

# Full parameter time complexity (FPTC): A method to evaluate the running time of machine learning classifiers for land use/land cover classification

Xiaorou Zheng<sup>a,b</sup>, Jianxin Jia<sup>a,c,d</sup>, Shanxin Guo<sup>a,c,\*</sup>, Jinsong Chen<sup>a,c</sup>, Luyi Sun<sup>a,c</sup>, Yingfei Xiong<sup>a,b</sup> and Wennan Xu<sup>a,b</sup>

**Abstract**—In emergency responses to natural disasters, actionable information provided by remote sensing images is crucial to help emergency managers become aware of the situation and assess the magnitude of the damage. Without the accurate prediction of time consumption, choosing an algorithm for land use/land cover (LULC) classification under these emergency circumstances could be blind and subjective. Here, we proposed a full parameter time complexity (FPTC) analysis and the corresponding coefficient  $\omega$  to estimate the actual running time of the LULC classification without actually running the code. The FPTC of five general algorithms is derived in this study. After derivation, the FPTC of  $k$ -nearest neighbors ( $k$ NN) is  $F(nv + n\log_2 u)$ , the FPTC of logistic regression (LR) is  $F(Qm^2vn)$ , the FPTC of classification and regression tree (CART) is  $F((m+1)nv\log_2 n)$ , the FPTC of random forest (RF) is  $F(s(m+1)nv\log_2 n)$ , and the FPTC of support vector machine (SVM) is  $F(m^2Qv(n+k))$ . The results show a strong linear relationship between the actual running time and FPTC (R-squared:  $k$ NN (0.991), LR (0.997), CART (0.999), RF (1.000), and SVM (0.999)), with different data size. The average root-mean-squared error between the real running time and the estimated running time is 3.34 s, which demonstrates the effectiveness of FPTC. Combining FPTC with the corresponding coefficient  $\omega$ , the running time of the classification can be precisely predicted, which will help emergency managers quickly choose algorithms in response to natural disasters with available remote sensing data and limited time.

**Index Terms**—full parameter time complexity (FPTC); traditional time complexity (TTC); algorithm running time; land use/land cover classification; Sentinel-2A

## I. INTRODUCTION

**A**SSESSMENTS of natural disasters and risk are the foundation of decision-making processes for a wide variety of actors from the public to government emergency managers. Quickly quantifying damage and expected future losses is usually the first step to becoming aware of the current situation [1]–[3]. The land use/land cover (LULC) products from remote sensing imagery can provide first-hand information for this purpose [4]–[6]. As a result that this decision-making process is usually urgent, choosing an appropriate classification

The authors are with<sup>a</sup> Center for Geo-Spatial Information, Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen 518055, China; <sup>b</sup> University of Chinese Academy of Sciences, Beijing 100049, China; <sup>c</sup> Shenzhen Engineering Laboratory of Ocean Environmental Big Data Analysis and Application, Shenzhen 518055, China; <sup>d</sup> Department of Photogrammetry and Remote Sensing, Finnish Geospatial Research Institute, Kruunukummi FI-02430, Finland. (Corresponding author: Shanxin Guo. Email:sx.guo@siat.ac.cn)

algorithm to achieve this goal with limited time and resources can be challenging [7]. In addition to classification accuracy, the actual time consumption of the algorithm is another aspect that needs to be carefully evaluated before running the task [4], [7], [8]. Without an accurate prediction of time consumption, choosing an algorithm for LULC classification under these emergency circumstances could be blind and subjective.

In general, the methods to estimate the time consumption of a classification task can be divided into two categories: 1) sampled-data-based methods and 2) time-complexity-based methods [9], [10]. The first category involves estimating based on running the program and launching the time-calculated function on the sampled dataset. These methods hold an assumption that the actual running time between the sample and the whole dataset can be simplified by a linear or nonlinear relationship. Although these methods are used in various studies [8], [10], [11], the drawbacks of these methods are that 1) this linear or nonlinear relationship highly depends on the hardware, which cannot be generalized across different computing environments, and 2) the effect of the different parameters of the algorithm (both operational and hidden parameters) is treated as a black box. The influence of these parameters remains unknown.

The second category involves estimating based on a time complexity analysis. The commonly used asymptotic time complexity belongs to this category, usually referred to as traditional time complexity (TTC) [9]. TTC is a function of input size (e.g., the number of samples), which measures the computational complexity with the input size increase under the iterations of unit operation (e.g., addition or multiplication). The time of each unit operation is assumed to have the same value, so the iterations can be estimated as proportional to the running time [9], [12]. In the process of calculating the TTC, many low-order details have been ignored. For instance, when we calculate TTC for a function  $f(n) = an^2 + bn + c$  (where  $n$  represents the data size), only  $f(n) = n^2$  is of concern [9], [12], [13]. As a result,  $O(n^2)$  is used to estimate the upper boundary of the time complexity of this function. The purpose of the TTC is to capture the acceleration of the running time as an increase in the data size ( $N$ ) at the theoretical level; this can hardly be used for the accurate prediction of the running time, especially for remote sensing LULC classification tasks, where the time consumption not only relates to the data size ( $N$ ) but also to other parameters (e.g., the number of available bands and the number of support

vectors). How to consider the influence of these parameters to predict the overall time consumption remains a challenge.

In fact, without considering the physical discrepancy between the different platforms (such as CPU and GPU), the time consumption of a classification algorithm can be influenced by 1) the date size, 2) the number of classes, 3) the number of bands/features, 4) the iteration structure of the algorithm, 5) the operational parameters of the algorithm, such as the number of trees in random forest, and 6) the hidden parameters of the algorithm, such as the number of support vectors. All these components affect the actual time consumption of the algorithm in different ways via unknown mechanisms. Determining how the contribution of each component can be quantified is key to predicting the actual time that will be consumed.

To address the above issues, we propose full parameter time complexity (FPTC), which takes all time-consuming parameters into account. The FPTC of five general algorithms—*k*-nearest neighbors (*k*NN) [14], logistic regression (LR) [15], classification and regression tree (CART) [16], random forest (RF) [17], and support vector machine (SVM) [18]—is derived in this study. We defined a coefficient  $\omega$  to model the physical discrepancy between different platforms for different classifiers. To test the effectiveness of FPTC and the corresponding coefficient  $\omega$ , the Xinjiang Uygur Autonomous Region, China, and the Sentinel-2A dataset were chosen as a case study. The results show that the running time of the classification task can be precisely predicted by combining FPTC with coefficient  $\omega$ . These will help emergency managers make quick decisions in response to natural disasters. Our contribution can be summarized as follows:

- 1) We propose a method to quantitatively evaluate the time efficiency of a machine learning classifier (FPTC) and derive the FPTC of five general algorithms: *k*NN, LR, CART, RF, and SVM.
- 2) To predict the time consumption, we propose the coefficient  $\omega$ , which is used to establish the relationship between running time and FPTC. The coefficient  $\omega$  can be easily obtained with the pre-experiment with a small sampled dataset under different computing environments.
- 3) For the parameters that cannot be estimated before running the algorithm, we analyze the relationship between these parameters and those easily obtained hyperparameters to predict the actual running time without actually running the algorithm.

The remainder of this study is organized as follows. In Section 2, FPTC is defined, and the corresponding mathematical derivation is described in detail. Sections 3 and 4 present the materials and the experimental result of the Xinjiang dataset. Section 5 concludes with a summary.

## II. FULL PARAMETER TIME COMPLEXITY

### A. Definition

FPTC is defined as two components. One is  $F(n, m, v, \theta')$ , the algorithm-related part, which can be derived based on the structural analysis for one particular classifier. This part is a

function of  $n, m, v$  and  $\theta'$ , where  $n$  represents the sample size,  $m$  represents the number of targeted classes,  $v$  represents the number of bands/features of the remote sensing image, and  $\theta'$  represents a collection of parameters related to the algorithm. It should be noticed that the  $\theta'$  can be different for different algorithms. For instance,  $\theta'$  of the FPTC in *k*NN consists of the number of nearest neighbors  $u$ , while the  $\theta'$  of the FPTC in SVM consists of the number of iterations  $Q$  and the number of support vectors  $k$ . The second component is the coefficient  $\omega$ , which is a physically related part reflecting the computing environmental factors, such as the speed of the CPU/GPU or the RAM. Therefore, the coefficient  $\omega$  may vary depending on the platform. Usually, a pre-experiment on a small part of the dataset can help us to evaluate this coefficient for a specific classifier. Combining these two parts, the definition of FPTC is as follows:

$$t^* = F(n, m, v, \theta') \quad (1)$$

$$t' = \omega \times t^* \quad (2)$$

where  $t'$  is the real running time,  $\omega$  is the coefficient, and  $t^*$  is the time estimated by structural analysis. In the next section, we will derive the algorithm-related part of FPTC for five classically and commonly used classifiers in the remote sensing field. We derive the FPTC of the selected algorithms in the following order: *k*NN, LR, CART, RF, and SVM.

### B. Deriving FPTC for *k*NN

The *k*NN classifier [14] memorizes the entire training data  $T_r$  and performs classification only if the test data  $x^{(e)} \in T_e$  are given, in which  $e \in \{n+1, n+2, \dots, n+b\}$ . The classifiers compute the distance or similarity between the test data  $x^{(e)}$  and each training sample. The classifier then found a group of  $u$  objects in the training set  $T_r$  that were closest to the test object  $x^{(e)}$ . In this process, the commonly used Manhattan distance is equivalent to the Minkowskian  $r$ -distance function with  $r = 1$ , adopted as the distance or similarity metric with the following [19]:

$$D^{(i)} = \sum_{j=1}^v |x_j^{(i)} - x_j^{(e)}| \quad (3)$$

where  $D^{(i)}$  is the distance between the training sample  $(x^{(i)}, y^{(i)})$  and the test data  $T_e$ , in which  $i \in \{1, 2, 3, \dots, n\}$ .

