysis for Change Detection in Multi-Temporal Remote Sensing Images

Bo Du, *Senior Member, IEEE*, Lixiang Ru, Chen Wu, *Member, IEEE*, and Liangpei Zhang, *Fellow, IEEE*

Abstract—Change detection has been a hotspot in remote sensing technology for a long time. With the increasing availability of multi-temporal remote sensing images, numerous change detection algorithms have been proposed. Among these methods, image transformation methods with feature extraction and mapping could effectively highlight the changed information and thus has better change detection performance. However, changes of multi-temporal images are usually complex, existing methods are not effective enough. In recent years, deep network has shown its brilliant performance in many fields including feature extraction and projection. Therefore, in this paper, based on deep network and slow feature analysis (SFA) theory, we proposed a new change detection algorithm for multi-temporal remote sensing images called Deep Slow Feature Analysis (DSFA). In DSFA model, two symmetric deep networks are utilized for projecting the input data of bi-temporal imagery. Then, the SFA module is deployed to suppress the unchanged components and highlight the changed components of the transformed features. The CVA pre-detection is employed to find unchanged pixels with high confidence as training samples. Finally, the change intensity is calculated with chi-square distance and the changes are determined by threshold algorithms. The experiments are performed on two real-world datasets and a public hyperspectral dataset. The visual comparison and quantitative evaluation have both shown that DSFA could outperform the other state-of-the-art algorithms, including other SFA-based and deep learning methods.

Index Terms—Change detection, Deep network, Slow feature analysis, Remote sensing images.

I. INTRODUCTION

CHANGE detection is defined as the process of identifying differences in the state of an object or phenomenon by observing it at different times [1]. With the rapid development of remote sensing technology, more remote sensing images of the earth surface are now available [2]–[4]. The multi-temporal remote sensing images covering the same area could help to detect land-cover and land-use changes, so that change

detection could be better applied to diverse real-world applications, such as deforestation monitoring, damage assessment, vegetation phenology variation study, and disaster monitoring [5]–[10].

Generally, change detection algorithms could be divided into the following categories: 1) Image algebra methods mainly include image difference, image ratio, image regression, and change vector analysis [11], [12]. These methods directly calculate the difference between multi-temporal remote sensing images; 2) Image transformation algorithms extract the effective features of multi-temporal remote sensing images by transforming and combining their feature bands, and mainly include Principle Component Analysis (PCA) [13], Multivariate Alteration Detection (MAD) [14], [15], Gram-Schmidt transformation (GS) [16] and Independent Component Analysis [17]; 3) Classification methods mainly include post-classification and compound classification, which are both based on classification to obtain land-use categories [18]–[21]; 4) Other advanced methods contains the algorithms based on wavelet, Markov random field, and local gradual descent, etc. [22]–[25]. Among all these kinds of change detection algorithms, image transformation methods have been widely studied and applied. The basic idea of image transformation is projecting the original multiband images into a new feature space to better separate changed and unchanged pixels. In this process, the most crucial work is to find an effective projecting algorithm to extract the determinative features.

Changed pixels in multi-temporal remote sensing images always have the feature differences with diverse change directions, while the features of unchanged pixels are supposed to be generally invariant [1]. However, owing to the atmospheric conditions, illumination and sensor calibration and so on, those unchanged pixels always have slight differences [26], [27]. Compared with changed pixels, changes of unchanged pixels usually have the consistent direction. By minimizing the feature variation of unchanged pixels, changed pixels could also be highlighted and separated. Inspired by this idea, slow feature analysis is proposed for detecting real changes and obtained satisfactory performance [28], [29].

SFA is a feature learning algorithm that extracts invariant and slowly varying features from input signals [30], [31]. And it has been successfully applied to solve diverse real-world problems, such as human action recognition, dynamic texture recognition and time series analysis, etc [32]–[35]. In change detection problems, changed and unchanged pixels correspond to quickly and slowly varying features in SFA, respectively.

Manuscript submitted December 2, 2018, revised June 12, 2019. This work was supported in part by the National Natural Science Foundation of China under Grants 61601333, 61822113, and 41871243. *Corresponding author: Chen Wu.*

B. Du is with the School of Computer Science, and Collaborative Innovation Center of Geospatial Technology, Wuhan University, Wuhan, P.R. China (e-mail: gunspace@163.com).

L. Ru is with the School of Computer Science, Wuhan University, Wuhan, P.R. China (e-mail: rulixiang@whu.edu.cn).

C. Wu and L. Zhang are with the State Key Laboratory of Information Engineering in Surveying, Mapping and Remote Sensing, and School of Computer Science, Wuhan University, Wuhan, P.R. China (e-mail: chen.wu@whu.edu.cn, zlp62@whu.edu.cn).

Based on this theory, *Wu, et al.* [28] used SFA to suppress the spectral difference between slowly varying unchanged pixels, so that the changed pixels can be highlighted and well detected. By solving SFA problems, the proposed algorithms in [28] could get the projecting matrices to map original data, so that the unchanged components could be suppressed. All these algorithms have shown their good performance in some real-world remote sensing images. However, limited by the feature representative ability, linear SFA algorithms are sometimes not able to separate the changed and unchanged pixels. The potential solutions include projecting original feature into a higher-dimensional complex feature space to improve the model's complexity and feature representation ability.

Actually, in [36], *Wu, et al.* proposed a kernel version slow feature analysis (KSFA) for scene change detection. And the results have also shown that nonlinear extension of SFA is effective. However, in this method, KSFA is only designed for computing the change probabilities of bi-temporal scene level features. Some of its details are not suitable for pixel-wise change detection of multi-spectral imagery. Besides, KSFA is sensitive to the selection of kernel function. Different kernel function could lead to very different performances.

Deep networks have been proved to have a powerful ability of representing non-linear functions, and thus can project original features into a more complex feature space [37], [38]. Due to the growing availability of both data and computing resources, deep neural networks have been resurging in these years. Numerous kinds of networks have been developed to complete different tasks, such as classification [39], detection [40], segmentation [41], and feature mapping [38], etc. Besides, in recent years, deep networks have also been applied to learn non-linear transformations of highly correlated datasets, and performed well [42].

Therefore, inspired by the idea of utilizing deep network learning non-linear transformations, we propose a new algorithm called Deep Slow Feature Analysis (DSFA) in this paper. In DSFA, two deep networks are used to extract and represent the features of remote sensing images obtained at different times, respectively. The transformed features by deep networks are then taken as the inputs of SFA to obtain the projecting matrix. The projecting matrix could extract the most invariant component of multi-temporal remote sensing images, so the changed pixels could be accentuated. We formulate the loss function for DSFA model to make sure that the transformed features can represent the original data better. The intention of DSFA is to extract the invariant components of input features, which means that utilizing unchanged pixels as the inputs will help accelerating the training process and improving the final performance. However, in fact, labeled data are usually rare in remote sensing problems. Therefore, in DSFA, we use CVA to make a pre-detection and find unchanged pixel pairs as the inputs for training process. When the deep network is converged, the transformed features will be calculated by passing original features through trained networks. Then the difference of transformed features in SFA space is calculated. Finally, the change intensity map is calculated with chi-square distance, and the binary change map is obtained with threshold algorithms.

The rest of this paper is organized as follows. Section II introduces the SFA theory and the details of SFA in change detection. Section III presents the algorithm details of proposed DSFA. In Section IV, we implement our proposed method and perform experiments on two real-world datasets and a public hyperspectral dataset. In Section V, some settings of our experiments are discussed. And Section VI draws the conclusion of this paper.

II. SLOW FEATUE ANALYSIS

In this section, we'll introduce the mathematical theory of SFA, and how SFA is extended to solve change detection problems. Mathematically, SFA is formulated as follows:

Given a multi-dimensional temporal signal $s(t) = [s_1(t), s_2(t), \dots, s_n(t)]$, where n represents the dimension and $t \in [t_0, t_1]$, the target of SFA is finding a set of transforming functions $[g_1(x), g_2(x), \dots, g_M(x)]$ to generate the output signal $z(t) = [g_1(s), g_2(s), \dots, g_M(s)]$ and ensuring that transformed signal is time invariant as possible. Mathematically, the objective function of SFA is

$$\min_{g_j} : \langle (\dot{g}_j(s))^2 \rangle_t, j \in [1, 2, \dots, M], \quad (1)$$

under the following constraints:

$$\langle g_j(s) \rangle_t = 0, \quad (2)$$

$$\langle g_j(s)^2 \rangle_t = 1, \quad (3)$$

$$\forall i < j : \langle g_i(s)g_j(s) \rangle_t = 0, \quad (4)$$

where $\langle g_j(s) \rangle_t$ denotes the mean signal of $g_j(s)$ over time t and $\dot{g}_j(s)$ is the first-order derivate of $g_j(s)$. Therefore, the objective of SFA is minimizing the mean value of the first-order derivate of transformed signal. Among these constraints, Constraint (2) is to simplify the process of solving the optimization problem. Constraint (3) ensures that each output signal could contain certain information. And Constraint (4) is presented to eliminate the correlation between output signals and force each signal carries different type of information.

In the linear case, the transforming function could be expressed as a mapping matrix:

$$g_j(s) = w_j^T s, \quad (5)$$

where w_j^T denotes the transposition of w_j . And the objective function and constraints could be reformulated as follows:

$$\langle (w_j^T \dot{s})^2 \rangle_t = w_j^T \langle \dot{s} \dot{s}^T \rangle_t w_j = w_j^T A w_j, \quad (6)$$

$$\langle (w_j^T s) \rangle_t = 0, \quad (7)$$

$$\langle (w_j^T s)(w_j^T s) \rangle_t = w_j^T \langle s s^T \rangle_t w_j = w_j^T B w_j = 1, \quad (8)$$

$$\langle (w_i^T s)(w_j^T s) \rangle_t = w_i^T \langle s s^T \rangle_t w_j = w_i^T B w_j = 0. \quad (9)$$

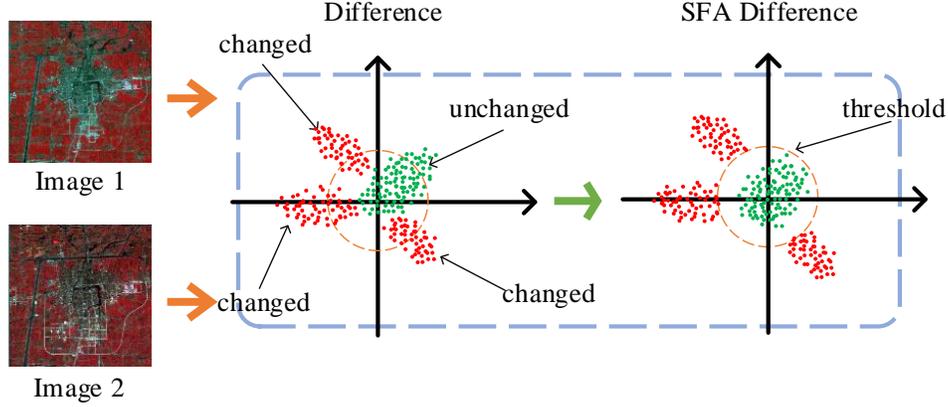


Fig. 1. A schematic of SFA in change detection.

In (6), $A = \langle \dot{s}\dot{s}^T \rangle_t$ is the expectation of the covariance matrix of the first-order derivative of input signals. (7) represents Constraint (2), and it can be implemented by pre-processing the input data. (8) and (9) denote Constraint (3) and (4), respectively. And $B = \langle ss^T \rangle_t$ is the expectation of covariance matrix of original input signals.

In SFA theory, (9) can be integrated to (6) as follows:

$$\langle (w_j^T \dot{s})^2 \rangle_t = w_j^T A w_j = \frac{w_j^T A w_j}{w_j^T B w_j} = \frac{\langle (w_j^T \dot{s})^2 \rangle_t}{\langle (w_j^T s)(w_j^T s) \rangle_t}. \quad (10)$$

And this optimization problem can be solved by the generalized eigenvalue problem:

$$AW = BW\Lambda, \quad (11)$$

where W and Λ is the generalized eigenvector matrix and a diagonal matrix of eigenvalues, respectively. According to (10) and (11), the most invariant component of the output signal has the smallest eigenvalue.

