



Semi-supervised regression via embedding space mapping and pseudo-label smearing

Liyan Liu^{1,2} · Jin Zhang^{1,2} · Kun Qian^{1,3} · Fan Min^{1,2,4}

Accepted: 14 July 2024 / Published online: 19 July 2024

© The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2024

Abstract

Co-training is a semi-supervised algorithm that aims to improve prediction effects by exchanging confident instances and pseudo-labels among multiple learners. One central issue is how to mitigate the negative impact of low-quality pseudo-labels during training. In this paper, we propose semi-supervised regression via embedding space mapping and pseudo-label smearing (S2RMS) to ensure that unlabeled data contribute positively to the prediction process. First, a Triplet neural network is trained using pairwise data generated from labeled data. This network maps the training data to the embedding space to better separate dissimilar instances. Second, the embedded data are randomly partitioned into different subsets to train corresponding regression models (a.k.a. regressors). These regressors are integrated into the prediction process. Third, these subsets are augmented using unlabeled data with high similarity to the labeled data and high-confidence pseudo-labels. Here, the similarity and confidence are calculated using the network and the smearing technique, respectively. Experiments are conducted on fourteen datasets, and the results are compared to those of three excellent algorithms. The results show that S2RMS outperforms other co-training and metric semi-supervised regression algorithms. The source code is available at: <https://github.com/BetaCatPro/S2RMS>.

Keywords Co-training · Deep metric learning · Embedding space · Semi-supervised regression

1 Introduction

Semi-supervised learning can enhance the performance of a learning algorithm by utilizing an extensive amount of unlabeled data in real-world tasks [1, 2]. Several semi-supervised

methods have been developed, including co-training [3, 4], self-training [5, 6], consistency regularization [7, 8], and holistic methods [9, 10]. The co-training approach can be described as a disagreement-based technique for semi-supervised learning [11]. During the co-training process, two separate learners are trained on different views, and they select the unlabeled instances about which they are most confident. The unlabeled instances of one learner are then given pseudo-labels and added to the training set of the other learner, which uses them to update its model parameters. Unlike graph and kernel methods [12, 13], co-training does not require additional complex computations or large amounts of memory.

Two main problems are encountered when conducting co-training. First, some studies have indicated that casually using unlabeled data may result in a decrease in model effectiveness [14]. This is because not all unlabeled data can be used for training purposes. It is difficult to assign accurate pseudo-label to some of these unlabeled instances [15, 16]. If unlabeled instances carrying incorrect pseudo-labels continue to be added to the training set, the predictive performance of the utilized learner can be seriously disturbed [17].

✉ Fan Min
minfan@swpu.edu.cn

Liyan Liu
liuliyang@swpu.edu.cn

Jin Zhang
zj20162325@163.com

Kun Qian
kqian2020@163.com

¹ School of Computer Science and Software Engineering, Southwest Petroleum University, Chengdu 610500, China

² Lab of Machine Learning, Southwest Petroleum University, Chengdu 610500, China

³ School of Sciences, Southwest Petroleum University, Chengdu 610500, China

⁴ Institute for Artificial Intelligence, Southwest Petroleum University, Chengdu 610500, China

The situation described above is comparable to the introduction of noisy data to the training set, which is beyond our control. Second, weaker learners struggle to assign accurate pseudo-labels to unlabeled data during the initial stages of training, resulting in a failure to satisfy the expected confidence level [18]. These unsuitable instances are also added to the training set, which leads to the accumulation of errors until the end of the training process.

In this paper, we present semi-supervised regression via embedding space mapping and pseudo-label smearing (S2RMS). This approach combines co-training with a deep metric neural network to handle the above issues. To address the first problem, we use a metric network to learn a decision function that can effectively utilize unlabeled data. To address the second problem, we employ a pseudo-label smearing technique to train robust regressors. Figure 1 shows the framework of S2RMS.

Figure 1 shows that S2MRS selects multiple regressors to build a co-training model. In this study, a particular set of regressors is chosen to test the validity of S2RMS. To better cope with different environments, S2RMS uses random forests [19] possessing varying numbers of evaluators

to form the base regressors. Each regressor also has different parameters to enhance the variability between different models.