The *k*NN classifier is a lazy learner, which means that the cost of building the model is cheap, but classifying the test samples is relatively expensive [20]. To calculate the complexity over testing samples, the FPTC of *k*NN classifiers is divided into two parts (Figure 1). First, the FPTC of calculating the distance between one training sample and test data  $x^{(e)}$  is  $F(v)$  (Equation (3)). When there are  $n$  testing samples that need to be classified, the FPTC becomes  $F(vn)$ . Second, training samples with a minimum distance should be selected from the training set. In this classic optimal searching algorithm, the FPTC becomes  $F(n \log_2 u)$ . Therefore, the total

FPTC of  $k$ NN is  $F(nv + n\log_2 u)$ . Considering the corresponding coefficient  $\omega_{kNN}$ , the FPTC of  $k$ NN is associated with the real running time  $t'_{kNN}$  as follows.

$$t'_{kNN} = \omega_{kNN} \times t^*_{kNN} = \omega_{kNN} \times F(nv + n\log_2 u) \quad (4)$$

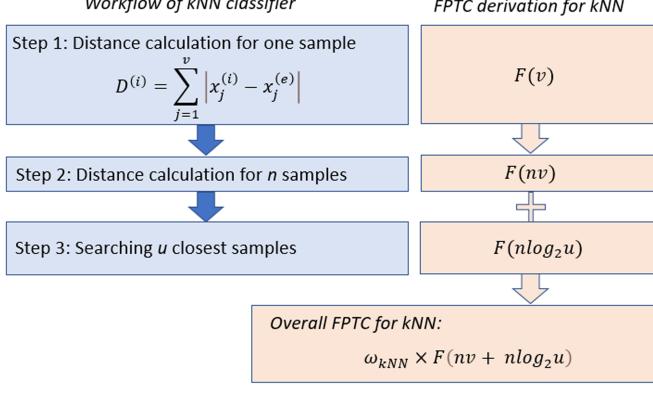


Fig. 1. Deriving Full parameter time complexit (FPTC) for  $k$ -nearest neighbors (kNN).

### C. Deriving FPTC for LR

LR [15] performs multiple classifications by replacing the posterior probabilities of sigmoid transformation with that of softmax transformation [21]. In our derivation, the L2 norm is added into the loss function as a regularization term, which improves the numerical stability and robustness of the LR classifier. According to our derivation, the loss function of LR with the L2 norm and softmax takes the following form [21]:

$$J(\boldsymbol{\theta}) = -\frac{1}{n} \cdot \sum_{i=1}^n \sum_{j=1}^m 1\{y^{(i)} = j\} \cdot \log \frac{\exp(\boldsymbol{\theta}_j^T \mathbf{x}^{(i)})}{\sum_{l=1}^m \exp(\boldsymbol{\theta}_l^T \mathbf{x}^{(i)})} + \frac{\alpha}{2} \cdot \|\boldsymbol{\theta}\|_2^2 \quad (5)$$

where  $\boldsymbol{\theta}$  is the  $v \times m$  matrix, in which elements are the parameters of LR,  $\theta_{pj}$  is the weight of the  $j$  category in the feature layer  $p$ , and  $\alpha$  controls the regularization strength.

The goal of the LR classifier is to find the parameters of  $\boldsymbol{\theta}$  with optimal values when the minimum value of loss function is obtained. The stochastic average gradient (SAG) is a common strategy to optimize the LR classifier [22]. The SAG is an improvement on stochastic gradient descent (SGD), and only proposed general equation in the original paper. According to our derivation,  $\boldsymbol{\theta}_j^T = \{\theta_{1j}, \theta_{2j}, \theta_{3j}, \dots, \theta_{vj}\}$  of LR with L2 norm and softmax transformation are updated according to (6)-(8) [22].

$$\boldsymbol{\theta}_j^{r+1} = \boldsymbol{\theta}_j^r - \frac{\lambda}{n} \sum_{i=1}^n z_i^r \quad (6)$$

$$z_i^r = \begin{cases} \nabla_{\boldsymbol{\theta}_j} J(\boldsymbol{\theta}) & \text{if } i = i_r \\ z_i^r & \text{otherwise} \end{cases} \quad (7)$$

$$\nabla_{\boldsymbol{\theta}_j} J(\boldsymbol{\theta}) = \mathbf{x}^{(i)} (1\{y^{(i)} = j\}) - \frac{\exp(\boldsymbol{\theta}_j^T \mathbf{x}^{(i)})}{\sum_{l=1}^m \exp(\boldsymbol{\theta}_l^T \mathbf{x}^{(i)})} + \alpha \boldsymbol{\theta}_j \quad (8)$$

where  $\lambda$  is the learning rate, and  $i_r$  is taken at random from the set  $\{1, 2, 3, \dots, n\}$  for the  $r$ th iteration.

For each iteration, (6)-(8) are updated once. Therefore, the loss function is calculated one at a time by (8), and  $\boldsymbol{\theta}_j$  with  $v$  elements is updated one at a time by (7). Based on the above analysis, the FPTC of each iteration is  $F(mvn)$ . In multiple LR classification, supposing that  $Q$  iterations (the number of iterations during the SAG process) are carried out for each category, the total FPTC of LR becomes  $F(Qmvn)$ . In this case, the FPTC of multiple LR in the  $m$  category can also be written as  $F(Qm^2vn)$ . Finally, the FPTC of LR is associated with a real running time  $t'_{LR}$  as follows, and the key steps for derivating FPTC for LR are shown in Figure 2:

$$t'_{LR} = \omega_{LR} \times t^*_{LR} = \omega_{LR} \times F(Qm^2vn) \quad (9)$$

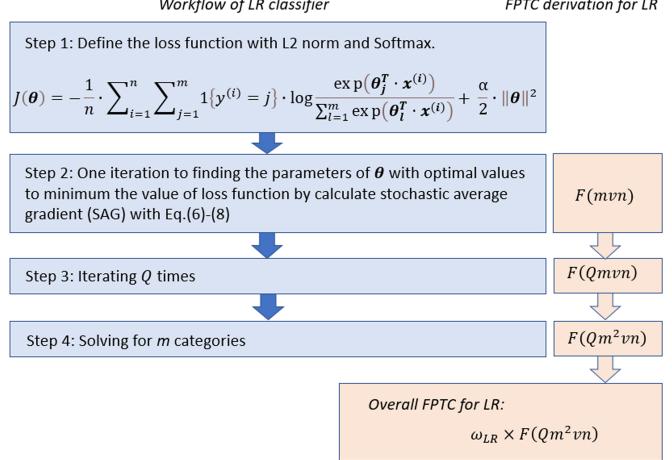


Fig. 2. Deriving full parameter time complexit (FPTC) for logistic regression (LR).

### D. Deriving FPTC for CART

CART focuses on compiling training data by the Gini splitting rule [16], [17]. The Gini index is similar to the entropy or information-gain criterion for constructing the  $IF(\mathbf{T}_r)$  impurity function [20]. The  $IF(\mathbf{T}_r)$  impurity is used to determine the split nodes of the binary tree and takes the following form [16]:

$$IF(\mathbf{T}_r) = 1 - \sum_{j=1}^m p(j/T_r)^2 = 1 - \sum_{j=1}^m \left( \frac{\sum_i^n 1\{y^{(i)} = j\}}{n} \right)^2 \quad (10)$$

where  $p(j/T_r)$  is the ratio of the number of samples in class  $j$  to the total number of samples in training data  $\mathbf{T}_r$ .

Suppose  $c$  is the split node for feature/band  $q$ . The binary tree is split into two subtrees  $\mathbf{T}_l$  and  $\mathbf{T}_g$ , in which  $\mathbf{T}_l = \{(\mathbf{x}, y) | x_q < c\}$  and  $\mathbf{T}_g = \{(\mathbf{x}, y) | x_q > c\}$ . The goal of

the CART classifier is to find the optimal  $c^*$  and  $q^*$ , which minimize the sum of  $IF(T_l)$  and  $IF(T_g)$ . The sum of the  $IF(T_l)$  and  $IF(T_g)$  is the optimization function. To make it easier to understand, we have made a simple change to the optimization function [16]:

$$J(q, c) = -\frac{n_l}{n} \sum_{j=1}^m p(j/T_l)^2 - \frac{n_g}{n} \sum_{j=1}^m p(j/T_g)^2 \quad (11)$$

$$c^*, q^* = \arg \min_{c, q} J(c, q) \quad (12)$$

where  $n_l$  is the number of samples in  $T_l$ , and  $n_g$  is the number of samples in  $T_g$ .

In the CART classifier, FPTC is divided into two parts (Figure 3). First,  $T_r$  is sorted  $v$  times by each feature in advance, which only needs to be done once. Ranking  $n$  samples is a classic computer problem, and the FPTC of its optimal algorithm is  $F(n \log_2 n)$ . The FPTC of  $v$  times ranking is then  $F(v \log_2 n)$ . Second, the binary tree is split according to (11) and (12). The FPTC of this process is  $F(mnv)$ . A binary decision tree is established when the training data are split repeatedly. According to previous studies, the average expected tree depth is known as  $\log_2 n$  [9]. Building every level in this binary tree costs  $F(mnv)$ . Thus, the FPTC of this part is  $F(mnv \log_2 n)$ . By adding the FPTC of the two parts, the FPTC of CART is  $F((m+1)nv \log_2 n)$ . In order to make the FPTC less redundant and more concise, when  $m$  is large enough, it can be simplified to  $F(mnv \log_2 n)$  with  $\frac{1}{m+1}$  error rate. The FPTC of CART is associated with the real running time  $t'_{CART}$  as follows, and the key steps for deriving FPTC for RF are shown in Figure 3:

$$t'_{CART} = \omega_{CART} \times t^*_{CART} = \omega_{CART} \times F((m+1)nv \log_2 n) \quad (13)$$

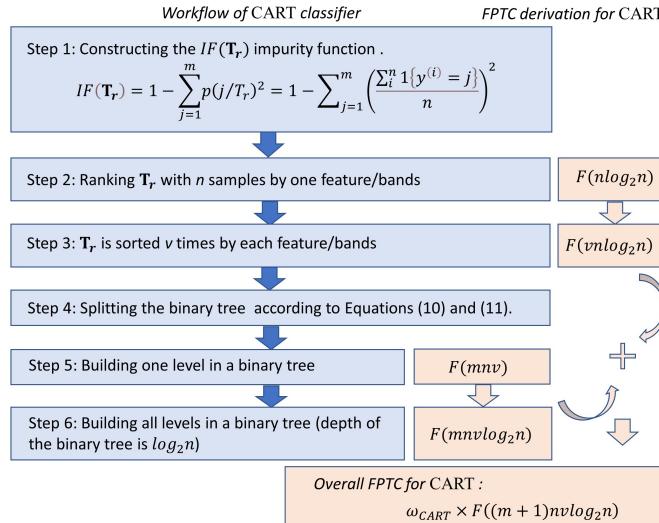


Fig. 3. Deriving full parameter time complexity (FPTC) for classification and regression tree (CART).