In pixel-based change detection problems, the input signals are raw pixels of remote sensing images, which are discrete. In consequence, SFA need to be reconstructed to cope with discrete cases. As shown in Figure 1, the objective of SFA in change detection problems is suppressing unchanged pixels to highlight changed ones, so that they could be separated much easier. Mathematically, let $x_i, y_i \in \mathbb{R}^m$ denote corresponding pixels in bi-temporal remote sensing images, where m is the number of bands. After normalizing the input data, the objective of SFA is reformulated as

$$\min_{w_j} : \frac{1}{n} \sum_{i=1}^n (w_j^T x_i - w_j^T y_i)^2, \quad (12)$$

where n is the total number of pixels. And constraints are rewritten as

$$\frac{1}{2n} \left[\sum_{i=1}^n w_j^T x_i + \sum_{i=1}^n w_j^T y_i \right] = 0, \quad (13)$$

$$\frac{1}{2n} \left[\sum_{i=1}^n (w_j^T x_i)^2 + \sum_{i=1}^n (w_j^T y_i)^2 \right] = 1, \quad (14)$$

$$\frac{1}{2n} \left[\sum_{i=1}^n (w_j^T x_i)(w_i^T x_i) + \sum_{i=1}^n (w_j^T y_i)(w_i^T y_i) \right] = 0. \quad (15)$$

In the generalized eigenvalue problem of SFA, A and B in (11) are reformulated as follows:

$$A = \frac{1}{n} \sum_{i=1}^n (x_i - y_i)(x_i - y_i)^T, \quad (16)$$

$$B = \frac{1}{2n} \left[\sum_{i=1}^n x_i x_i^T + \sum_{i=1}^n y_i y_i^T \right]. \quad (17)$$

When A and B are obtained, the eigenvector matrix W will be solved. By normalizing W , the final mapping matrix is obtained.

$$\hat{w}_j = \frac{w_j}{\sqrt{w_j^T B w_j}}. \quad (18)$$

Then the change detection result, the difference between transformed bi-temporal images, is calculated as $D_j = \hat{w}_j^T x_j - \hat{w}_j^T y_j$.

III. METHODOLOGY

As mentioned above, those existing SFA-based change detection algorithms are all linear. In order to improve the representing ability of features and final change detection performance, in this section, we propose Deep Slow Feature Analysis (DSFA). The main structure of DSFA is shown in Figure 2.

As we can see in Figure 2, the input of DSFA is pairwise pixels of multi-temporal imagery. Then DSFA could be roughly divided to two parts: Deep Network module and SFA constraint. In the Deep Network module, two symmetric networks, whose layers are all Fully Connected Layer, are

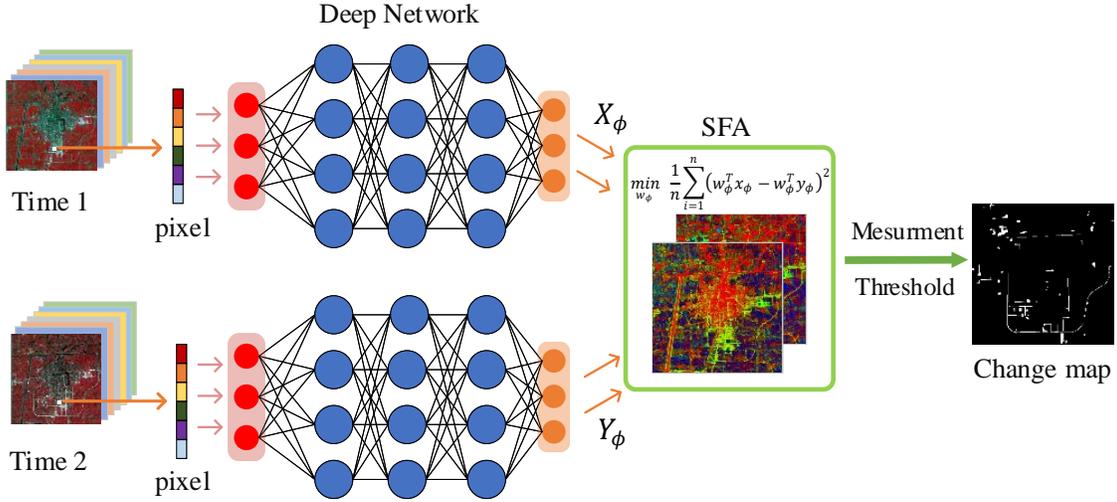


Fig. 2. A schematic of DSFA, consisting two deep networks.

used to project original input data into a new complex high-dimensional feature space. In Figure 2, the red nodes denote the nodes of input layers, the blue nodes represent the nodes of hidden layers and the yellow nodes are used to represent output layers. Each hidden layer of the Deep Network module has the same number of nodes. After the original data is transformed, we use the SFA constraint to suppress the invariant components and highlight the changed components of transformed features. We formulate the loss function of DSFA so that the parameters of deep networks could be solved based on gradient-based optimization algorithms.

A. Formulation

Mathematically, DSFA is defined as follows: Assuming the original bi-temporal remote sensing images are $X, Y \in \mathbb{R}^{m \times n}$, where m and n respectively denote the number of feature bands and pixels. For clarity, let h_i denotes the number of nodes of the i -th hidden layer of the networks, and o is the number of nodes of the output layer. Given an instance X , the output of the first hidden layer could be formulated as

$$f_1^1(X) = s(w_1^1 X + b_1^1), \quad (19)$$

where $w_1^1 \in \mathbb{R}^{h_1 \times m}$ and $b_1^1 \in \mathbb{R}^{h_1}$ denote the weight matrix and the bias vector, respectively. And $s(\cdot)$ represents the activation function. The output of the subsequent layers is calculated in the same way. For a network with l hidden layers, the output of the last hidden layer is $f_l^1(X) = s(w_l^1 f_{l-1}^1(X) + b_l^1)$, where $w_l^1 \in \mathbb{R}^{h_l \times h_{l-1}}$ and $b_l^1 \in \mathbb{R}^{h_l}$. After that, $f_l^1(X)$ will be mapped by the output layer.

Finally, the final transformed feature of this network is

$$X_\phi = f(\theta_1, X) = s(w_o^1 f_l^1(X) + b_o^1), \quad (20)$$

where $w_o^1 \in \mathbb{R}^{o \times h_l}$ and $b_o^1 \in \mathbb{R}^o$ are the weight matrix and bias vector, respectively. And θ_1 is the set of all the parameters

in the network, including $w_1^1, \dots, w_l^1, w_o^1$ and $b_1^1, \dots, b_l^1, b_o^1$. And for another instance Y , Y_ϕ has a symmetric expression and meaning.

$$Y_\phi = f(\theta_2, Y) = s(w_o^2 f_l^2(X) + b_o^2). \quad (21)$$

When the original given data is mapped into a new high dimensional feature space by deep networks, let $\hat{X}_\phi = X_\phi - \frac{1}{n} \mathbf{1} X_\phi$ and $\hat{Y}_\phi = Y_\phi - \frac{1}{n} \mathbf{1} Y_\phi$ denote the centralized X_ϕ and Y_ϕ , respectively, where $\mathbf{1} \in \mathbb{R}^{o \times o}$ is a matrix whose elements are all 1. Then the covariance matrix of transformed data will be calculated.

$$\Sigma_{XX} = \hat{X}_\phi \hat{X}_\phi^T + r * I, \quad (22)$$

$$\Sigma_{YY} = \hat{Y}_\phi \hat{Y}_\phi^T + r * I, \quad (23)$$

$$\Sigma_{XY} = (\hat{X}_\phi - \hat{Y}_\phi)(\hat{X}_\phi - \hat{Y}_\phi)^T. \quad (24)$$

where I denotes the identity matrix and r is a regularization constant. Assume that $r > 0$, so that Σ_{XX} and Σ_{YY} are both positive definite and invertible. Therefore, in DSFA problem, the generalized eigenvalue problem to be solved is formulated as:

$$A_\phi W = B_\phi W \Lambda \Leftrightarrow B_\phi^{-1} A_\phi W = W \Lambda, \quad (25)$$

where $A_\phi = \Sigma_{XY}$ and $B_\phi = \frac{1}{2}(\Sigma_{XX} + \Sigma_{YY})$. According to (22 – 24), the final form of this problem is

$$\left[\frac{1}{2}(\Sigma_{XX} + \Sigma_{YY}) \right]^{-1} \Sigma_{XY} W = W \Lambda. \quad (26)$$

Based on SFA theory, the most invariant component has the smallest eigenvalue. Thus, the objective of DSFA could be designed as minimizing the total square of all eigenvalues, so that the variance of unchanged pixels can be suppressed

and changed pixels are much easier to be detected. The loss function of DSFA then could be formulated as follows:

$$\mathcal{L}(\theta_1, \theta_2) = \text{tr}[(B_\phi^{-1}A_\phi)^2], \quad (27)$$

where $\text{tr}(\cdot)$ denotes the trace of a matrix. Utilizing (27), the loss value of DSFA could be calculated and the parameters of networks θ_1 and θ_2 can be obtained with gradient-based optimization algorithm.

B. Optimization

To calculate the gradient of $\mathcal{L}(\theta_1, \theta_2)$ with respect to all the w_l^v and b_l^v , we could use the back-propagation algorithm, which requires computing the gradient of $\mathcal{L}(\theta_1, \theta_2)$ with respect to \hat{X}_ϕ and \hat{Y}_ϕ .

According to the reference [43], and using the fact that A_ϕ and B_ϕ are both symmetric, we could then have:

$$\nabla_A = \frac{\partial \mathcal{L}(\theta_1, \theta_2)}{\partial A_\phi} = 2B_\phi^{-1}A_\phi B_\phi^{-1}, \quad (28)$$

$$\nabla_B = \frac{\partial \mathcal{L}(\theta_1, \theta_2)}{\partial B_\phi} = -2B_\phi^{-1}A_\phi B_\phi^{-1}A_\phi B_\phi^{-1}. \quad (29)$$

Utilizing the derivation in [42], we could have the gradient of A_ϕ with respect to each element of \hat{X}_ϕ :

$$\begin{aligned} \frac{\partial A_\phi^{ab}}{\partial \hat{X}_\phi^{ij}} &= \frac{1}{n}(\xi_{(a=i)}\hat{X}_\phi^{bj} + \xi_{(b=i)}\hat{X}_\phi^{aj}) \\ &- \frac{1}{n}(\xi_{(a=i)}\hat{Y}_\phi^{bj} + \xi_{(b=i)}\hat{Y}_\phi^{aj}), \end{aligned} \quad (30)$$

where $\xi_{(e)}$ represents the indicator function. If e is true, then $\xi_{(e)} = 1$, otherwise $\xi_{(e)} = 0$. Similarly, the gradient of B_ϕ with respect to each element of \hat{X}_ϕ is computed as follows:

$$\frac{\partial B_\phi^{ab}}{\partial \hat{X}_\phi^{ij}} = \frac{1}{2n}(\xi_{(a=i)}\hat{X}_\phi^{bj} + \xi_{(b=i)}\hat{X}_\phi^{aj}), \quad (31)$$

Integrating (28)-(31), the gradient of $\mathcal{L}(\theta_1, \theta_2)$ with respect to \hat{X}_ϕ^{ij} is:

$$\begin{aligned} \frac{\partial \mathcal{L}(\theta_1, \theta_2)}{\partial \hat{X}_\phi^{ij}} &= \sum_{ab} \nabla_A^{ab} \frac{\partial A_\phi^{ab}}{\partial \hat{X}_\phi^{ij}} + \sum_{ab} \nabla_B^{ab} \frac{\partial B_\phi^{ab}}{\partial \hat{X}_\phi^{ij}} \\ &= \frac{2}{n}(\nabla_A \hat{X}_\phi - \nabla_A \hat{Y}_\phi)_{ij} + \frac{1}{n}(\nabla_B \hat{X}_\phi)_{ij}. \end{aligned} \quad (32)$$

The derivation process isn't straight and its details are presented in Appendix A. Finally, it's obvious that the gradient of $\mathcal{L}(\theta_1, \theta_2)$ with respect to \hat{X}_ϕ could be computed as:

$$\frac{\partial \mathcal{L}(\theta_1, \theta_2)}{\partial \hat{X}_\phi} = \frac{2}{n}(\nabla_A \hat{X}_\phi - \nabla_A \hat{Y}_\phi) + \frac{1}{n}\nabla_B \hat{X}_\phi. \quad (33)$$

And for another instance Y_ϕ , the expression of $\mathcal{L}(\theta_1, \theta_2)/\partial Y_\phi$ is symmetric. We then could utilize Gradient Descent algorithms to minimize the loss to obtain the parameters of deep network module of DSFA.

According to loss function, the objective of DSFA is projecting the difference of pairwise pixels into an invariant difference feature space. Therefore, if we utilize unchanged pairwise pixels as training samples, the learned non-linear

projection of deep network will have better performance in extracting the invariant components. However, in practice, priori labeled information in change detection is always hard to get. To select unchanged pairwise pixels for training process, in this paper, we use the CVA method to make a pre-detection. In this process, CVA and Kmeans method are employed to obtain the difference map and the binary change map of input multi-temporal imagery, respectively. Training samples are then randomly selected from the detected unchanged areas.