Under the assumption of smoothing [1], instances located in high-density areas that are close to each other typically have similar labels, while those that are further apart are more likely to be distinct. Based on this hypothesis, a regression model trained with a small amount of labeled data will have a higher prediction accuracy when confronted with unlabeled data that behave similarly to the labeled data.

First, a pairwise dataset is generated for all anchor instances in the labeled dataset, and this dataset comprises anchor-positive and anchor-negative instance pairs. Subsequently, a deep metric Triplet neural network is trained using this pairwise dataset to learn embedding mappings. Similar instances are located close to each other in this embedding space. The Triplet neural network is used to map labeled and unlabeled data into an embedding space. Second, k subsets are created by randomly sampling the transformed labeled data. Each subset is divided into two parts, with one part used for training the regressors and the other used to calculate the confidence level. The regressors are initialized with

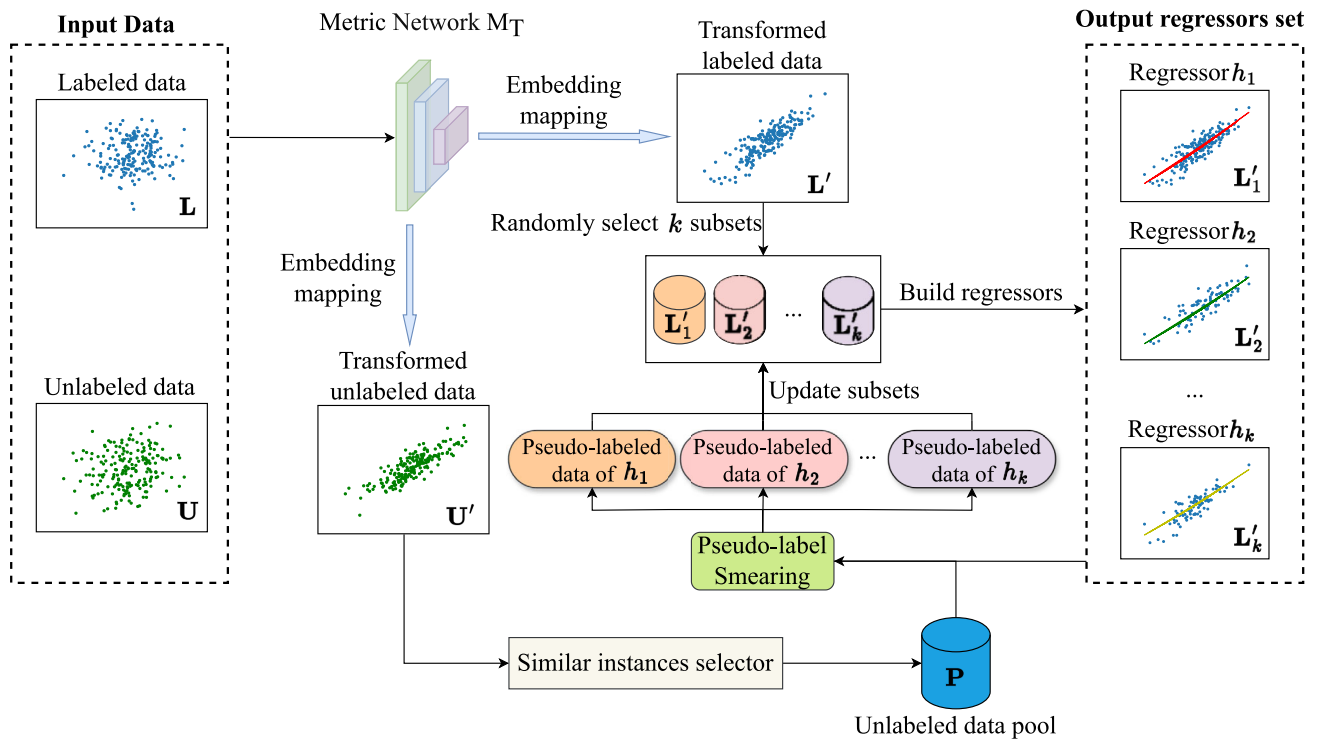


Fig. 1 The framework of S2RMS. The deep metric network M_T used in S2RMS maps the labeled data L and unlabeled data U into a d -dimensional embedding space, which facilitates subsequent operations. To establish the co-training framework, k subsets $\{L'_i\}_{i=1}^k$ are randomly and repeatedly selected from the transformed L' and used to initialize k regressors $\{h_j\}_{j=1}^k$. During the co-training iterations, instances are