### E. Deriving FPTC for RF

RF is a type of ensemble learning [17]. The goal of the ensemble method is to combine several weak learners in order to improve the performance of the classifier. The ensemble methods used in RF include bagging, boosting, bootstrap, etc. The bagging, as the best ensemble method with a strong and complex model, is used in our study. Without special requirement, the bagging method randomly extracts the  $n$  sample size subset with a replacement from the training set. Generally, CART is often used as a weak learner in RF. When we build an RF with an  $s$  tree,  $s$  subsets are extracted and used to build an  $s$  CART.

When we build a CART tree with  $n$  samples, the average expected tree depth is  $\log_2 n$ . Every level cost is  $F(mnv)$ , and the FPTC of this part is  $F(mnv \log_2 n)$  [9]. When adding the sorting time  $F(vn \log_2 n)$ , the FPTC of building one tree is  $F((m+1)nv \log_2 n)$ , so the FPTC of RF with  $s$  trees is  $F(s(m+1)nv \log_2 n)$ . Similarly, the FPTC of RF can be simplified to  $F(smnv \log_2 n)$  with a  $\frac{1}{m+1}$  error rate when  $m$  is large enough. The FPTC of RF is associated with the real running time  $t'_{RF}$  as follows, and the key steps for deriving FPTC for RF are shown in Figure 4:

$$t'_{RF} = \omega_{RF} \times t^*_{RF} = \omega_{RF} \times F(s(m+1)nv \log_2 n) \quad (14)$$

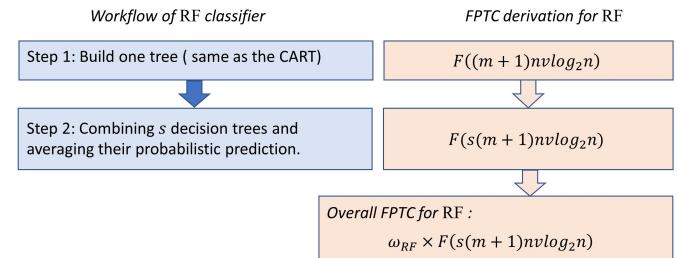


Fig. 4. Deriving full parameter time complexity (FPTC) for random forest (RF).

### F. Deriving FPTC for SVM

SVM [18] aims to find the hyperplane with the maximum distance from the support vectors, which is often treated as a linear programming problem to transform it into a dual problem by the Lagrangian multiplier method [23], [24]. The function of hyperplane  $g(\mathbf{x}, \boldsymbol{\theta})$  takes the following form [18]:

$$g(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\theta}^T \mathbf{x} + b = \sum_{j=1}^v \theta_j x_j + b. \quad (15)$$

In general, a soft margin hyperplane [24] is often used to avoid overfitting, which is done by introducing the variable  $\xi$ . Based on the Lagrangian multiplier method, the loss function with a soft margin takes the following form [18], [24]:

$$\begin{aligned}
 L(\boldsymbol{\theta}, b, \boldsymbol{a}, \boldsymbol{\xi}, \boldsymbol{\mu}) = & \frac{1}{2} \|\boldsymbol{\theta}\|^2 + C \sum_{i=1}^n \xi_i \\
 & + \sum_{i=1}^n \alpha_i (1 - \xi_i - y^{(i)} (\boldsymbol{\theta}^T \boldsymbol{x}^{(i)} + b)) \\
 & - \sum_{i=1}^n \mu_i \xi_i
 \end{aligned} \tag{16}$$

where  $C$  is the regularization parameter and  $\alpha_i \geq 0, \mu_i \geq 0$ . The dual question of finding the hyperplane with the maximum distance in SVM with soft margin is as follows [18], [24]:

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} K_{ij} \tag{17}$$

$$\text{s.t. } \sum_{i=1}^n \alpha_i y^{(i)} = 0, \alpha_i \geq 0, i = 1, 2, \dots, n \tag{18}$$

where  $K_{i,j} = K(\boldsymbol{x}_i, \boldsymbol{x}_j) = \exp(-\gamma \times \|\boldsymbol{x}_i - \boldsymbol{x}_j\|^2)$  ( $\gamma > 0$ ) is the radial basis function (RBF). In our study, the RBF, as the most commonly kernel function, is used in our research. The kernel choice can be a key issue for the SVM performance, since the different kernel could lead to different results. From the perspective of time complexity, they are the same. For instance, the time complexity of RBF is  $F(v)$  (where  $v$  represents the number of the features of inputs), which is identical to the time complexity of the linear kernel. From this perspective, the “kernel trick” is efficient by measuring the difference between samples in the higher dimensional space without actually projecting them into it. Based on this reason, the different kernel function selection, such as linear, polynomial, sigmoid, or RBF, will not have a significant effect on the FPTC in SVM.

Karush–Kuhn–Tucker (KKT) conditions are as follows [18], [24]:

$$\alpha_i \geq 0, \mu_i \geq 0 \tag{19}$$

$$y^{(i)} \cdot g(\boldsymbol{x}^{(i)}) - 1 + \xi_i \geq 0 \tag{20}$$

$$\alpha_i \cdot (y^{(i)} \cdot g(\boldsymbol{x}^{(i)}) - 1 + \xi_i) = 0 \tag{21}$$

$$\xi_i \geq 0, \mu_i \xi_i = 0 \tag{22}$$

It is difficult and costly to solve the dual problem directly. In our study, sequential minimal optimization (SMO) [25] is considered a decomposition method to conquer this difficulty. At each iteration, we solve a simple two-variable problem ( $\alpha_i$  and  $\alpha_j$ ) without needing any optimization software [25].

$$\alpha_j^{r+1} = \alpha_j^r + \frac{y^{(j)}((g(\boldsymbol{x}^{(i)}) - y^{(i)}) - (g(\boldsymbol{x}^{(j)}) - y^{(j)}))}{K_{i,i} + K_{j,j} - 2K_{i,j}} \tag{23}$$

$$\alpha_i^{r+1} = \alpha_i^r + y_i y_j (\alpha_j^r - \alpha_j^{r+1}) \tag{24}$$

when  $\alpha_i^{r+1}$  is not at the bounds,  $b$  is calculated as follows [25]:

$$\begin{aligned}
 b^{r+1} = & g(\boldsymbol{x}^{(i)}) - y^{(i)} + y_i (\alpha_i^{r+1} - \alpha_i^r) K_{i,i} \\
 & + y_j (\alpha_j^{r+1} - \alpha_j^r) K_{i,j} + b^r
 \end{aligned} \tag{25}$$

when  $\alpha_j^{r+1}$  is not at the bounds,  $b$  is calculated as follows [25]:

$$\begin{aligned}
 b^{r+1} = & g(\boldsymbol{x}^{(j)}) - y^{(j)} + y_i (\alpha_i^{r+1} - \alpha_i^r) K_{i,j} \\
 & + y_j (\alpha_j^{r+1} - \alpha_j^r) K_{j,j} + b^r
 \end{aligned} \tag{26}$$

when neither  $\alpha_i^{r+1}$  nor  $\alpha_j^{r+1}$  is at bounds, (25) and (26) are equal. When both  $\alpha_i^{r+1}$  and  $\alpha_j^{r+1}$  are at bounds,  $b^{r+1}$  is halfway between (25) and (26). The  $k$  is the number of support vectors and  $\boldsymbol{x}^{(i)} \in \{\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \dots, \boldsymbol{x}^{(k)}\}$  is the support vector. Notably, the SMO optimizer [26] has been replaced with the SMO-type optimizer since LIBSVM V2.8. The SMO-type optimizer may speed up convergence and reduce the number of iterations. When we consider the FPTC in one iteration, the difference in times of the FPTC between SMO and the SMO-type optimizer can be neglected.

In each iteration, all training samples are calculated and judged whether they are in KKT conditions (19)–(22). We selected the one that does not comply with KKT ( $(x^{(i)}, y^{(j)})$  (21) and the one ( $x^{(i)}, y^{(j)}$ ) with the greatest distance from it. We then solved  $\alpha_i, \alpha_j, \theta$ , and  $b$ . The FPTC of each iteration is  $F(nv + vk)$ . KKT conditions are updated with every iteration, and the loop is stopped when all training samples fit the KKT conditions or we reach the maximum number of iterations we set. Thus, the FPTC of iteration  $Q$  (the number of iterations during the SMO process) is  $F(Qv(n + k))$ .

Originally, SVM was designed for binary classification, and it is not good at multiple classifications. There are still several multi-classification methods for SVM; one-versus-one (OVO) or pairwise is one of them. In OVO, an SVM classifier needs to be built between each of the two classes, that is,  $m(m-1)/2$  classifiers need to be built for  $m$  categories. Thus, the FPTC of SVM for multi-classification is  $F(m^2Qv(n + k))$ . It is associated with the real running time  $t'_{SVM}$  as follows, and the key steps for deriving FPTC for SVM are shown in Figure 5:

$$t'_{SVM} = \omega_{SVM} \times t^*_{SVM} = \omega_{SVM} \times F(m^2Qv(n + k)) \tag{27}$$

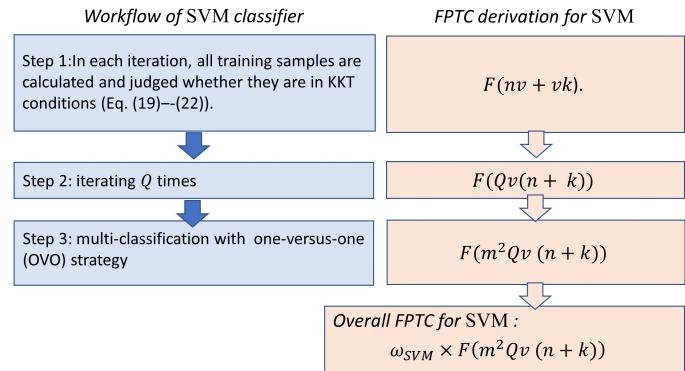


Fig. 5. Deriving full parameter time complexity (FPTC) for support vector machine (SVM).