After obtained the training set and trained the network, the original data will be passed through the deep network to get the transformed features X_ϕ and Y_ϕ . Then, the generalized eigenvalue problem will be solved to obtain the projecting matrix w_ϕ and the difference between mapped features is calculated as follows:

$$D_\phi = w_\phi^T X_\phi - w_\phi^T Y_\phi. \quad (34)$$

Then the change intensity of bi-temporal images could be calculated. In order to eliminate the differences in the scale of each feature bands, in this paper, we use chi-square distance to measure the intensity of changes, which is calculated as

$$\text{chi2} = \sum_{i=1}^m \frac{(D_\phi^i)^2}{\sigma_i^2}. \quad (35)$$

In (35), m is the number of feature bands, and σ^2 is variance of each bands obtained by statistically analyzing. Threshold algorithms, such as OSTU method and Kmeans method, are then employed to get the final binary change map. The whole detailed process of training and generating binary change map for DSFA is summarized in Algorithm 1.

IV. EXPERIMENT

To evaluate the performance of DSFA, in this section, we implement DSFA on TensorFlow and perform experiments on three multi-temporal remote sensing image datasets. Datasets used in our experiment include two Enhanced Thematic Mapper (ETM) datasets and a public hyperspectral change detection dataset. The first one is Taizhou dataset, covering the city of Taizhou, China, acquired in 2000 and 2003. And the second is Nanjing dataset, which are respectively acquired in 2000 and 2002. Both datasets were obtained by the Landsat 7 Enhanced Thematic Mapper Plus (ETM+) sensor with a spatial resolution of 30 m. And 6 spectral bands (1-5 and 7) are selected for our experiments. The band 6 has a spatial resolution of 60m, so it's dropped and not used in our experiments. The third dataset is River dataset¹, and consists of two hyperspectral images with a size of 463×241 , which are respectively obtained in May, 2013 and December, 2013, Jiangsu Province, China. Each image in this dataset contains 198 spectral bands after noisy bands removal.

A. Experiment settings

In the DSFA model, the weight and bias matrices of each layer are initialized randomly, and need to be optimized. The other values, including the number of layers and nodes in each

¹Avaliable: <http://crabwq.github.io/>

Algorithm 1 Process of training and generating binary change map for DSFA.

Input:

Multi-temporal input images I^1 and I^2 ;

Output:

The binary change map D ;

- 1: Standardize I^1 and I^2 using $z - score$ method;
- 2: Employ CVA pre-detection to generate training samples X and Y ;
- 3: Initialize the network's parameters $\{\theta_1, \theta_2\}$;
- 4: **while** $i < max_epochs$ **do**
- 5: Calculate the projected features $X_\phi = f_1(X, \theta_1)$ and $Y_\phi = f_2(Y, \theta_1)$;
- 6: Calculate the loss value:

$$\mathcal{L}(\theta_1, \theta_2) = tr[(B_\phi^{-1} A_\phi)^2];$$

- 7: Calculate the gradient: $\partial\mathcal{L}(\theta_1, \theta_2)/\partial\theta_1$ and $\partial\mathcal{L}(\theta_1, \theta_2)/\partial\theta_2$;
- 8: Updating the parameters using Gradient Descent algorithm;
- 9: $i++$;
- 10: **end while**
- 11: Calculate the mapped features I_ϕ^1 and I_ϕ^2 of I^1 and I^2 ;
- 12: Solve SFA problem to obtain projecting matrix w_ϕ ;
- 13: Calculate the difference map:

$$\Delta I = w_\phi^T I_\phi^1 - w_\phi^T I_\phi^2;$$

- 14: Threshold to get the binary change map D ;
- 15: **return** D ;

view and the DSFA regularization parameter in (22-24) are hyperparameters. As for the DSFA regularization parameter, we tuned it over the range $[10^{-8}, 10^{-1}]$, and eventually selected 10^{-4} as the value for our proposed model. The influence of the regularization parameter r is discussed in the Section V.

Some other conventional and SFA-based change detection algorithms are also implemented for comparison, including CVA, PCA [13], MAD [14], IRMAD [44], USFA [28], ISFA [28], PCANet [45] and SDPCANet [46]. All of them are unsupervised algorithms. Before calculating the difference map, PCA uses Principal Component Analysis method to project original data into a new lower dimensional feature space. MAD is a change detection method based on the established theory Canonical Correlation Analysis (CCA), which is firstly proposed in [47]. It utilizes CCA to maximize the correlation between the features of multi-temporal images. IRMAD is an iteratively weighted extension of MAD. It firstly calculates the original MAD variates. And in the following iterations, it applies different weights to each pixels or regions to emphasize the changed parts of images. USFA and ISFA are proposed in [28]. Based on the SFA theory, USFA computes a projecting matrix to suppress the unaltered components of input data to highlight changed components. And ISFA is an

iteratively weighted extension of USFA, and has the same way to calculate weights as IRMAD. PCANet method firstly takes gabor wavelets and fuzzy c-means as the pre-detection method to select the training samples. Then, a PCANet [48] model is trained with the image patches centered at the interested pixels. Finally, the change map is obtained by classifying the remain patches with the trained model. SDPCANet developed PCANet by using a context-aware saliency detection method [49] to select more robust and confident training samples in the pre-detection process.

For all these algorithms, we choose all of the output feature bands to calculate the change intensity.

B. Experiments on Taizhou ETM dataset

The study area of the first dataset is Taizhou city, Jiangsu Province, China. The image size is 400×400 . Figure 3 shows the pseudo color and ground truth images of this dataset. (a) and (b) are the pseudo color images acquired in 2002 and 2003, respectively. And (c) is the sampled ground truth image of changed and unchanged regions of Taizhou city, where the green pixels represent the unchanged regions, red pixels represent changed regions. The background of image (c) is the gray scale image of (a), and they denote the unsampled regions. The changed area contains 4227 pixels, and unchanged area contains 17163 pixels.

In the experiment of DSFA on Taizhou dataset, 4000 pixels, which are about 2.5% of the total number of pixels, are randomly selected from the unchanged region of CVA pre-detected image for training to get the parameters of the networks and the projecting matrix of SFA. Due to the use of random initialization, for DSFA, we take the sum of change intensity of 10 independent runs as the final change intensity map, and the presented values of evaluation criteria are the results of the summed intensity map.

Figure 4 shows the change intensity maps of Taizhou dataset by (a) CVA, (b) PCA, (c) MAD, (d) IRMAD, (e) USFA, (f) ISFA, (g) DSFA-64-2, (h) DSFA-128-2, and (i) DSFA-256-2. Since PCANet and SDPCANet are both classification-based methods, there're no intensity maps of them. DSFA- $h-l$ refers to a DSFA model with l hidden layers and each hidden layer has h nodes. All of these change intensity maps are calculated with all the output feature bands. In this figure, brighter regions have bigger change probabilities. As the Figure 4 shows, visually, PCA, ISFA and DSFA-128-2 have the best performance in differentiating the changed and unchanged pixels. The unchanged regions of MAD and IRMAD are grey, which means they could not suppress the unchanged background from changed pixels very well. Similarly, CVA and USFA have bad performance in extracting changed pixels from unchanged background. As for other DSFA-based methods, DSFA-64-2 and DSFA-256-2, they have a moderate performance in change intensity map among all these methods. Though DSFA-based methods visually have some noise points, actually, these noise points probably represent truly changed pixels of unsampled region.

In Table I, we present the accurate evaluation of binary change results segmented by OTSU method. PCANet and SDPCANet are both classification-based methods, so their results

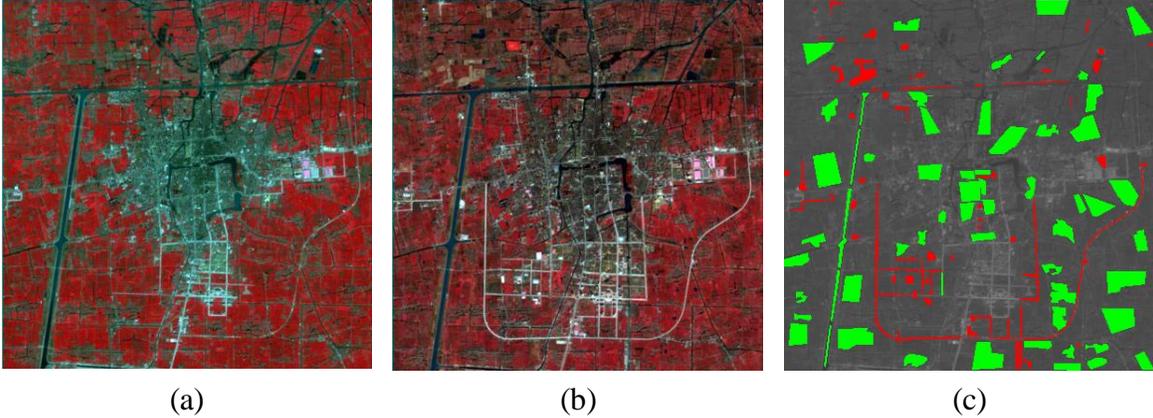


Fig. 3. The pseudo-color images of Taizhou obtained in (a) 2000, (b) 2003, and (c) ground truth.

presented here are their classification results and needn't to be processed by OTSU. The evaluation criteria include the overall accuracy of sampled changed area (OA_CHG), the overall accuracy of sampled unchanged area (OA_UN), the overall accuracy of all sampled regions (OA), Kappa coefficient, and F1 score. The best values of each evaluation criteria are highlighted with bold.

TABLE I
CHANGE DETECTION RESULTS OF TAIZHOU DATASET USING OTSU.

OTSU	OA_CHG	OA_UN	OA	Kappa	F1
CVA	0.8439	0.9970	0.9667	0.8890	0.9093
PCA	0.7755	0.9961	0.9525	0.8374	0.8658
MAD	0.8855	0.9474	0.9352	0.8030	0.8148
IRMAD	0.9056	0.9818	0.9667	0.8942	0.9150
USFA	0.7093	0.9922	0.9363	0.7773	0.8148
ISFA	0.8077	0.9991	0.9612	0.8684	0.8918
PCANet	0.8469	0.9992	0.9691	0.8967	0.9155
SDPCANet	0.9151	0.9863	0.9722	0.9115	0.9287
DSFA-64-2	0.8294	0.9982	0.9648	0.8819	0.9032
DSFA-128-2	0.8985	0.9954	0.9763	0.9227	0.9372
DSFA-256-2	0.8450	0.9966	0.9667	0.8888	0.9090

As the Table I shows, SDPCANet and PCANet have the best performance on OA_CHG and OA_UN, respectively. On detecting unchanged pixels, IRMAD has the second best performance. On the contrary, ISFA performs bad on detection changed regions. And it is worth noting that DSFA-128-2 outperforms the other algorithms on OA, which indicates that it has a higher accuracy in both changed and unchanged part of remote sensing images. And other DSFA-based methods also have very good performance on OA, especially compared with USFA and ISFA. Besides, on Kappa coefficient, all DSFA-based methods have better performance than USFA and ISFA. The Kappa coefficient and F1 score of DSFA-128-2 are respectively 0.9227 and 0.9372, which are also much

better than the other change detection methods. Considering the total detection accuracy of all changed and unchanged pixels, Kappa coefficient, and F1 score, DSFA-128-2 is the best method, and SDPCANet is the second best method and only slightly worse than DSFA-128-2.

The change detection results obtained by Kmeans method are presented in Table II. PCANet and SDPCANet's results presented here also needn't to be processed by Kmeans. As we can see from this table, all of these methods don't show obvious differences in performance when using different threshold algorithms. And this suggests that these methods, including our proposed DSFA-based algorithms, are robust to different threshold methods. The results in Table II are very similar to those in Table I. SDPCANet has the best performance on OA_CHG, but shows lower accuracy on OA_UN. On the contrary, PCANet is the best method in detecting unchanged regions, but has low accuracy in detecting changed pixels. For both changed and unchanged regions, DSFA-128-2 has a detection accuracy of 97.64%, which is still the highest among all methods. DSFA-128-2 also has the highest Kappa coefficient and F1 score. Generally, all of DSFA-based algorithms have pretty good performance. And among all the methods, DSFA-128-2 is still the best one.