identified from the transformed U' using a similar instances selector and added to the data pool P . The co-regressor then employs a label smearing technique within P to ensure representation stability and assigns appropriate pseudo-labels for the unlabeled instances. Instances with pseudo-labels are added to the training set of the corresponding regressor, which is then retrained

different parameters to enhance the diversity and ensure the efficacy of S2RMS. In each iteration, the Triple neural network functions as an instance selector, choosing two similar unlabeled instances for each labeled instance. These similar instances are then added to the pool. The other regressors use pseudo-label smearing to select stable instances from the pool for another iteration. Afterward, these stable instances are assigned pseudo-labels with confidence. Finally, all k training subsets are augmented with these stable instances. All regressors are trained again with the updated data before moving on to the next iteration. After completing all training processes, the unknown data are predicted through averaging.

The effectiveness of the S2RMS algorithm is evaluated through experiments conducted on 14 benchmark datasets. In addition, these datasets are normalized to better visualize the experimental results. We compare the S2RMS algorithm with three excellent semi-supervised regression algorithms. These algorithms also use co-training or metric networks. The results indicate that the S2RMS algorithm outperforms these algorithms.

The remainder of this paper is organized as follows. In Section 2, popular metric learning and semi-supervised learning approaches are reviewed. Section 3 presents the proposed S2RMS algorithm. The spatial mapping method and the smearing technique are described in this section. Section 4 provides the experimental results and analysis. In Section 5, we conclude the paper and outline further work directions.

2 Related works

This section reviews two fundamental techniques that are used in S2RMS: deep metric learning and co-training.

2.1 Deep metric learning

Metric learning aims to measure the similarity between instances using an optimal distance metric for performing the learning task. The main objective is to learn a new metric that reduces the distances between instances belonging to the same class and increases the distances between instances belonging to different classes [20]. Traditional metric learning methods typically use a linear projection strategy, which may not be suitable for real-world problems with nonlinear structures. In deep metric learning, this problem is solved by activation functions with nonlinear structures [21].

The development of deep metric learning is primarily reflected in the available loss functions. The two most basic loss functions in deep metric learning are Contrastive Loss [22] and Triplet Loss [23]. Many variants have been derived from these two functions, such as the Quadruple Loss [24], ArcFace Loss [25], and Ranked List Loss [26]. Furthermore,

deep metric learning has been primarily applied in three areas: 1) image retrieval, 2) person re-identification (ReID), and 3) face recognition. In the field of image retrieval, some studies [27, 28] have examined tasks such as visual sentiment analysis and massive image retrieval by designing different similarity losses. Constrained by the limitations of manual annotation, TSSML [29] improves the hard instance mining process through curriculum learning for semi-supervised ReID tasks. In the facial recognition task, DDFMs [30] use a novel deep metric network approach based on the triplet loss for facial recognition/verification and visual image classification purposes.

2.2 Co-training

Co-training is a divergence-based approach for semi-supervised learning that is designed for multi-view data [3]. This method assumes that data can be effectively represented through various perspectives and that distinct learners can be trained on each view. These learners then label unlabeled instances and incorporate their most confident unlabeled instances into the training sets of other learners. Following each iteration, the learners undergo retraining with the augmented training set.

In the field of co-training, research has focused primarily on three areas: 1) classification, 2) regression, and 3) image segmentation and recognition. In classification tasks, Tri-training [31] utilizes three classifiers with a single view and assigns pseudo-labels to a classifier only if the other two classifiers agree on the label of the example. SIVLC [32] addresses the multi-view challenges encountered in co-training by utilizing sufficient irrelevant views and label consistency. Deep Co-Training [33] successfully constructs deep collaborative classification networks on two views by imposing view difference constraints. Tri-net [34] enhances disparities by generating three other distinct classification models from a shared model with a single view.