### III. STUDY AREA AND DATASETS

To test the accuracy of the FPTC derived above, one study area and the remote sensing data were selected.

#### A. Study area

The study area is located in the middle of the Xinjiang Uygur Autonomous Region, China ( $42^{\circ}7' - 42^{\circ}13'$  N,  $86^{\circ}13' - 86^{\circ}22'$  E) (Figure 6). This area covers an extensive area of around 1,660,000 km<sup>2</sup>, mainly covered by grassland and sandy desert with a typical continental climate. Forest areas are sparsely scattered within the high mountains and along the rivers. Oasis landscapes have developed within inland river deltas, alluvial–diluvial plains, and along the edges of diluvial–alluvial fans. Agricultural land and human settlements are distributed around these areas.

With the rapid growth of the population in recent years, the ecosystem in this area faces a great challenge in terms of the dramatic change of land cover combined with changing precipitation patterns [27]. The frequency of the occurrence of natural disasters, such as sandstorms, floods, and snowstorms, increases as the global climate changes. Under these circumstances, a rapid assessment of land cover change after a disaster not only limits the loss of life and property but also provides data for emergency managers to optimize the emergency response procedure.

#### B. Datasets

Sentinel-2 provides continuous high-resolution images with a multispectral instrument (MSI). With the high revisit frequency (five days combined with Sentinel 2A and 2B), Sentinel-2 imagery has been widely used for land cover mapping, change detection, and emergency response [28]. Compared with other public and free multispectral products, Sentinel-2 contains bands covering the red edge [29], which can provide indispensable information for land-cover mapping, land change detection, and the retrieval of other geophysical variables [28], [29]. In our study, a Sentinel-2A image acquired on 8 September 2016 was downloaded from the United States Geological Survey (<https://earthexplorer.usgs.gov/>). This image corresponds to Level-1C products, which are radiometrically and geometrically corrected Top-of-Atmosphere (TOA) products with subpixel multispectral registration [29]. This image is cloud-free, so the atmospheric correction procedure has little influence on classification in this study [30], [31].

We selected an area ( $2048 \times 2048$  pixels) from the original Sentinel-2A image (Figure 6(c)). The spatial resolutions of B5, B6, B7, B8a, B11, and B12 were resampled to 10 m by nearest-neighbor resampling [29].

Based on the field investigation, the land cover classification system in the study area was established. Eight typical land cover types, farmland, orchard, forest, grassland, water, residential area, roadway, and idle land (Figure 7) were selected as land cover types at the level-1 category, which remains the same as the study from other groups in this area, such as Gong and Howarth [32] and Gong et al. [33]. The image was manually interpreted to create a digital land cover map through

image interpretation with intensive field samples (done in October 2016) over this area (Figure 7b). Both imageries with 12 bands and the land cover map provide sufficient training samples for land cover classification in this area. To understand how time consumption changes under the different bands or the different training sample size, the training sample is divided into different groups, with samples randomly selected over the land cover map. Considering the slow training speed of SVM, only nine training groups for 1000 to 10,000 samples were prepared for SVM. Table I shows the details of the sub-training sample settings.

TABLE I  
ALGORITHM PARAMETER SET UP AND SOURCE OF CODES.

Classifier	Number of groups	Sample size
kNN	29	1000, 2000, 3000, 4000, 5000,
		6000, 7000, 8000, 9000, 10,000,
		20,000, 30,000, 40,000, 50,000,
		60,000, 70,000, 80,000, 90,000,
LR		100,000, 110,000, 120,000, 130,000,
		140,000, 150,000, 160,000, 170,000
CART		180,000, 190,000, 200,000
		1000, 2000, 3000, 4000, 5000,
RF		6000, 7000, 8000, 9000, 10,000
SVM	10	1000, 2000, 3000, 4000, 5000,
		6000, 7000, 8000, 9000, 10,000

All classifiers are programmed by Python 3.6 (Python Software Foundation. Python Language Reference, version 3.6. Available at <https://www.python.org/>). The SVM classifier is programmed using LIBSVM [34], while the rest of the classifiers are programmed using Scikit-learn [35]. All experiments ran on the Ubuntu 14.0 (<https://ubuntu.com/>) platform, which was equipped with an Intel Xeon e5-2620 CPU and four TITAN XP graphics cards.

#### C. Assessment

To verify the accuracy of the FPTC, three assessments were applied: 1) a 1:1 plot was used to compare the real running time to the FPTC, 2) we estimated the running time by FPTC and calculated the root-mean-squared error (RMSE) between the estimated and the observed running time, and 3) we compared the FPTC and the TTC with respect to the real observed running time under different feature selections.

## IV. RESULTS

#### A. The validation of the FPTC and correction coefficient

We selected sub-training samples randomly from the training dataset and constructed sets of sub-training samples (Table I). These samples were classified, and the real running time was recorded. As we can see from Figure 8, the linear relationship between the FPTC and the real running time indicates the effectiveness of FPTC. The R-squared values for the 1:1 plot of the five classifiers were all larger than 0.99 (kNN: 0.991, LR: 0.997, CART: 0.999, RF: 1.000, and SVM: 0.999), which indicates that the linear relationship between the algorithm part of the FPTC, and the real running time is extremely strong ( $p < 0.001$ ).

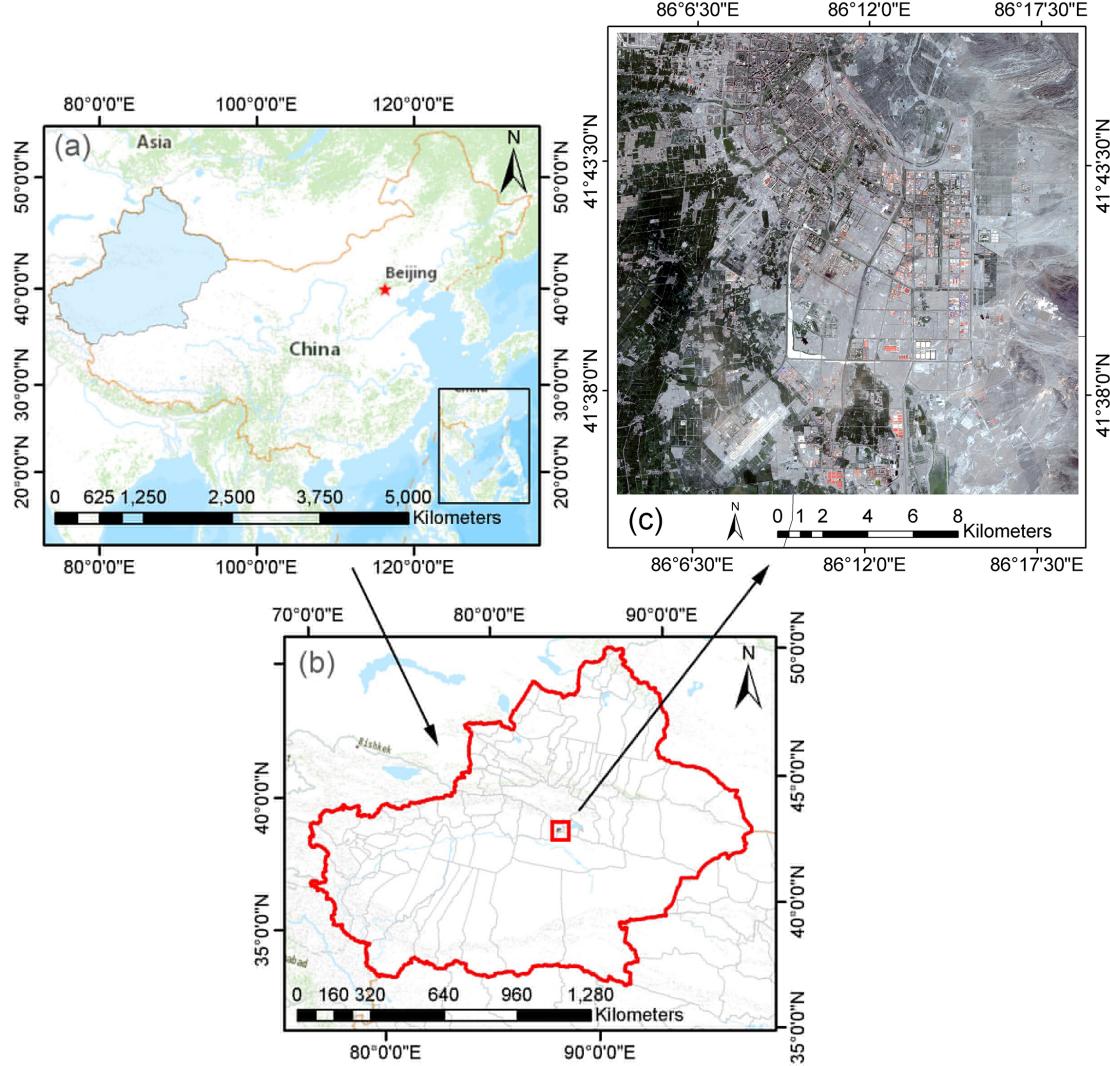


Fig. 6. The study area: (a) China; (b) Xinjiang Uygur Autonomous Region; (c) location of the study area.

The coefficient  $\omega$  of each algorithm can be found from the slope of the regression line. For  $k$ NN, the correction coefficient  $\omega_{kNN} = 6.90E - 8$ , which is obtained from the linear relationship ( $t'_{kNN} = 6.90E - 8t_{kNN}^*$ ,  $R^2 = 0.991$ ,  $p \leq 0.001$ ). This relationship shows that the actual running time (in seconds) is equal to FPTC  $kNN \times \omega_{kNN}$  (Figure 8a). Similarly, for LR, the coefficient  $\omega_{LR} = 1.70E - 9$  has a linear relationship ( $t'_{LR} = 1.70E - 09t_{LR}^*$ ,  $R^2 = 0.997$ ,  $p \leq 0.001$ ) (Figure 8b). For CART, the correction coefficient  $\omega_{CART} = 1.06E - 08$ , which has a linear relationship ( $t'_{CART} = 1.06E - 08t_{CART}^*$ ,  $R^2 = 0.999$ ,  $p \leq 0.001$ ) (Figure 8c). For RF, the correction coefficient  $\omega_{RF} = 2.35E - 9$ , which has a linear relationship ( $t'_{RF} = 2.35E - 9t_{RF}^*$ ,  $R^2 = 1.000$ ,  $p \leq 0.001$ ) (Figure 8d). For SVM, the correction coefficient  $\omega_{SVM} = 2.66E - 9$ , which has a linear relationship ( $t'_{SVM} = 2.66E - 9t_{SVM}^*$ ,  $R^2 = 0.999$ ,  $p \leq 0.001$ ) (Figure 8e).