In Table III, we present the best change detection results of Taizhou dataset by traversing of all thresholds. Since SDPCANet and PCANet don't need to be post-processed by threshold methods, their presented results are still based on their classification results. In this table, we could see that all DSFA-based methods could outperform the other algorithms exclude CVA and ISFA. And among all DSFA-based methods, DSFA-128-2 has best performance in all evaluation criteria. ISFA has almost the same performance with DSFA-128-2. Besides, it's worth noting that the best change detection results of USFA and ISFA are much better than those obtained with OTSU and Kmeans method, while DSFA-based methods' best results are very close to those using OTSU and Kmeans method. We can conclude that though the best results of ISFA are very close to DSFA, the latter has much better

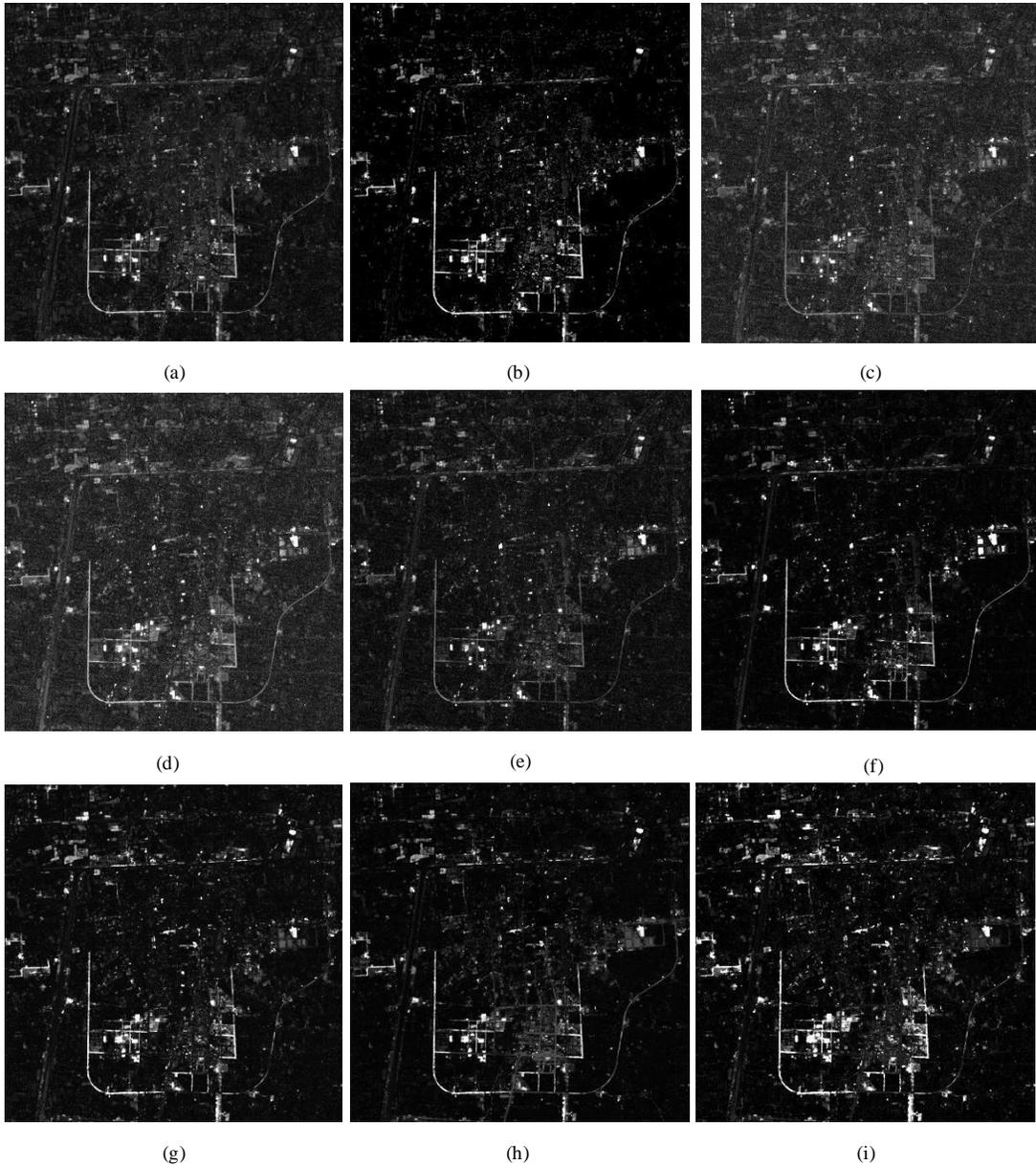


Fig. 4. Change intensity maps of Taizhou dataset by (a) CVA, (b) PCA, (c) MAD, (d) IRMAD, (e) USFA, (f) ISFA, (g) DSFA-64-2, (h) DSFA-128-2, and (i) DSFA-256-2.

discriminability than the former.

In Figure 5, we present the binary change maps obtained by OTSU method of (a) CVA, (b) PCA, (c) MAD, (d) IRMAD, (e) USFA, (f) ISFA, (g) PCANet, (h) SDPCANet, (i) DSFA-64-2, (j) DSFA-128-2 and (k) DSFA-256-2. In this figure, green, red, white, and purple regions represent unchanged pixels that are detected as unchanged, changed pixels that are detected as changed, changed pixels that are detected as unchanged, and unchanged pixels that are detected as changed, respectively. And we could refer them as true positive, false negative, and false positive samples. As Figure 5 presents, intuitively, DSFA-128-2 have the best performance.

And compared with DSFA-128-2, the results of MAD-based methods have more false positive pixels than other algorithms. CVA, PCA and two SFA-based methods tend to classify changed pixels as unchanged. Compared with DSFA-128-2, PCANet has more false negative regions and SDPCANet has more false positive regions. The other DSFA-based methods, DSFA-64-2 and DSFA-256-2, are prone to judge some specific changed regions as unchanged.

C. Experiments on Nanjing ETM dataset

The second experiment is carried on the Nanjing ETM dataset. Nanjing dataset includes two 6 spectral bands remote

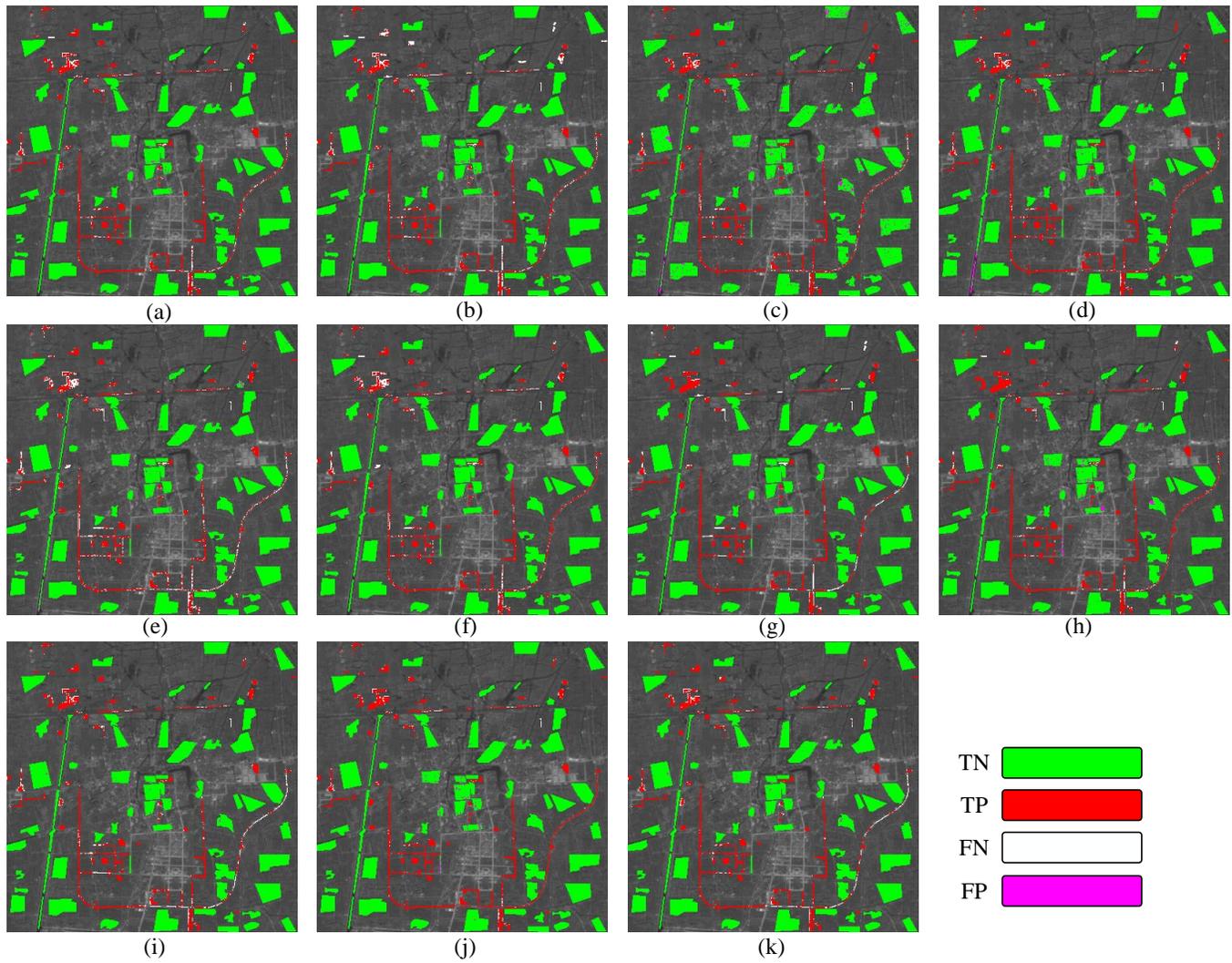


Fig. 5. The binary change maps of (a) CVA, (b) PCA, (c) MAD, (d) IRMAD, (e) USFA, (f) ISFA, (g) PCANet, (h) SDPCANet, (i) DSFA-64-2, (j) DSFA-128-2 and (k) DSFA-256-2.

TABLE II
CHANGE DETECTION RESULTS OF TAIZHOU DATASET USING KMEANS.

Kmeans	OA_CHG	OA_UN	OA	Kappa	F1
CVA	0.8453	0.9970	0.9670	0.8900	0.9102
PCA	0.7731	0.9964	0.9523	0.8365	0.8649
MAD	0.8827	0.9500	0.9367	0.8066	0.8464
IRMAD	0.9054	0.9818	0.9667	0.8942	0.9149
USFA	0.7166	0.9915	0.9372	0.7814	0.8185
ISFA	0.8074	0.9991	0.9612	0.8683	0.8916
PCANet	0.8469	0.9992	0.9691	0.8967	0.9155
SDPCANet	0.9151	0.9863	0.9722	0.9115	0.9287
DSFA-64-2	0.8316	0.9981	0.9652	0.8830	0.9042
DSFA-128-2	0.9006	0.9951	0.9764	0.9232	0.9377
DSFA-256-2	0.8457	0.9966	0.9668	0.8892	0.9094

TABLE III
BEST CHANGE DETECTION RESULTS OF TAIZHOU DATASET.

BEST	OA	Kappa	F1
CVA	0.9756	0.9222	0.9373
PCA	0.9633	0.8810	0.9041
MAD	0.9472	0.8298	0.8626
IRMAD	0.9669	0.8945	0.9150
USFA	0.9476	0.8315	0.8640
ISFA	0.9776	0.9287	0.9426
PCANet	0.9691	0.8967	0.9155
SDPCANet	0.9722	0.9115	0.9287
DSFA-64-2	0.9715	0.9070	0.9254
DSFA-128-2	0.9783	0.9304	0.9439
DSFA-256-2	0.9713	0.9072	0.9250

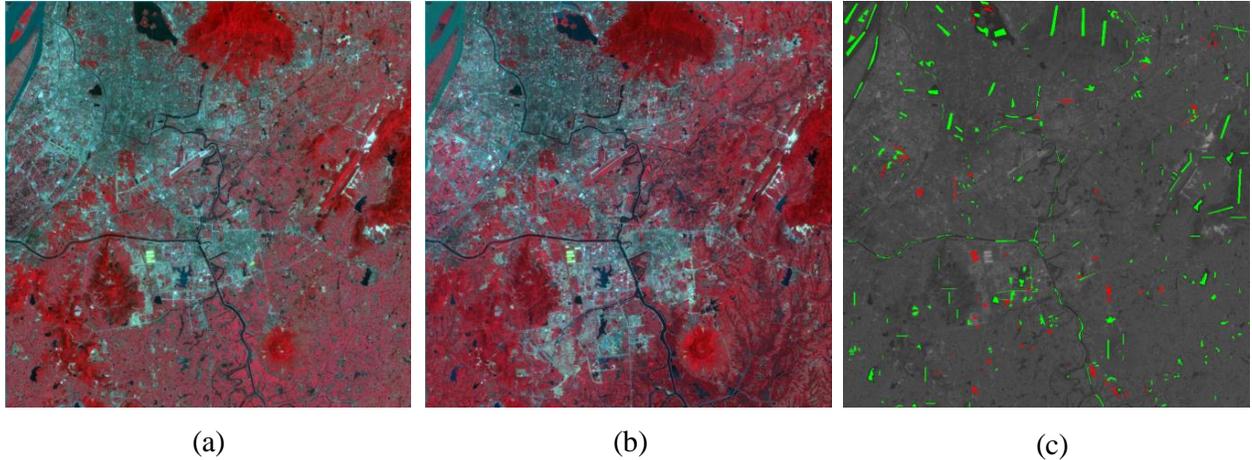


Fig. 6. The pseudo-color images of Nanjing city obtained in (a) 2000, (b) 2002, and (c) ground truth.

sensing images with a size of 800×800 , which are acquired in 2000 and 2002, respectively. Figure 6 presents the pseudo color images of Nanjing city obtained in (a) 2000, (b) 2002, and (c) is the ground truth of sampled changed and unchanged areas. The red part of (c) represents the sampled changed area of Nanjing city, which includes 2363 pixels. And the green part is the sampled unchanged area and includes 12393 pixels.

In the experiment on Nanjing dataset, we randomly select 8000 pixels from unchanged area pre-detected by CVA to train our DSFA model. Like the experiment on Taizhou dataset, the presented results of each evaluation criteria of DSFA are based on the total change intensity map of 10 runs.

Figure 7 shows the change intensity maps of Nanjing dataset by (a) CVA, (b) PCA, (c) MAD, (d) IRMAD, (e) USFA, (f) ISFA, (g) DSFA-64-2, (h) DSFA-128-2, and (i) DSFA-256-2. In this figure, brighter regions have bigger change probabilities. As we can see from this figure, USFA and ISFA have less bright area, which means that they tend to detect much less changed pixels than other change detection algorithms. And CVA, MAD and IRMAD have more bright area which indicates that these methods are prone to categorize these unchanged pixels to changed. DSFA-128-2 and DSFA-256-2 have very close results to each other. Both them have very good discriminability of changed and unchanged pixels. In addition, the result of PCA is also very close to DSFA-64-2. But the distinction between their changed and unchanged regions is not very obvious. On the whole, visually, the result of DSFA-128-2 is the best in calculating the change intensity.

In Table IV, we present the change detection results of Nanjing dataset utilizing OTSU method. The best values of each evaluation criteria are highlighted with bold in this table. As we can see, in general, DSFA-based methods, especially DSFA-128-2, have the best performance among all these methods. DSFA-128-2 could outperform other algorithms on OA_UN, OA, Kappa coefficient and F1 score. And in these criteria, all DSFA-based methods are much better than others. MAD and IRMAD have the best performance on OA_CHG,

which is consistent with their change intensity results. Similar to MAD and IRMAD, CVA and PCA have very high values on OA_CHG, but are far worse than DSFA-based methods on OA_UN, OA and Kappa coefficient. The results of PCANet and SDPCANet are also very similar to the results of PCA method. On the contrary, USFA and ISFA do well in detecting unchanged pixels, but have the lowest accuracy on OA_CHG.

TABLE IV
CHANGE DETECTION RESULTS OF NANJING DATASET USING OTSU.

OTSU	OA_CHG	OA_UN	OA	Kappa	F1
CVA	0.8595	0.9168	0.9076	0.6933	0.7487
PCA	0.8625	0.9363	0.9244	0.7398	0.7853
MAD	0.9534	0.8530	0.8691	0.6236	0.6999
IRMAD	0.9530	0.8922	0.9019	0.6987	0.7568
USFA	0.5959	0.9680	0.9084	0.6234	0.6757
ISFA	0.6416	0.9760	0.9224	0.6816	0.7260
PCANet	0.8680	0.9334	0.9229	0.7367	0.7829
SDPCANet	0.8426	0.9397	0.9242	0.7351	0.7806
DSFA-64-2	0.7288	0.9817	0.9412	0.7647	0.7987
DSFA-128-2	0.7465	0.9806	0.9431	0.7747	0.8078
DSFA-256-2	0.7360	0.9793	0.9403	0.7633	0.7980

Table V shows the evaluation results of the experiment on Nanjing dataset using Kmeans method. Similar to the results of OTSU, compared to MAD-based and SFA-based methods, DSFA is still better in detecting unchanged and changed areas, respectively. On the whole, DSFA-based algorithms have higher overall accuracies, Kappa values and F1 score than others. PCANet-based methods have higher OA_CHG than DSFA-based methods, but worse performances on the other criteria. In general, PCANet-based methods is the second best.

In Table VI, we present the best threshold result of each changed detection methods by traversing all values. We could

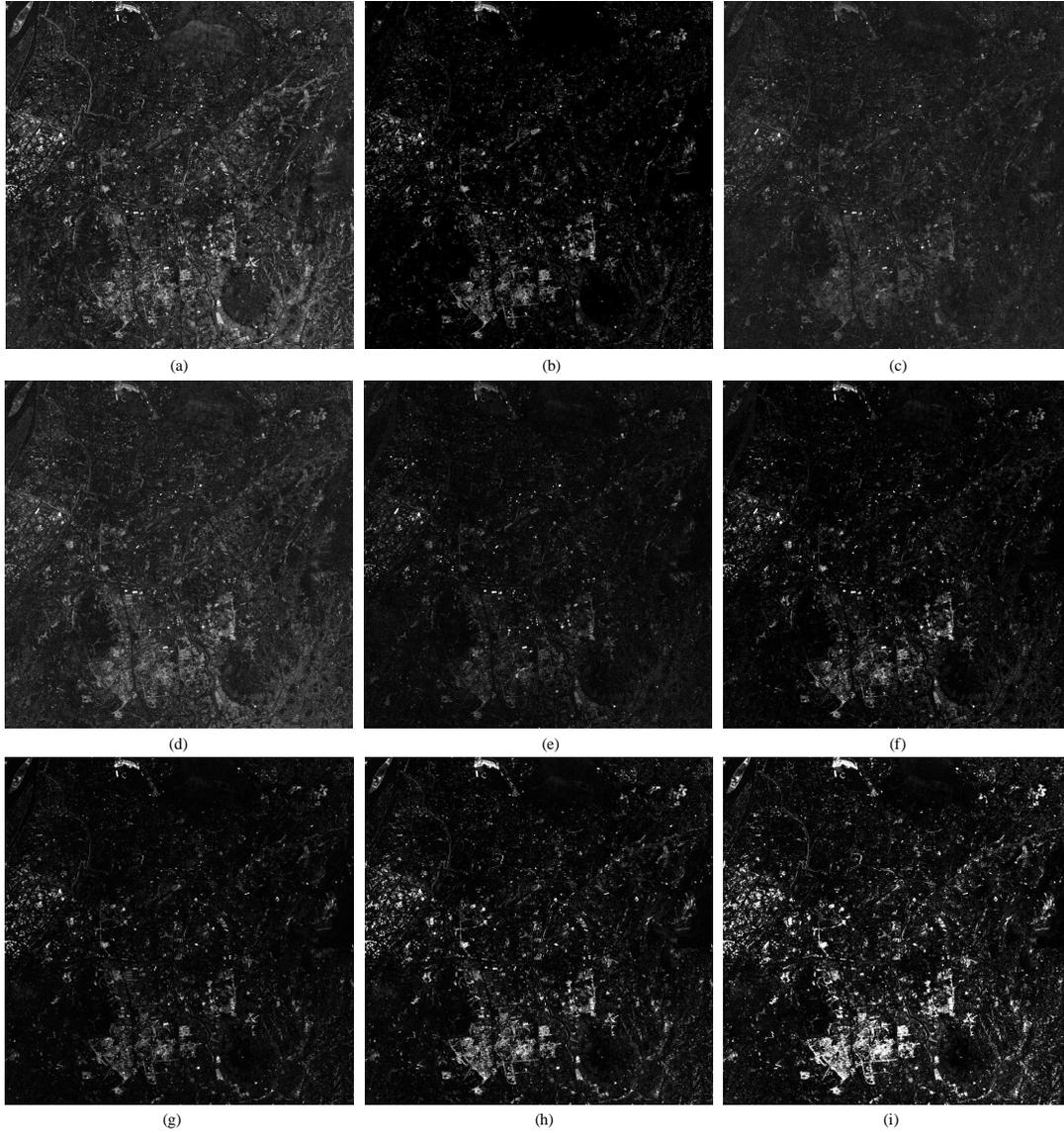


Fig. 7. Change intensity maps of Nanjing dataset by (a) CVA, (b) PCA, (c) MAD, (d) IRMAD, (e) USFA, (f) ISFA, (g) DSFA-64-2, (h) DSFA-128-2, and (i) DSFA-256-2.

TABLE V
CHANGE DETECTION RESULTS OF NANJING DATASET USING KMEANS.

Kmeans	OA_CHG	OA_UN	OA	Kappa	F1
CVA	0.8578	0.9184	0.9087	0.6958	0.7506
PCA	0.8650	0.9352	0.9240	0.7390	0.7846
MAD	0.9518	0.8557	0.8711	0.6276	0.7028
IRMAD	0.9564	0.8882	0.8991	0.6924	0.7523
USFA	0.5832	0.9692	0.9074	0.6159	0.6685
ISFA	0.6437	0.9760	0.9227	0.6833	0.7275
PCANet	0.8680	0.9334	0.9229	0.7367	0.7829
SDPCANet	0.8426	0.9397	0.9242	0.7351	0.7806
DSFA-64-2	0.7290	0.9817	0.9412	0.7647	0.7987
DSFA-128-2	0.7463	0.9807	0.9432	0.7748	0.8079
DSFA-256-2	0.7361	0.9792	0.9403	0.7632	0.7980

see from this table that DSFA-based methods are still the best on all the criteria. PCA, IRMAD and ISFA have high values on F1 score, but are much worse on OA and Kappa than DSFA-based methods. Besides, it's also worth noting that the best results of DSFA-based methods are very close to the results obtained by OTSU and Kmeans, which could be an evidence of the good discriminability of DSFA's results. On the contrary, threshold results and the best results of USFA and ISFA have a sensible difference. And the best results of CVA, PCA and MAD-based methods are also much better than their threshold results in both OA and Kappa coefficient.

Figure 8 shows the binary change maps of (a) CVA, (b) PCA, (c) MAD, (d) IRMAD, (e) USFA, (f) ISFA, (g) PCANet, (h) SDPCANet, (i) DSFA-64-2, (j) DSFA-128-2 and (k) DSFA-256-2, which are segmented by OTSU method. According to this figure, we could see that the binary change