The slope can be roughly estimated based on the two available datasets, regardless of the magnitude of the training data. This means that this value can be obtained by prerunning

the algorithm under two small parts of the total dataset. As a result that the coefficient  $\omega$  represents the physical part of the FPTC, this value only varies when the algorithm is applied to different computational environments.

#### B. Estimating the real running time with the FPTC

Based on the strong linear relationship provided by Figure 8, the running time can be estimated accurately. To calculate the RMSE between the estimated and observed running time, we fixed samples at 100,000 and recorded the real running time. Table II shows the estimated and actual running times for each algorithm.

Regarding the FPTC of each classifier, the highest FPTC comes from SVM (3.05E+11), followed by RF (1.20E+10), LR (1.09E+09), and CART (1.49E+08); the lowest FPTC, produced by  $k$ NN, was 1.35E+06 (shown in Table II).

Table II shows that the real running time keeps a certain consistency with the estimated running time. From Table II, we can see that the highest real running time was achieved by SVM (803.800 s), while the estimated running time of

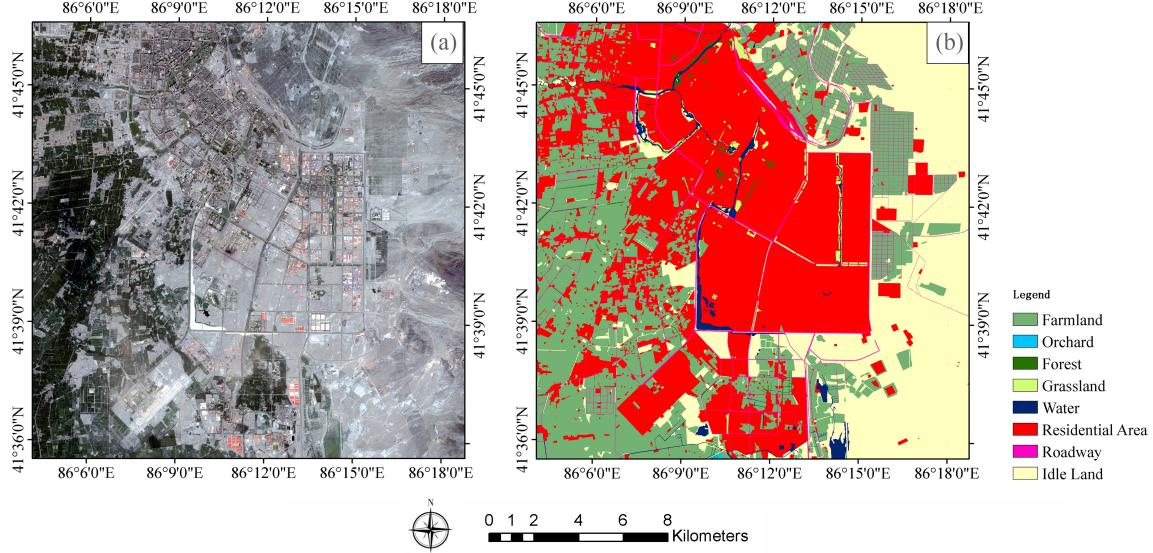


Fig. 7. Location of the selected image: (a) image; (b) ground truth.

TABLE II  
ESTIMATED AND REAL RUNNING TIMES FOR EACH ALGORITHM.

Classifier	Parameter choice	FPTC	$\omega$	Average estimated running time (s)	Real running time (s)	RMSE (s)
<i>k</i> NN	$v = 10; n = 100,000; k = 11$	$1.35E + 06$	$6.90E - 08$	0.093	0.081	0.01
LR	$v = 10; n = 100,000; C = 1E - 05; Q = 17$	$1.09E + 09$	$1.70E - 09$	1.850	1.741	0.18
CART	$v = 10; n = 100,000;$	$1.49E + 08$	$1.06E - 08$	1.585	1.632	0.16
RF	$v = 10; n = 100,000; s = 80$	$1.20E + 10$	$2.35E - 09$	28.104	28.445	0.57
SVM	$v = 10; n = 100,000; C = 100; \gamma = 8; Q = 44355; k = 7290$	$3.05E + 11$	$2.66E - 09$	810.055	803.800	15.76

SVM was 810.055 s. The next highest real running time was achieved by RF (28.445 s), while the estimated running time of RF was 28.104 s. LR (1.741 s), CART (1.632 s), and *k*NN (0.081 s) produced the lowest real running times, while the estimated running times of LR, CART, and *k*NN were 1.850 s, 1.585 s, and 0.093 s, respectively (Table II). Regarding the RMSE of each classifier, the highest RMSE came from SVM (15.76 s), followed by RF (0.57 s), LR (0.18 s), CART (0.16 s), and *k*NN (0.01 s) (shown in Table II). The average RMSE between the real running time and the estimated running time was 3.34 s, which shows that the real running time can be estimated accurately by FPTC.

### C. Comparing FPTC and TTC

How does the performance of FPTC compare to TTC? In this section, we compare FPTC to TTC using both theoretical and experimental methods.

Firstly, at a theoretical level, TTC has a limited ability to show the difference in running time between the different algorithms. For instance, the TTC of *k*NN and LR is  $O(n)$ , while the TTC of CART and RF is  $O(n \log_2 n)$  (Table III). If we use TTC for the evaluation of running time, there is no difference between the running time of *k*NN and LR or between CART and RF. From our experiments with 100,000

samples, the real running time of *k*NN was 0.081 s, and the real running time of LR was 1.741 s. Similarly, a noticeable difference exists between CART and RF. The real running time of CART was 1.632 s, while the real running time of RF was 28.445 s (Table III).

The TTC of LR is  $O(n)$ , while the TTC of CART is  $O(n \log_2 n)$ . In terms of TTC, the running time of CART should be higher than that of LR. The reality is just the opposite. The real running time of CART was 1.632 s, which is lower than that of LR (1.741 s). The reason is that, in TTC, many low-order details and key parameters have been ignored. These details and parameters may be negligible when comparing algorithm time complexity at a theoretical level, but they are essential to estimating running time at the operational level.

The FPTC we proposed can make up for the above shortcomings of TTC. FPTC is derived by examining the overall structure of the classifier program and its mathematical principles. As we can see from Table III, the FPTC of the *k*NN is different from the FPTC of the LR.

In addition, FPTC can also reflect changes in running time with different parameters. When  $n$  training samples changed from low to high and other influencing parameters were kept the same, the FPTC of SVM was most vulnerable to the

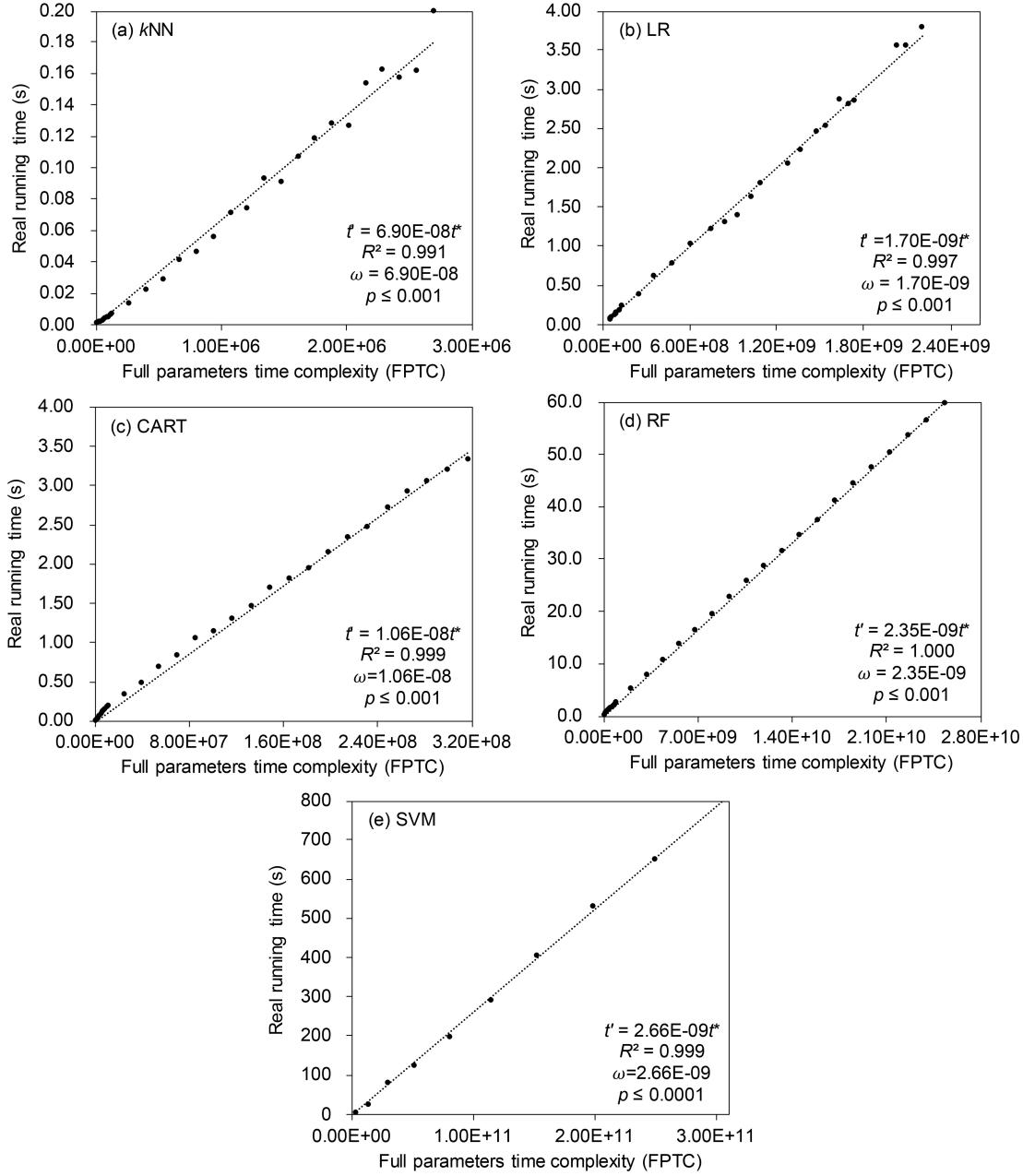


Fig. 8. The relationship between full parameter time complexity (FPTC) and the real running time.

effects of the change, followed by RF, CART, *k*NN, and LR. If  $n$  changed from 1 to 128, the FPTC of LR increased 128-fold, *k*NN a little more than 128-fold, CART and RF 896-fold, and SVM more than 16,384-fold. When the number of categorizations  $m$  changed from low to high and other influencing parameters were kept the same, the FPTC of SVM and LR was most vulnerable to the effects of the change, followed by CART and RF; *k*NN was unaffected. For example, if  $m$  changed from two classes to 200 classes, the FPTC of SVM and LR increased 10,000-fold, and that of CART and RF increased 100-fold. When the number of features or bands  $v$  changed from low to high and other influencing parameters were kept the same, the time complexities of SVM, LR, CART, and RF were changed within a polynomial time, while that of

*k*NN was less affected.