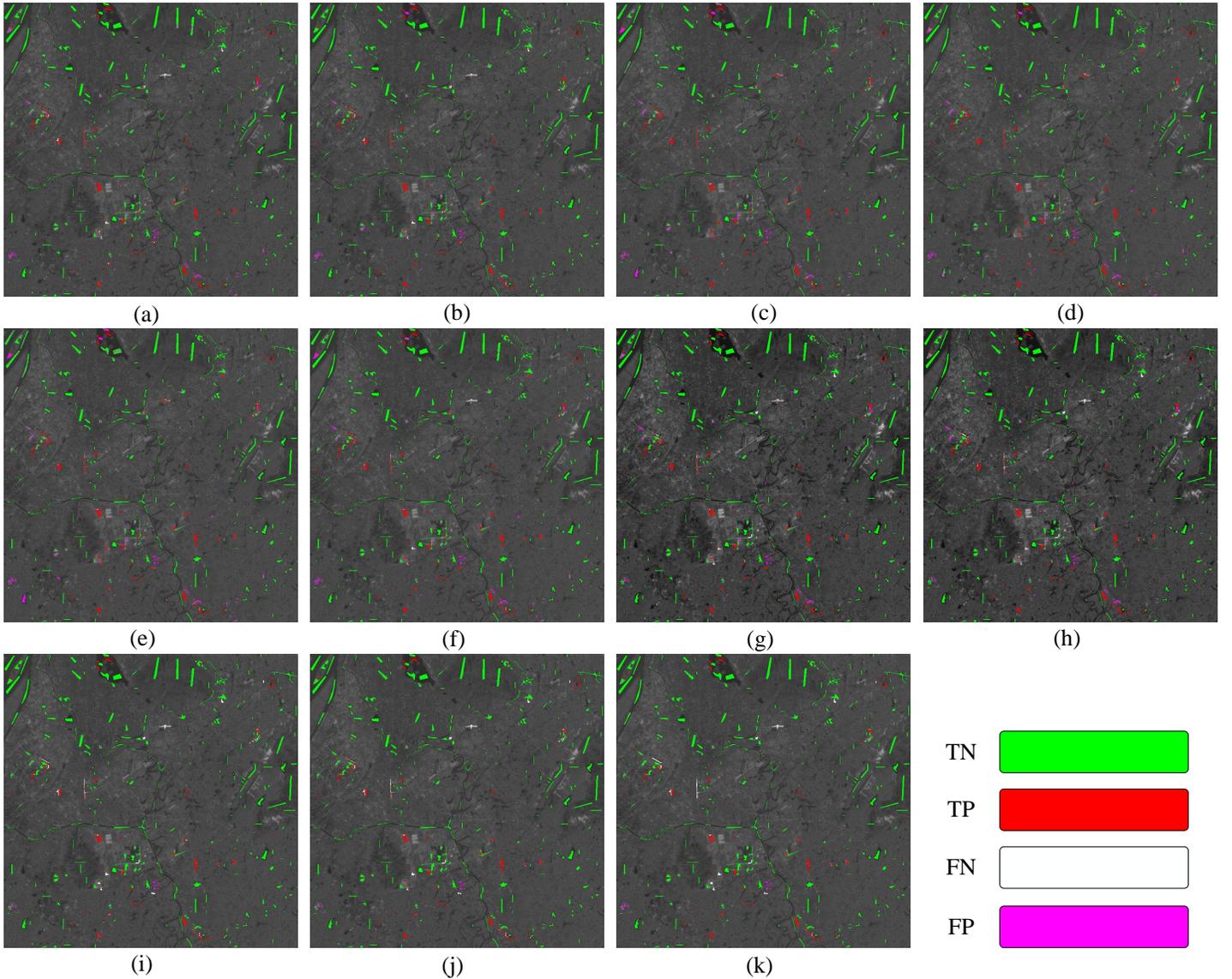


Fig. 8. The binary change maps of (a) CVA, (b) PCA, (c) MAD, (d) IRMAD, (e) USFA, (f) ISFA, (g) PCANet, (h) SDPCANet, (i) DSFA-64-2, (j) DSFA-128-2 and (k) DSFA-256-2.

TABLE VI
BEST CHANGE DETECTION RESULTS OF NANJING DATASET.

BEST	OA	Kappa	F1
CVA	0.9248	0.7178	0.7652
PCA	0.9341	0.7518	0.7925
MAD	0.9227	0.7244	0.7725
IRMAD	0.9229	0.7340	0.7815
USFA	0.9164	0.6997	0.7517
ISFA	0.9336	0.7578	0.7984
PCANet	0.9229	0.7367	0.7829
SDPCANet	0.9242	0.7351	0.7806
DSFA-64-2	0.9450	0.7915	0.8244
DSFA-128-2	0.9439	0.7850	0.8195
DSFA-256-2	0.9409	0.7664	0.8015

result of DSFA with different net structure are almost the same. Obviously, compared with DSFA's results, results of MAD and IRMAD have much more purple pixels, which represent the false positive samples. On the contrary, results of USFA and ISFA contain more false negative pixels, which are colored with white. The results of CVA and PCA are close to DSFA's results, but still has less true negative and more false positive samples than the latter. Besides, PCANet and SDPCANet also have a higher false positive rate than DSFA-based methods.

D. Experiments on River dataset

The River dataset consists of two 198 bands images with a spatial size of 463×241 . The changed regions of this dataset contain 12566 pixels, while the unchanged regions contain 99017 pixels. Figure 9 presents the bi-temporal images

and ground truth map of River dataset. In Figure 9, changed regions are white and unchanged regions are black.

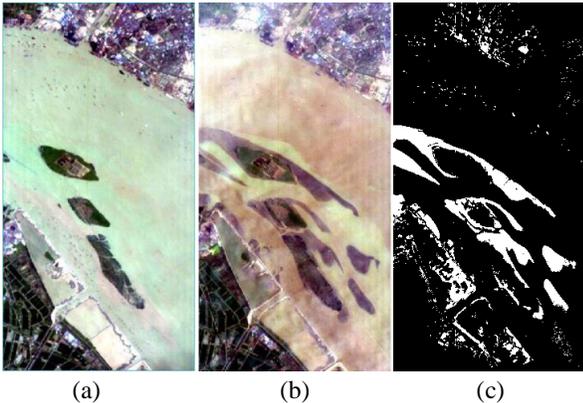


Fig. 9. The bi-temporal imagery of River dataset obtained in (a) May, 2013, (b) Dec, 2013, and (c) ground truth map.

In Figure 10, we present the change intensity maps of our proposed methods and all control methods. PCANet and SDPCANet are based on classification, so there's no relevant intensity maps in this figure. As can be observed from Figure 10, intuitively, all DSFA-based methods have better discriminability than CVA, PCA, and methods based on MAD and USFA. CVA, PCA and ISFA also have a better performance in separating the changed and unchanged regions than MAD, IRMAD and USFA. Visually, compared to the ground truth map, DSFA-based methods have relatively high false negative rate in the upper-right area of imagery. And other methods have brighter upper-right and lower-left area, which suggests that these methods are prone to detect these areas as changed, while most of them are unchanged actually.

We then use OTSU and Kmeans method to obtain the results with different criteria using the aforementioned methods. The obtained numerical results, along with the results of PCANet and SDPCANet, are presented in Table VII. The best value of each column is highlighted in bold in this table.

As could be observed from Table VII, DSFA-based method could achieve better performance on OA_UN, OA, Kappa and F1 score. Among all these methods, DSFA-128-2 has the best performance on OA, Kappa and F1 score, and the third best performance on OA_UN. DSFA-64-2 and SDPCANet both have the highest accuracy on OA_UN. Though PCANet have high performance on OA_CHG and F1 score, its performance on OA_UN, OA and Kappa are much worse than DSFA-based methods. In addition, it's also worth noting that the results using Kmeans and OTSU of our proposed methods still show very slight differences, which indicates that our proposed DSFA method are robust to different threshold methods.

The best results of each method are obtained by traversing all possible thresholds, and are presented in Table VIII. DSFA-based methods still have the best performance. Specifically, DSFA methods have much better performance on OA, Kappa and F1 score than other methods. Actually, DSFA-128-2 could outperform all other methods on all criteria. DSFA-64-2 and DSFA-256-2 respectively have the second and third best OA

and Kappa value, and they're very close to DSFA-128-2 on F1 score. In addition, the best values of DSFA methods are only slightly better than the results obtained with threshold methods, which also suggests that the transformed features of DSFA have a better discriminability.

In Figure 11, the binary change maps obtained by different methods are presented. Consistent with the results in Figure 10, DSFA algorithms have lower accuracies in detecting the changes in the upper-right regions of the original images, but have much better performance in other regions. The changes in the upper-right regions are not apparent and the background is complex, which we think is the main reason of DSFA's lower accuracy. On the contrary, CVA, PCA, MAD-based and SFA-based methods have a relatively high false positive rate in both the upper-right regions and lower-left regions. It's also noticed that SDPCANet also has a high false negative rate in the upper-right region, and PCANet tends to categorize the unchanged pixels in the lower-left regions as changed. On the whole, DSFA methods have the best performance visually and numerically.

E. Runtime Analysis

Though our proposed DSFA is based on fully connected networks, it's actually not very time consuming compared with other methods. We present the comparison of the runtime of IRMAD, ISFA, DSFA-128-2, PCANet and SDPCANet on three datasets in Figure 12. IRMAD and ISFA are implemented with MATLAB and run on CPU. PCANet and SDPCANet are also implemented with MATLAB but accelerated with 12 threads. DSFA-128-2 is implemented with Python and runs on CPU and GPU separately, which are respectively denoted by DSFA-CPU and DSFA-GPU in Figure 12. The CPU used is Intel Xeon E5 with a clock rate of 2.2 GHz. The GPU used is a single NVIDIA 1080Ti card.

As presents in this figure, ISFA and IRMAD are the two fastest methods, followed by DSFA-GPU and DSFA-CPU. Two PCANet-based method are the most time consuming. Besides, DSFA-GPU and DSFA-CPU are both faster than IRMAD and ISFA on River dataset, due to the smaller image size and more spectral bands of this dataset. On Taizhou and Nanjing dataset, the runtime of DSFA-GPU is very close to ISFA and IRMAD. DSFA-CPU is a little more time consuming, but it's still acceptable considering its improvements than IRMAD and ISFA.

V. DISCUSSION

A. Hyperparameter Analysis

In our experiments, we take 10^{-4} as the value of the regularization parameter r in Equation (22-23). However, in fact, r does not have significant influence on the final results when it's small enough.

In Figure 13, we present the relationship curves between the final change detection accuracy and r on three datasets. The network used is DSFA-128-2. It can be observed that when $r < 10^{-4}$, the accuracy curves on three dataset only have ignorable changes. On the contrary, when $r > 10^{-4}$, the accuracies are much lower because a larger r may corrupt the characteristic of the covariance matrices in Equation (22-23).

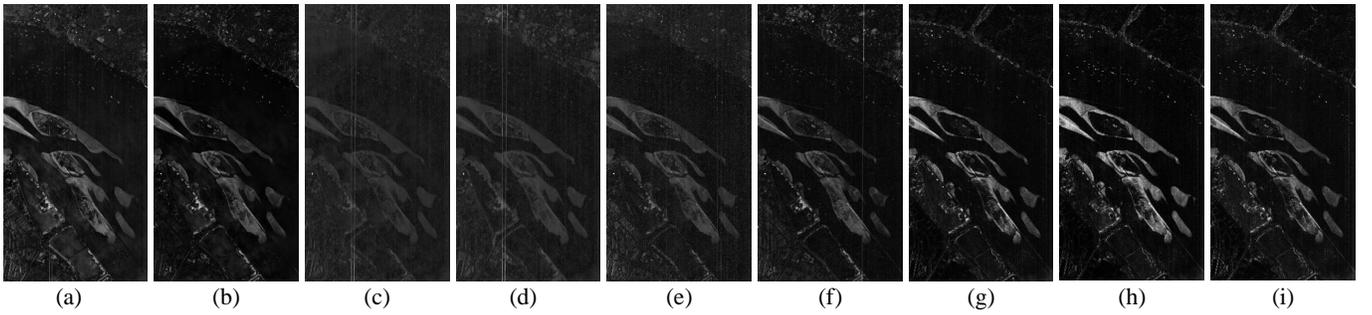


Fig. 10. The change intensity maps of (a) CVA, (b) PCA, (c) MAD, (d) IRMAD, (e) USFA, (f) ISFA, (g) DSFA-64-2, (h) DSFA-128-2, and (i) DSFA-256-2.

TABLE VII
CHANGE DETECTION RESULTS OF RIVER DATASET.

	Kmeans					OTSU				
	OA_CHG	OA_UN	OA	Kappa	F1	OA_CHG	OA_UN	OA	Kappa	F1
CVA	0.8168	0.9082	0.8979	0.5868	0.6432	0.8712	0.8770	0.8764	0.5474	0.6135
PCA	0.5899	0.9532	0.9123	0.5531	0.6024	0.5734	0.9560	0.9129	0.5484	0.5971
MAD	0.8022	0.9142	0.9016	0.5927	0.6474	0.8563	0.8864	0.8830	0.5591	0.6223
IRMAD	0.8093	0.9130	0.9013	0.5940	0.6488	0.8271	0.9059	0.8970	0.5872	0.6440
USFA	0.8297	0.8953	0.8879	0.5638	0.6250	0.8400	0.8871	0.8818	0.5514	0.6155
ISFA	0.6127	0.9377	0.9011	0.5267	0.5826	0.6377	0.9314	0.8984	0.5281	0.5856
PCANet	0.8024	0.9487	0.9322	0.6889	0.7273	0.8024	0.9487	0.9322	0.6889	0.7273
SDPCANet	0.5393	0.9850	0.9348	0.6166	0.6507	0.5393	0.9850	0.9348	0.6166	0.6507
DSFA-64-2	0.6164	0.9848	0.9434	0.6796	0.7293	0.6134	0.9851	0.9432	0.6780	0.7102
DSFA-128-2	0.6877	0.9812	0.9482	0.7207	0.7508	0.6864	0.9815	0.9483	0.7206	0.7494
DSFA-256-2	0.6622	0.9777	0.9422	0.6888	0.7283	0.6615	0.9778	0.9422	0.6884	0.7207

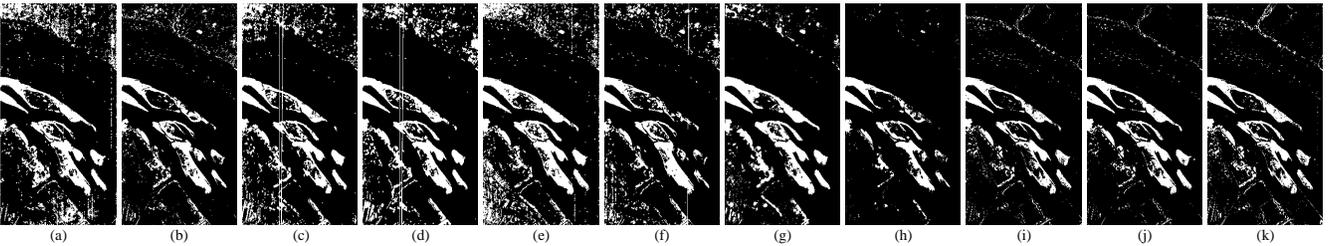


Fig. 11. The binary change maps of (a) CVA, (b) PCA, (c) MAD, (d) IRMAD, (e) USFA, (f) ISFA, (g) PCANet, (h) SDPCANet, (i) DSFA-64-2, (j) DSFA-128-2 and (k) DSFA-256-2.

B. Selection of Training Samples

In the Figure 14, we present the final accuracies using difference training sample selection strategies. This experiment is performed on the River dataset using DSFA-128-2. In Figure 14, Negative and Ground Truth strategy respectively mean that training samples are selected from the changed and unchanged regions of the ground truth image. Random strategy means training samples are absolutely randomly selected from the original imagery. And CVA strategy denotes that the training samples are selected from the unchanged regions of the change detection results of CVA.

As shown in this figure, Negative strategy leads to a very bad result, since the learned projection from changed pixel

pairs conflicts with the main idea of SFA and DSFA. Random strategy is very slightly better than CVA and Ground Truth strategy on OA_UN, but much worse on the other criteria. This because Random strategy will take quite a few changed pixel pairs as training samples, which would mislead the training process of DSFA. In addition, the results of CVA are almost the same with results of Ground Truth strategy, which indicates that DSFA with a simple pre-detection step to generate training samples could also achieve the same valid performance with using the Ground Truth. And in the field of change detection, labeling ground truth are usually hard and time consuming in both research and practical problems. Therefore, CVA is taken as the pre-detection method in our proposed algorithm.

TABLE VIII
 BEST CHANGE DETECTION RESULTS OF RIVER DATASET.

BEST	OA	Kappa	F1
CVA	0.9264	0.6242	0.6841
PCA	0.9204	0.6075	0.6641
MAD	0.9140	0.5972	0.6481
IRMAD	0.9095	0.5984	0.6510
USFA	0.9180	0.6098	0.6590
ISFA	0.9098	0.5285	0.5879
PCANet	0.9322	0.6889	0.7273
SDPCANet	0.9348	0.6166	0.6507
DSFA-64-2	0.9454	0.7109	0.7419
DSFA-128-2	0.9483	0.7270	0.7566
DSFA-256-2	0.9423	0.7007	0.7344

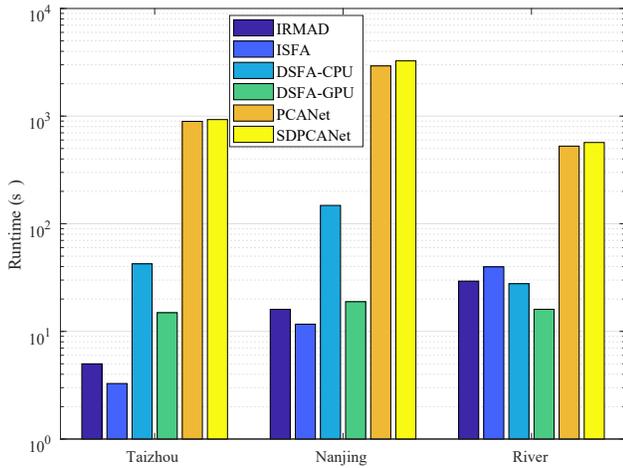


Fig. 12. The comparison of runtime of different change detection methods.

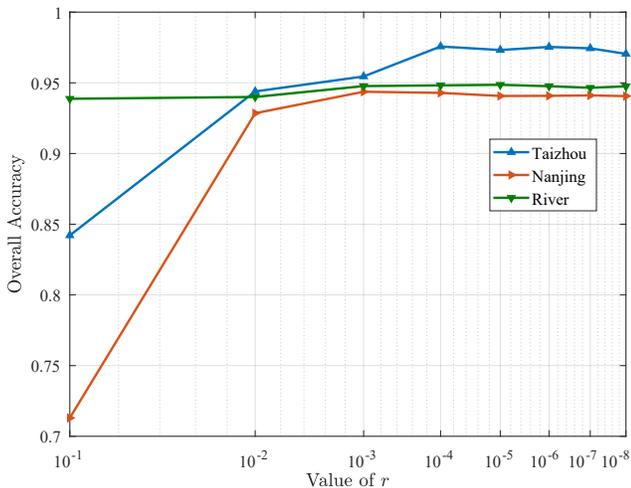
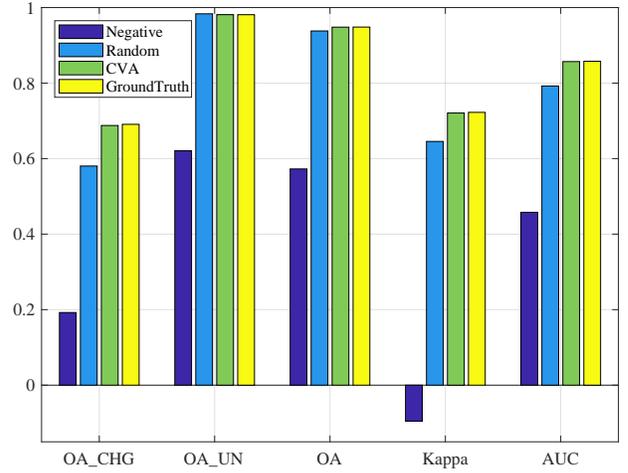

 Fig. 13. The relationship between r and OA on three datasets.


Fig. 14. Comparison of different training sample selection strategies on River dataset.

VI. CONCLUSION

In this paper, we proposed a novel change detection algorithm called DSFA for multi-temporal remote sensing images. In the DSFA model, two deep networks are used to project the bi-temporal original input data into a new feature space. Then, SFA is used to extract the most invariant components of unchanged pixels and suppress them in changed regions to highlight changed components. We formulated the SFA process and loss function of DSFA model, and presented the derivation of computing gradient of loss. Our proposed algorithm is unsupervised, which means it doesn't need priori labeled pixels for the training process.

We implemented our algorithm and performed experiments on two multi-spectral datasets and a public hyperspectral dataset. The visual and quantitative results have both shown that our method could outperform the other state-of-the-art methods, including other SFA-based and deep network algorithms.

Our proposed method currently focuses on differentiating the changed and unchanged regions in bi-temporal remote sensing imagery. The future work is required to explore DSFA's potential in detecting multi-classes changes. And in consideration of that SFA is originally designed for solving the problems of continuous signals, it will be promising to develop a specific DSFA model for change detection of sequent or video imagery.

APPENDIX A

DERIVATION OF GRADIENT OF LOSS

Here we will present the detailed deduction process of computing the gradient of $\mathcal{L}(\theta_1, \theta_2)$ with respect to \hat{X}_ϕ . Based on the reference [43], we have the following equations.

$$\frac{\partial \text{tr}(ABA^T C)}{\partial A} = CAB + C^T AB^T, \quad (36)$$

$$\frac{\partial (X^{-1})_{kl}}{\partial X_{ij}} = -(X^{-1})_{ki} (X^{-1})_{jl}. \quad (37)$$

Based on (36) and the fact that A_ϕ and B_ϕ are both symmetric, we could obtain:

$$\nabla_A = \frac{\partial \mathcal{L}(\theta_1, \theta_2)}{\partial A_\phi} = 2B_\phi^{-1} A_\phi B_\phi^{-1}, \quad (38)$$

$$\begin{aligned} \frac{\partial \mathcal{L}(\theta_1, \theta_2)}{\partial B_\phi^{-1}} &= 2A_\phi B_\phi^{-1} A_\phi \\ \Leftrightarrow \left(\frac{\partial \mathcal{L}(\theta_1, \theta_2)}{\partial B_\phi^{-1}} \right)_{kl} &= 2(A_\phi B_\phi^{-1} A_\phi)_{kl}. \end{aligned} \quad (39)$$

Then, combining (37), $\nabla_B = \partial \mathcal{L}(\theta_1, \theta_2) / \partial B_\phi$ is calculated as the following equation:

$$\begin{aligned} \nabla_B &= -2 \sum_{kl} (A_\phi B_\phi^{-1} A_\phi)_{kl} (B_\phi^{-1})_{ki} (B_\phi^{-1})_{jl} \\ &= -2 \sum_{kl} (B_\phi^{-1})_{ik} (A_\phi B_\phi^{-1} A_\phi)_{kl} (B_\phi^{-1})_{lj} \\ &= -2 (B_\phi^{-1} A_\phi B_\phi^{-1} A_\phi B_\phi^{-1}). \end{aligned} \quad (40)$$

We could expand the expression of A_ϕ out:

$$\begin{aligned} A_\phi &= \Sigma_{XY} = \frac{1}{n} (\hat{X}_\phi - \hat{Y}_\phi) (\hat{X}_\phi - \hat{Y}_\phi)^T \\ &= \frac{1}{n} (\hat{X}_\phi \hat{X}_\phi^T + \hat{Y}_\phi \hat{Y}_\phi^T - \hat{X}_\phi \hat{Y}_\phi^T - \hat{Y}_\phi \hat{X}_\phi^T). \end{aligned} \quad (41)$$

First, based on the derivation in the appendix of [42], we have:

$$\begin{aligned} \frac{\partial (\hat{X}_\phi \hat{X}_\phi^T)^{ab}}{\partial \hat{X}_\phi^{ij}} &= \begin{cases} \frac{2}{n} (\hat{X}_\phi^{ij} - \frac{1}{n} \sum_k \hat{X}_\phi^{ik}), a = i, b = i \\ \frac{1}{n} (\hat{X}_\phi^{bj} - \frac{1}{n} \sum_k \hat{X}_\phi^{bk}), a = i, b \neq i \\ \frac{1}{n} (\hat{X}_\phi^{aj} - \frac{1}{n} \sum_k \hat{X}_\phi^{ak}), a \neq i, b = i \\ 0, a \neq i, b \neq i \end{cases} \\ &= \frac{1}{n} (\xi_{(a=i)} \hat{X}_\phi^{bj} + \xi_{(b=i)} \hat{X}_\phi^{aj}). \end{aligned} \quad (42)$$

Also,

$$\frac{\partial (\hat{X}_\phi \hat{Y}_\phi^T)^{ab}}{\partial \hat{X}_\phi^{ij}} = \frac{1}{n} (\hat{Y}_\phi^{bj} - \frac{1}{n} \sum_k \hat{Y}_\phi^{bk}) = \frac{1}{n} \xi_{(a=i)} \hat{Y}_\phi^{bj}. \quad (43)$$

Integrating (42) and (43) into (41):

$$\begin{aligned} \frac{\partial A_\phi^{ab}}{\partial \hat{X}_\phi^{ij}} &= \frac{\partial (\hat{X}_\phi \hat{X}_\phi^T)^{ab}}{\partial \hat{X}_\phi^{ij}} - \frac{\partial (\hat{Y}_\phi \hat{X}_\phi^T)^{ab}}{\partial \hat{X}_\phi^{ij}} - \frac{\partial (\hat{X}_\phi \hat{Y}_\phi^T)^{ab}}{\partial \hat{X}_\phi^{ij}} \\ &= \frac{1}{n} (\xi_{(a=i)} \hat{X}_\phi^{bj} + \xi_{(b=i)} \hat{X}_\phi^{aj}) \\ &\quad - \frac{1}{n} (\xi_{(b=i)} \hat{Y}_\phi^{aj} + \xi_{(a=i)} \hat{Y}_\phi^{bj}). \end{aligned} \quad (44)$$