Secondly, to further illustrate the difference between FPTC and TTC from experiments, we analyzed the changing trends of FPTC and TTC under all combinations of different bands ( $v = 3, 4, 5, \dots, 10$ ) and different sample sizes ( $n = 10, 20, 30, \dots, 100,000$ ) and compared them with real running time trends. The values of TTC, FPTC, and the real running times are mapped from low to high in red to green in Figure 9. The results show that TTC does not respond to changes in  $v$ , while FPTC can better reflect the variation of  $v$ .

As we can see, FPTC shows a similar pattern to the real running time under different bands and data sizes. The pattern of the TTC is different since the impact of the different features/bands is ignored by TTC.

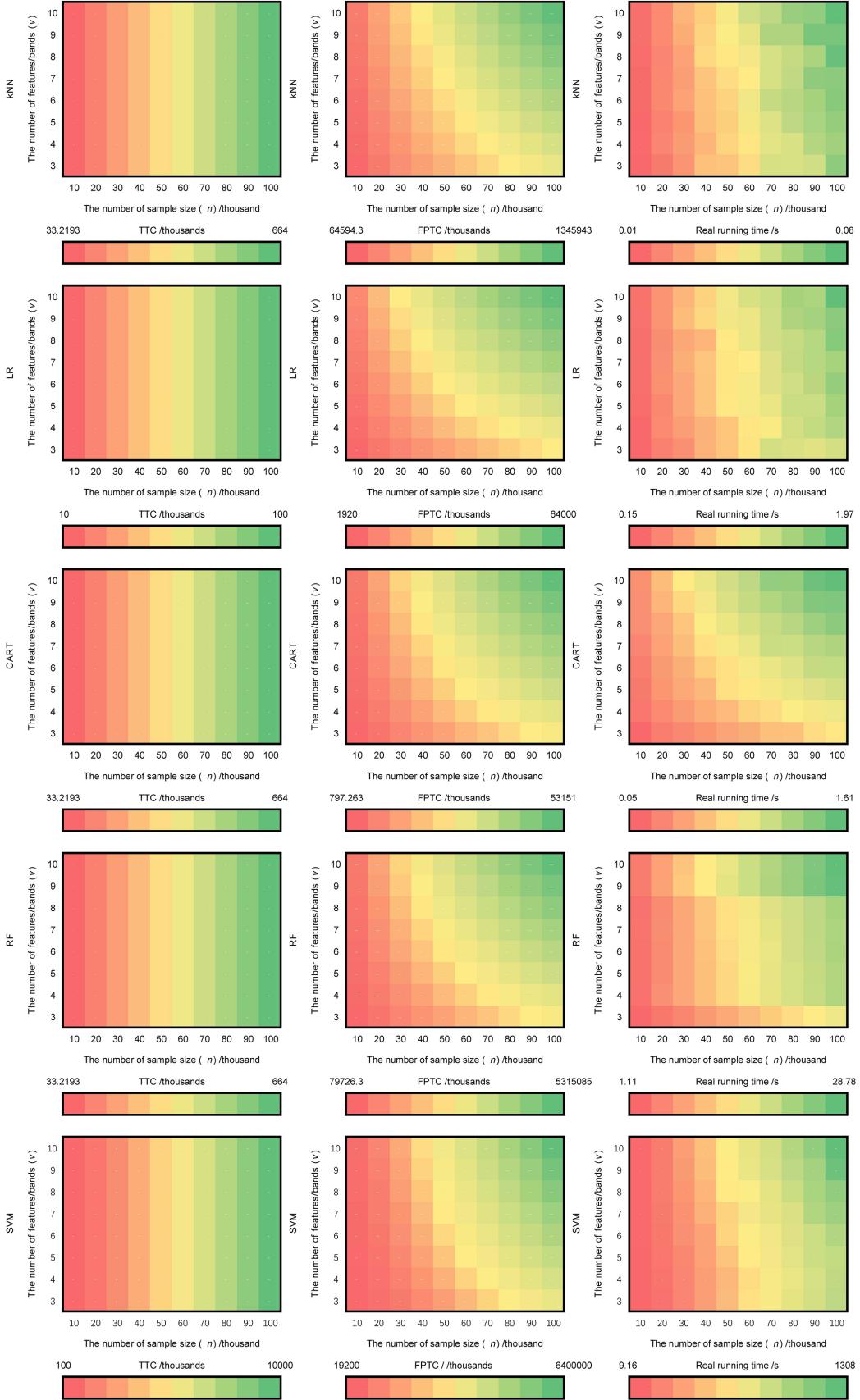


Fig. 9. The comparison between TTC (left column) and FPTC (middle column) to the real running time (right column).

TABLE III  
COMPARISON OF TRADITIONAL TIME COMPLEXITY (TTC) AND FULL PARAMETER TIME COMPLEXITY (FPTC).

Classifier	$\omega$	FPTC	TTC	Real running time (s) (100,000 samples)
kNN	$6.90E - 08$	$F(nv + n\log_2 u)$	$O(n)$	0.081
LR	$1.70E - 09$	$F(m^2vn)$	$O(n)$	1.741
CART	$1.06E - 08$	$F((m + 1)nv\log_2 n)$	$O(n\log_2 n)$	1.632
RF	$2.35E - 09$	$F(s(m + 1)nv\log_2 n)$	$O(n\log_2 n)$	28.445
SVM	$2.66E - 09$	$F(m^2vn^2)$	$O(n^2)$	803.800

#### D. A simple application of FPTC

Without estimating the running time accurately, choosing a classifier for a time-limited LULC classification task could be blind and subjective. In this section, we create an urgent classification task that should be done within 6 h. The natural question is of how many training samples we should prepare for different classifiers to fit this time limit.

Through FPTC and  $\omega$ , we can estimate the maximum sample size (MSS) that can be processed within the time threshold. A larger maximum training sample means that the algorithm can handle more samples, which will improve the overall classification accuracy. Figure 10 shows the relationship between sample size and the running time calculated by FPTC with corresponding coefficient  $\omega$ . The exact number of training samples for each hour is also shown in these figures.

Taking a 1-h training limit as an example, the results show that the algorithm with the smallest MSS is SVM (0.21 million) (Figure 10e), followed by RF (10 million) (Figure 10d), CART (140 million) (Figure 10c), and LR (200 million) (Figure 10b). The algorithm with the maximum MSS is kNN (3.84 billion) (Figure 10a).

In our study, the main goal was to provide a quantitative measurement for emergency managers to compare and filter different classifiers under different time and resource limits. Threshold analysis can help us to perform classifier screening. For instance, suppose we need to process a Sentinel-2A image ( $5000 \times 5000$  pixels). It would require 25 million pixels to construct the model. With SVM, it is impossible to execute the current task. If we only consider the running time, kNN could be the optimal classifier.

## V. DISCUSSION

### A. The effect of training parameters on FPTC.

The parameters that affect the running time of an algorithm can be classified into two major categories. One is a hyperparameter, which is set before starting the learning process (e.g., the number of trees  $s$  in RF). The other is a training parameter, which can only be obtained after running (e.g., the total iterations  $Q$  and the number of support vectors  $k$  in SVM). These training parameters can only be obtained after the program runs. Thus, it will be difficult to estimate the running time of these classifiers without running them. Fortunately, the characteristics of the training parameter can be estimated by other pre-obtained parameters such as the sample size. Figure 11 shows that the total iterations  $Q$  in SVM has a linear correlation with the number of samples ( $R^2 = 1.00$ ) (Figure 11a). The hyperparameter  $k$  also has the

same significant linear correlations ( $R^2 = 1.00$ ) (Figure 11b). According to these linear correlations, parameters  $Q$  and  $k$  can be removed from  $F(m^2Qv(n+k))$  in SVM. The real running time  $t'_{SVM}$  in (27) can then be rewritten as (28), where  $\omega_{SVM}$  is the new coefficient. The total FPTC of SVM is  $F(m^2vn^2)$ . After this, the estimated running time is no longer influenced by the training parameters.

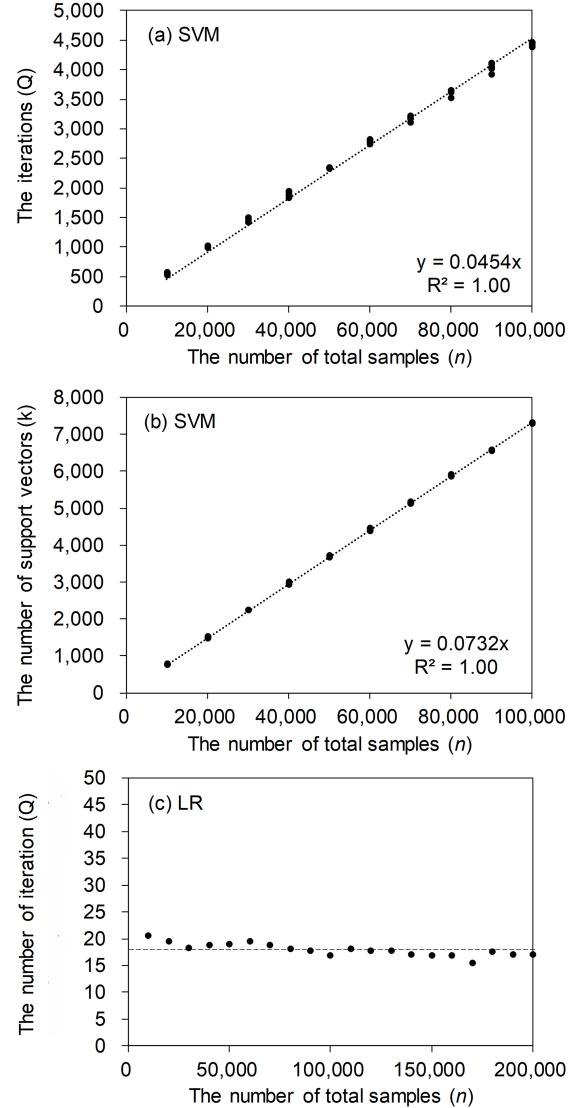


Fig. 11. The correlations between the training parameters and the number of samples.