Similarly, with respect to B_ϕ , we have:

$$\begin{aligned} \frac{\partial B_\phi^{ab}}{\partial \hat{X}_\phi^{ij}} &= \frac{\partial \Sigma_{XX}^{ab} + \partial \Sigma_{YY}^{ab}}{2 \partial \hat{X}_\phi^{ij}} \\ &= \frac{1}{2n} (\xi_{(a=i)} \hat{X}_\phi^{bj} + \xi_{(b=i)} \hat{X}_\phi^{aj}). \end{aligned} \quad (45)$$

Putting (44) and (45) together, the gradient of $\mathcal{L}(\theta_1, \theta_2)$ with respect to \hat{X}_ϕ^{ij} is then computed as:

$$\begin{aligned} \frac{\partial \mathcal{L}(\theta_1, \theta_2)}{\partial \hat{X}_\phi^{ij}} &= \sum_{ab} \nabla_A^{ab} \frac{\partial A_\phi^{ab}}{\partial \hat{X}_\phi^{ij}} + \sum_{ab} \nabla_B^{ab} \frac{\partial B_\phi^{ab}}{\partial \hat{X}_\phi^{ij}} \\ &= \frac{1}{n} (\sum_b \nabla_A^{ib} \hat{X}_\phi^{bj} + \sum_a \nabla_A^{ai} \hat{X}_\phi^{aj}) \\ &\quad - \frac{1}{n} (\sum_b \nabla_A^{ib} \hat{Y}_\phi^{bj} + \sum_a \nabla_A^{ai} \hat{Y}_\phi^{aj}) \\ &\quad + \frac{1}{2n} (\sum_b \nabla_B^{ib} \hat{X}_\phi^{bj} + \sum_a \nabla_B^{ai} \hat{X}_\phi^{aj}) \\ &= \frac{1}{n} (\nabla_A \hat{X}_\phi + \nabla_A^T \hat{X}_\phi - \nabla_A \hat{Y}_\phi - \nabla_A^T \hat{Y}_\phi)_{ij} \\ &\quad + \frac{1}{2n} (\nabla_B \hat{X}_\phi + \nabla_B^T \hat{X}_\phi)_{ij}. \end{aligned} \quad (46)$$

Obviously, ∇_A and ∇_B are both symmetric matrices. Therefore,

$$\frac{\partial \mathcal{L}(\theta_1, \theta_2)}{\partial \hat{X}_\phi^{ij}} = \frac{2}{n} (\nabla_A \hat{X}_\phi - \nabla_A \hat{Y}_\phi)_{ij} + \frac{1}{n} (\nabla_B \hat{X}_\phi)_{ij}. \quad (47)$$

Finally, we could obtain the gradient of $\mathcal{L}(\theta_1, \theta_2)$ with respect to \hat{X}_ϕ :

$$\frac{\partial \mathcal{L}(\theta_1, \theta_2)}{\partial \hat{X}_\phi} = \frac{2}{n} (\nabla_A \hat{X}_\phi - \nabla_A \hat{Y}_\phi) + \frac{1}{n} \nabla_B \hat{X}_\phi. \quad (48)$$

REFERENCES

- [1] A. Singh, "Review article digital change detection techniques using remotely-sensed data," *International journal of remote sensing*, vol. 10, no. 6, pp. 989–1003, 1989.
- [2] B. Du, L. Zhang, D. Tao, and D. Zhang, "Unsupervised transfer learning for target detection from hyperspectral images," *Neurocomputing*, vol. 120, pp. 72–82, 2013.
- [3] L. Zhang, L. Zhang, D. Tao, and X. Huang, "Tensor discriminative locality alignment for hyperspectral image spectral–spatial feature extraction," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 51, no. 1, pp. 242–256, 2013.
- [4] B. Du, L. Zhang, L. Zhang, T. Chen, and K. Wu, "A discriminative manifold learning based dimension reduction method for hyperspectral classification," *International Journal of Fuzzy Systems*, vol. 14, no. 2, pp. 272–277, 2012.
- [5] G. Xian, C. Homer, and J. Fry, "Updating the 2001 national land cover database land cover classification to 2006 by using landsat imagery change detection methods," *Remote Sensing of Environment*, vol. 113, no. 6, pp. 1133–1147, 2009.
- [6] G. Xian and C. Homer, "Updating the 2001 national land cover database impervious surface products to 2006 using landsat imagery change detection methods," *Remote Sensing of Environment*, vol. 114, no. 8, pp. 1676–1686, 2010.
- [7] P. R. Coppin and M. E. Bauer, "Digital change detection in forest ecosystems with remote sensing imagery," *Remote sensing reviews*, vol. 13, no. 3–4, pp. 207–234, 1996.
- [8] R. E. Kennedy, P. A. Townsend, J. E. Gross, W. B. Cohen, P. Bolstad, Y. Wang, and P. Adams, "Remote sensing change detection tools for natural resource managers: Understanding concepts and tradeoffs in the design of landscape monitoring projects," *Remote sensing of environment*, vol. 113, no. 7, pp. 1382–1396, 2009.
- [9] M. A. Wulder, C. R. Butson, and J. C. White, "Cross-sensor change detection over a forested landscape: Options to enable continuity of medium spatial resolution measures," *Remote Sensing of Environment*, vol. 112, no. 3, pp. 796–809, 2008.
- [10] H. Luo, C. Liu, C. Wu, and X. Guo, "Urban change detection based on dempster–shafer theory for multitemporal very high-resolution imagery," *Remote Sensing*, vol. 10, no. 7, p. 980, 2018.

- [11] M. K. Ridd and J. Liu, "A comparison of four algorithms for change detection in an urban environment," *Remote sensing of environment*, vol. 63, no. 2, pp. 95–100, 1998.
- [12] L. Bruzzone and D. F. Prieto, "Automatic analysis of the difference image for unsupervised change detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 38, no. 3, pp. 1171–1182, 2000.
- [13] J. Deng, K. Wang, Y. Deng, and G. Qi, "Pca-based land-use change detection and analysis using multitemporal and multisensor satellite data," *International Journal of Remote Sensing*, vol. 29, no. 16, pp. 4823–4838, 2008.
- [14] A. A. Nielsen, K. Conradsen, and J. J. Simpson, "Multivariate alteration detection (mad) and maf postprocessing in multispectral, bitemporal image data: New approaches to change detection studies," *Remote Sensing of Environment*, vol. 64, no. 1, pp. 1–19, 1998.
- [15] A. A. Nielsen and K. Conradsen, "Multivariate alteration detection (mad) in multispectral, bi-temporal image data: A new approach to change detection studies," 1997.
- [16] J. B. Collins and C. E. Woodcock, "An assessment of several linear change detection techniques for mapping forest mortality using multitemporal landsat tm data," *Remote Sensing of Environment*, vol. 56, no. 1, pp. 66–77, 1996.
- [17] S. Marchesi and L. Bruzzone, "Ica and kernel ica for change detection in multispectral remote sensing images," in *Geoscience and Remote Sensing Symposium, 2009 IEEE International, IGARSS 2009*, vol. 2. IEEE, 2009, pp. II–980.
- [18] F. Bovolo, L. Bruzzone, and M. Marconcini, "A novel approach to unsupervised change detection based on a semisupervised svm and a similarity measure," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 46, no. 7, pp. 2070–2082, 2008.
- [19] Z. Huang, X. Jia, and L. Ge, "Sampling approaches for one-pass land-use/land-cover change mapping," *International Journal of Remote Sensing*, vol. 31, no. 6, pp. 1543–1554, 2010.
- [20] B. Demir, F. Bovolo, and L. Bruzzone, "Detection of land-cover transitions in multitemporal remote sensing images with active-learning-based compound classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 50, no. 5, pp. 1930–1941, 2012.
- [21] O. Ahlqvist, "Extending post-classification change detection using semantic similarity metrics to overcome class heterogeneity: A study of 1992 and 2001 us national land cover database changes," *Remote Sensing of Environment*, vol. 112, no. 3, pp. 1226–1241, 2008.
- [22] T. Celik and K.-K. Ma, "Multitemporal image change detection using undecimated discrete wavelet transform and active contours," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 49, no. 2, pp. 706–716, 2011.
- [23] D. Liu, K. Song, J. R. Townshend, and P. Gong, "Using local transition probability models in markov random fields for forest change detection," *Remote Sensing of Environment*, vol. 112, no. 5, pp. 2222–2231, 2008.
- [24] Z. Yetgin, "Unsupervised change detection of satellite images using local gradual descent," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 50, no. 5, pp. 1919–1929, 2012.
- [25] L. Gueguen, P. Soille, and M. Pesaresi, "Change detection based on information measure," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 49, no. 11, pp. 4503–4515, 2011.
- [26] P. Coppin, I. Jonckheere, K. Nackaerts, B. Muys, and E. Lambin, "Review article digital change detection methods in ecosystem monitoring: a review," *International journal of remote sensing*, vol. 25, no. 9, pp. 1565–1596, 2004.
- [27] D. Lu, P. Mausel, E. Brondizio, and E. Moran, "Change detection techniques," *International journal of remote sensing*, vol. 25, no. 12, pp. 2365–2401, 2004.
- [28] C. Wu, B. Du, and L. Zhang, "Slow feature analysis for change detection in multispectral imagery," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 52, no. 5, pp. 2858–2874, 2014.
- [29] C. Wu, B. Du, X. Cui, and L. Zhang, "A post-classification change detection method based on iterative slow feature analysis and bayesian soft fusion," *Remote Sensing of Environment*, vol. 199, pp. 241–255, 2017.
- [30] L. Wiskott and T. J. Sejnowski, "Slow feature analysis: Unsupervised learning of invariances," *Neural computation*, vol. 14, no. 4, pp. 715–770, 2002.
- [31] L. Wiskott, P. Berkes, M. Franzius, H. Sprekeler, and N. Wilbert, "Slow feature analysis," *Scholarpedia*, vol. 6, no. 4, p. 5282, 2011.
- [32] Z. Zhang and D. Tao, "Slow feature analysis for human action recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, no. 3, pp. 436–450, 2012.
- [33] L. Sun, K. Jia, T.-H. Chan, Y. Fang, G. Wang, and S. Yan, "Dl-sfa: deeply-learned slow feature analysis for action recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 2625–2632.
- [34] M. Franzius, N. Wilbert, and L. Wiskott, "Invariant object recognition with slow feature analysis," in *International Conference on Artificial Neural Networks*. Springer, 2008, pp. 961–970.
- [35] —, "Invariant object recognition and pose estimation with slow feature analysis," *Neural computation*, vol. 23, no. 9, pp. 2289–2323, 2011.
- [36] C. Wu, L. Zhang, and B. Du, "Kernel slow feature analysis for scene change detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 4, pp. 2367–2384, 2017.
- [37] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle, "Greedy layer-wise training of deep networks," in *Advances in neural information processing systems*, 2007, pp. 153–160.
- [38] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [39] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [40] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in neural information processing systems*, 2015, pp. 91–99.
- [41] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.
- [42] G. Andrew, R. Arora, J. Bilmes, and K. Livescu, "Deep canonical correlation analysis," in *International Conference on Machine Learning*, 2013, pp. 1247–1255.
- [43] K. B. Petersen and M. S. Pedersen, "The matrix cookbook," nov 2012, version 20121115. [Online]. Available: <http://www2.imm.dtu.dk/pubdb/p.php?3274>
- [44] A. A. Nielsen, "The regularized iteratively reweighted mad method for change detection in multi-and hyperspectral data," *IEEE Transactions on Image Processing*, vol. 16, no. 2, pp. 463–478, 2007.
- [45] F. Gao, J. Dong, B. Li, and Q. Xu, "Automatic Change Detection in Synthetic Aperture Radar Images Based on PCANet," *IEEE Geoscience and Remote Sensing Letters*, vol. 13, no. 12, pp. 1792–1796, dec 2016.
- [46] M. Li, M. Li, P. Zhang, Y. Wu, W. Song, and L. An, "SAR Image Change Detection Using PCANet Guided by Saliency Detection," *IEEE Geoscience and Remote Sensing Letters*, vol. 16, no. 3, pp. 402–406, mar 2019.
- [47] D. R. Hardoon, S. Szedmak, and J. Shawe-Taylor, "Canonical correlation analysis: An overview with application to learning methods," *Neural computation*, vol. 16, no. 12, pp. 2639–2664, 2004.
- [48] T. H. Chan, K. Jia, S. Gao, J. Lu, Z. Zeng, and Y. Ma, "PCANet: A Simple Deep Learning Baseline for Image Classification?" *IEEE Transactions on Image Processing*, vol. 24, no. 12, pp. 5017–5032, dec 2015.
- [49] S. Goferman, L. Zelnik-Manor, and A. Tal, "Context-aware saliency detection," *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 10, pp. 1915–1926, 2011.