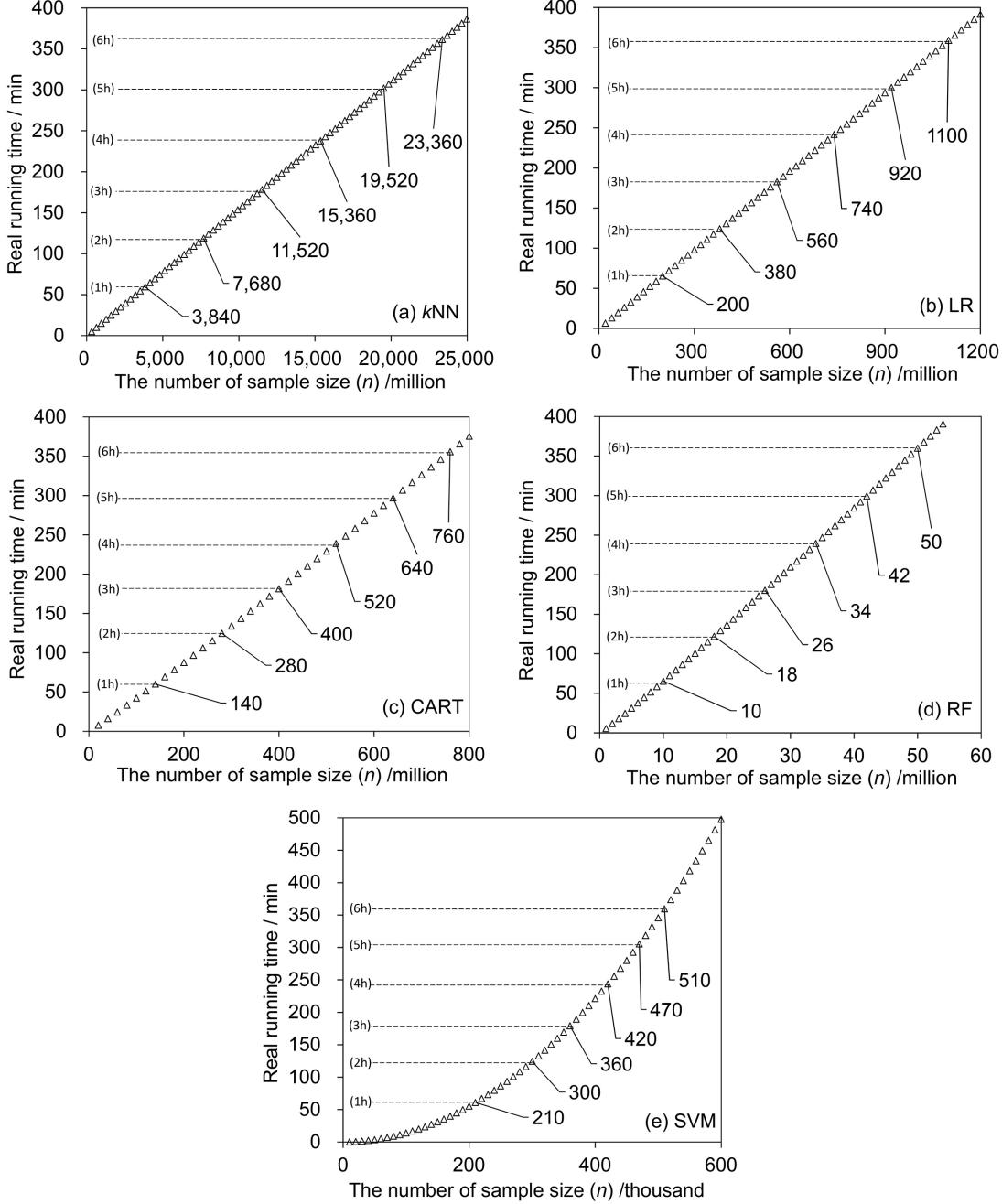


Fig. 10. The capacity of each algorithm.

$$\begin{aligned}
 t'_{SVM} &= \omega_{SVM} \times t^*_{SVM} = \omega_{SVM} \times F(m^2 Q v(n+k)) \\
 &= \omega'_{SVM} \times F(m^2 v n^2)
 \end{aligned} \tag{28}$$

Another classifier with training parameters is LR. The total iterations  $Q$  in LR is generated in the SAG algorithm, which can only be obtained after iteration. Figure 11c shows that the total iterations  $Q$  is stable with an average of  $18.02 \pm 1.16$ ; it changes as the sample size  $n$  increases (Figure 11c). Therefore, in the derivation of FPTC,  $Q$  can be approximated as a constant and combined with coefficient  $\omega$ . After this correction, the total FPTC of LR is corrected to  $F(m^2 v n)$  (29). Thus, the

estimated running time is no longer influenced by the training parameters.

$$\begin{aligned}
 t'_{LR} &= \omega_{LR} \times t^*_{LR} = \omega_{LR} \times F(Q m^2 v n) \\
 &= \omega'_{LR} \times F(m^2 v n)
 \end{aligned} \tag{29}$$

#### B. Determining coefficient $\omega$ with a pre-experiment

As mentioned before, the coefficient  $\omega$  is the key to estimating the actual running time of classifiers. This parameter is more related to the physical computational environment, such as the compiling programs and computer hardware. Figure 8 shows that the coefficient  $\omega$  (the slope) will not change as the

input data size increases, even when the data size is small. Based on this finding, the coefficient  $\omega$  can be calculated through two small-scale pre-experiments. The coefficient  $\omega$  can then be reused in other estimates with a larger sample size. Once the physical computational environment changes, the coefficient  $\omega$  should be reevaluated under the new circumstances. This may limit the application of FPTC to a cloud computing environment since the physical computational environment may change during the realization process. How the coefficient  $\omega$  is related to the computational resource (CPU or GPU frequency) still needs to be quantitatively evaluated in the future.

## VI. CONCLUSION

In emergency response to natural disasters, accurate time predictions help emergency managers to choose a classification algorithm with limited time and resources. In this study, we proposed FPTC and the coefficient  $\omega$  to estimate the running time of each classifier. The FPTC of five common classifiers (*k*NN, LR, CART, RF, and SVM) was derived by examining the overall structure of the classifier program and its mathematical principles. A linear regression model was built based on the relationship between the real running time and the FPTC, and the coefficient  $\omega$  was obtained from the linear regression. We then accurately predicted the running time of each classifier and filtered out the appropriate classifier. The results can be summarized as follows:

- 1) We proposed a method to quantitatively evaluate the time efficiency of machine learning classifiers called FPTC. We derived the FPTC of five general classifiers. The results show that the FPTC of *k*NN is  $F(nv + n\log_2 u)$ , the FPTC of LR is  $F(Qm^2vn)$ , the FPTC of CART is  $F((m + 1)nv\log_2 n)$ , the FPTC of RF is  $F(s(m + 1)nv\log_2 n)$ , and the FPTC of SVM is  $F(m^2Qv(n + k))$ .
- 2) A strong linear relationship between the FPTC and the running time was found in our study ( $R^2 \geq 0.991$ ,  $p \leq 0.001$ ). This linear relationship verifies the correctness of the FPTC derivation process. The correction coefficient  $\omega$  of each algorithm can be found from a strong linear regression.
- 3) The running time of each classifier was estimated by coefficient  $\omega$  with FPTC. Our study showed that the average RMSE between the real running time and the estimated running time is 3.34 s, which shows the feasibility and accuracy of using FPTC to predict the running time of algorithms.
- 4) Our study showed that training parameter  $Q$  in SVM had significant linear correlations with the number of samples ( $R^2 = 1.00$ ), and  $Q$  in LR was stable and did not alter with  $n$ . According to the above rules, the total FPTC of SVM was corrected to  $F(m^2vn^2)$  and the total FPTC of LR was corrected to  $F(m^2vn)$ . The updated FPTC was not affected by whether the program was run in advance.

In a future study, we plan to derive more FPTC values for algorithms. A suitable algorithm with good accuracy and low

FPTC can be quickly filtered for an emergency task, which helps emergency managers make quick decisions in response to natural disasters based on the amount of remote sensing data available.

TABLE IV  
NOMENCLATURE.

Notations	Descriptions
$T_r \in R^{n \times v}$	Training dataset
$T_e \in R^{z \times v}$	Test dataset
$n$	The number of total samples for $T_r$
$v$	The number of bands/features
$m$	The number of categorizations
$b$	The number of total samples for $T_e$
$t'$	The running time
$C$	The regularization parameter
$\omega$	The constant coefficient
$\theta'$	The model parameters
$u$	The number of object group in <i>k</i> NN
$\lambda$	The learning rate
$s$	The number of trees in the random forest
$Q$	The total iterations
$k$	The number of support vectors
$\gamma$	The kernel parameters in SVM

## ACKNOWLEDGMENT

This research was supported by the Strategic Priority Research Program of the Chinese Academy of Sciences (Project No. XDA19030301). The National Key Research and Development Program of China (Project No. 2017YFB0504203). Natural Science Foundation of China project (Project No. 41601212, 41801360, 41771403, 41801358). The Fundamental Research Foundation of Shenzhen Technology and Innovation Council (Project No. KCXFZ202002011006298). The authors thank two reviewers for suggestions and thank Yong.Zhang and Jincheng.Jiang from SIAT for valuable suggestions.

## REFERENCES

- [1] C. N. Koyama, H. Gokon, M. Jimbo, S. Koshimura, and M. Sato, "Disaster debris estimation using high-resolution polarimetric stereo-SAR," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 120, pp. 84–98, 2016.
- [2] S. Liu and M. E. Hodgson, "Satellite image collection modeling for large area hazard emergency response," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 118, pp. 13–21, 2016.
- [3] X. Zheng, F. Wang, and Z. Li, "A multi-UAV cooperative route planning methodology for 3D fine-resolution building model reconstruction," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 146, pp. 483–494, 2018.
- [4] R. Khatami, G. Mountrakis, and S. V. Stehman, "A meta-analysis of remote sensing research on supervised pixel-based land-cover image classification processes: General guidelines for practitioners and future research," *Remote Sensing of Environment*, vol. 177, pp. 89–100, 2016.
- [5] D. Marmanis, J. D. Wegner, S. Galliani, K. Schindler, M. Datcu, and U. Stilla, "Semantic segmentation of aerial images with an ensemble of CNNs," in *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, 2016, pp. 473–480.

- [6] L. Yu, L. Liang, J. Wang, Y. Zhao, Q. Cheng, L. Hu, S. Liu, L. Yu, X. Wang, P. Zhu, X. Li, Y. Xu, C. Li, W. Fu, X. Li, W. Li, C. Liu, N. Cong, H. Zhang, F. Sun, X. Bi, Q. Xin, D. Li, D. Yan, Z. Zhu, M. F. Goodchild, and P. Gong, "Meta-discoveries from a synthesis of satellite-based land-cover mapping research," *International Journal of Remote Sensing*, vol. 35, no. 13, pp. 4573–4588, 2014.
- [7] X. Huang, J. Wang, J. Shang, C. Liao, and J. Liu, "Application of polarization signature to land cover scattering mechanism analysis and classification using multi-temporal C-band polarimetric RADARSAT-2 imagery," *Remote Sensing of Environment*, vol. 193, pp. 11–28, 2017.
- [8] D. L. Donoho, Y. Tsaig, I. Drori, and J. L. Starck, "Sparse solution of underdetermined systems of linear equations by stagewise orthogonal matching pursuit," *IEEE Transactions on Information Theory*, vol. 58, no. 2, pp. 1094–1121, 2012.
- [9] T. H. Cormen, C. E. Leiserson, R. L., and C. Stein, *Introduction to algorithms*. MIT Press, 2009.
- [10] F. A. Mianji and Y. Zhang, "Robust hyperspectral classification using relevance vector machine," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 49, no. 6, pp. 2100–2112, 2011.
- [11] A. Plaza and C.-I. Chang, "An improved N-FINDR algorithm in implementation," in *Algorithms and Technologies for Multispectral, Hyperspectral, and Ultraspectral Imagery XI*, vol. 5806, Orlando, Florida, 2005.
- [12] W. Yan and W. Wu, *Data structure*. Beijing: Tsinghua university press., 1992. (in Chinese).
- [13] D. E. Knuth, *The art of computer programming, volume vol.2: Seminumerical algorithm*. Addison-Wesley Educational Publishers Inc, 1997.
- [14] T. M. Cover and P. E. Hart, "Nearest neighbor pattern classification," *IEEE Transactions on Information Theory*, vol. 13, no. 1, pp. 21–27, 1967.
- [15] S. Mor-Yosef, A. Samueloff, B. Modan, D. Navot, and J. G. Schenker, "Ranking the risk factors for cesarean: Logistic regression analysis of a nationwide study," *Obstetrics and Gynecology*, vol. 75, pp. 944–947, 1990.
- [16] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone, *Classification and regression trees*. Boca Raton,FL: Chapman and Hall/CRC, 1984.
- [17] L. Breiman, "Random forests," *Machine Learning*, vol. 45, pp. 5–32, 2001.
- [18] V. N. Vapnik, *Statistical Learning Theory*. New York: Wiley, 1998.
- [19] B. G. Batchelor, *Pattern recognition: Ideas in practice*. New York: Plenum Press, 1978.
- [20] X. Wu, V. Kumar, Q. J. Ross, J. Ghosh, Q. Yang, H. Motoda, G. J. McLachlan, A. Ng, B. Liu, P. S. Yu, Z. H. Zhou, M. Steinbach, D. J. Hand, and D. Steinberg, "Top 10 algorithms in data mining," *Knowledge and Information Systems*, vol. 14, pp. 1–37, 2008.
- [21] C. Bishop, *Pattern Recognition and Machine Learning*. Springer, 2006.
- [22] M. Schmidt, N. Le Roux, and F. Bach, "Minimizing finite sums with the stochastic average gradient," *Mathematical Programming*, vol. 162, pp. 83–112, 2017.
- [23] I. Guyon, B. Boser, and V. Vapnik, "Automatic capacity tuning of very large VC-dimension classifiers," in *Advances in Neural Information Processing Systems: Morgan Kaufmann*, 5, C. E. Hanson, S., Cowan, J., Giles, Ed., vol. 5. San Mateo, CA: Morgan Kaufmann Publishers Inc., 1993, pp. 147–155.
- [24] C. Cortes and V. Vapnik, "Support-vector networks," *Machine Learning*, vol. 20, pp. 273–297, 1995.
- [25] J. C. Platt, "Sequential minimal optimization: A fast algorithm for training support vector machines," Tech. Rep., 1998.
- [26] R. E. Fan, P. H. Chen, and C. J. Lin, "Working set selection using second order information for training support vector machines," *Journal of Machine Learning Research*, vol. 6, pp. 1889–1918, 2005.
- [27] T. Wang, C. Z. Yan, X. Song, and J. L. Xie, "Monitoring recent trends in the area of aeolian desertified land using Landsat images in China's Xinjiang region," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 68, pp. 184–190, 2012.
- [28] M. Drusch, U. Del Bello, S. Carlier, O. Colin, V. Fernandez, F. Gascon, B. Hoersch, C. Isola, P. Laberinti, P. Martimort, A. Meygret, F. Spoto, O. Sy, F. Marchese, and P. Bargellini, "Sentinel-2: ESA's optical high-resolution mission for GMES operational services," *Remote Sensing of Environment*, vol. 120, pp. 25–36, 2012.
- [29] G. Navarro, I. Caballero, G. Silva, P. C. Parra, Á. Vázquez, and R. Caldeira, "Evaluation of forest fire on Madeira Island using Sentinel-2A MSI imagery," *International Journal of Applied Earth Observation and Geoinformation*, vol. 58, pp. 97–106, 2017.
- [30] C. Li, J. Wang, L. Wang, L. Hu, and P. Gong, "Comparison of classification algorithms and training sample sizes in urban land classification with landsat thematic mapper imagery," *Remote Sensing*, vol. 6, pp. 964–983, 2014.
- [31] C. Song, C. E. Woodcock, K. C. Seto, M. P. Lenney, and S. A. Macomber, "Classification and change detection using Landsat TM data: When and how to correct atmospheric effects?" *Remote Sensing of Environment*, vol. 75, pp. 230–244, 2001.
- [32] P. Gong and P. J. Howarth, "Land-use classification of SPOT HRV data using a cover-frequency method," *International Journal of Remote Sensing*, vol. 13, pp. 1459–1471, 1992.
- [33] P. Gong, D. J. Marceau, and P. J. Howarth, "A comparison of spatial feature extraction algorithms for land-use classification with SPOT HRV data," *Remote Sensing of Environment*, vol. 40, pp. 137–151, 1992.
- [34] C. C. Chang and C. J. Lin, "LIBSVM: A Library for support vector machines," *ACM Transactions on Intelligent Systems and Technology*, vol. 2, pp. 1–27, 2011.
- [35] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.



**Xiaorou Zheng** is currently pursuing the joint Ph.D. degree in computer application technology with the Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen, China, and the Shenzhen College of Advanced Technology, University of Chinese Academy of Sciences, Shenzhen. Her current research interests include land use and land cover classification, semi-supervised learning, weakly supervised learning, fully convolutional networks and generative adversarial networks.



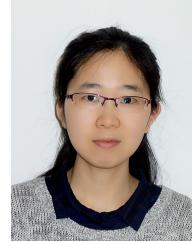
**Jianxin Jia** received the Ph.D. degree in electronic science and technology from the University of Chinese Academy of Sciences, Beijing, China, in 2018. He is currently a Research Scientist with the Finnish Geospatial Research Institute. His research interests include remote sensing imaging system design, hyperspectral image processing and application.



**Shanxin Guo** received the Ph.D. and master's degree in photogrammetry and remote sensing from Wuhan University, Hubei, China. He joined a visiting Ph.D. program with the Department of Geography, University of Wisconsin-Madison, USA. His current research interests include digital soil mapping with remote sensing data, deep learning for land cover and land use mapping, and the gaussian process for forest traits mapping.



**Jinsong Chen** received the Ph.D. degree from the Chinese Academy of Sciences, Shenzhen, China, in 2004. From 2004 to 2011, he was a Post-Doctoral Researcher and a Professor with the Chinese University of Hong Kong, Hong Kong. He joined the Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, as the Director of the Center for Spatial Information Science and Systems. His current research interests include spatial and environmental big data processing methods and multisource spatial and environmental data assimilation and information fusion methods. His current research interests include digital soil mapping with remote sensing data, deep learning for land cover and land use mapping, and the gaussian process for forest traits mapping.



**Luyi Sun** received her B.S. and M.Sc degrees from Beijing Institute of Technology, Beijing, China, in 2009 and 2012, respectively, and the Ph.D. degree from University College London, London, United Kingdom, in 2017. She is currently with the Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, Shenzhen, China. Her research interests include deformation monitoring of the earth surface using microwave remote sensing, as well as land cover mapping based on integration of Synthetic Aperture Radar (SAR) and multi-spectral imagery.



**Yingfei Xiong** received the B.E. degree in remote sensing science and technology from China University of Mining and Technology (Beijing), Beijing, China, in 2018. He is currently pursuing the master's degree with the Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen, China. His research interests include image fusion and super-resolution.



**Wennia Xu** received the master's degree in Shenzhen Institutes of Advanced Technology, Chinese Academy of Science, Shenzhen, China.

Her current research interests include computer vision, digital image process and intelligent speech recognition.