In regression tasks, COREG [35] uses two K -nearest neighbor (K -NN) [36] models with different distance measures and K values for regression analysis. CoBCReg [37] uses three RBFNN [38] regressors based on COREG to perform confident instance selection through a committee mechanism. In addition, by integrating a safer and self-paced learning process [39], the collaborative regression model can learn security pseudo-labels with higher quality, resulting in improved model performance. In real-world scenarios, E-CoGRF [40] combines co-training and a random forest model with information entropy grouping to assess the severity levels of depressive symptoms.

In image processing tasks, particularly in the field of medical images, Co-DeepSVS [41] applies support vector regression to multi-view deep co-training networks for survival time estimation. UCMT [42] uses the mean teacher model as a base co-learner while incorporating uncertainty-

guided region mixing for medical image segmentation. Some studies [43] have improved the semantic segmentation performance achieved for medical images by fusing collaborative training with adversarial training. In addition, in the field of remote sensing images, RSCoTr [44] combines multi-task learning and co-training to construct a transformer that can simultaneously handle three tasks: classification, segmentation and detection.

In addition to the previously mentioned research areas, co-training is also commonly employed in subfields such as multi-label learning [45], time series learning [46], and class-imbalanced learning [47]. Furthermore, the amalgamation of co-training with other semi-supervised techniques can significantly enhance the performance achieved in various tasks [48, 49].

3 The proposed algorithm

In this section, we present the implementation details of the S2RMS algorithm. The core techniques of the S2RMS model include embedding space mapping, a similarity-based instance selection strategy and pseudo-label smearing. Finally, the process details of the S2RMS algorithm are given. Table 1 details the notations used in the paper.

Table 1 Notations

Notation	Meaning
\mathbf{L}	The labeled dataset
\mathbf{U}	The unlabeled dataset
n_l	The number of labeled instances
n_u	The number of unlabeled instances
h_i	The i -th regressor
k	The number of regressors
\mathbf{L}'	The transformed labeled dataset
\mathbf{U}'	The transformed unlabeled dataset
x	The instance
d	The dimensionality of the embedding space
f^d	The Triplet neural network
\mathcal{L}	The loss function
\mathbf{P}	The unlabeled data pool
τ	The similarity threshold
$\mathbf{L}'_i, \mathbf{C}'_i$	The subsets of the transformed labeled dataset
μ	The standard deviation threshold
δ	The validation error
Δ_x	The confidence of x
Ω_i	The i -th set of instances with highest confidence

3.1 The model

Let $\mathbf{L} = \{(x_i, y_i)\}_{i=1}^{n_l}$ denote the labeled dataset, and $\mathbf{U} = \{x_j\}_{j=n_l+1}^{n_l+n_u}$ the unlabeled dataset, where n_l and n_u represent the numbers of labeled and unlabeled instances, respectively. Let $\mathbf{H} = \{h_j\}_{j=1}^k$ denote a set of k regressors and \mathcal{L} be the loss function using mean square error (MSE). A general approach using co-training is to train k learners using the original labeled data \mathbf{L} . Then, the unlabeled data \mathbf{U} are selected to facilitate mutual learning among the learners and improve the performance of the algorithm. In our algorithm, \mathbf{L} and \mathbf{U} are mapped to \mathbf{L}' and \mathbf{U}' , respectively, by f^d to better select unlabeled data for training. f^d represents a metric-based neural network that estimates similarity and maps data to the embedding space, with d denoting the dimensionality of the embedding space. Furthermore, the label smearing technique is proposed as a means of obtaining stable instances. The regressors then pass through these stable instances to smooth the predicted output. Then, the objective function can be expressed as follows:

$$\min_{\mathbf{H}} E = \sum_{j \in \{1, \dots, k\}} \left(\sum_{x_l \in \mathbf{L}'_j} \mathcal{L}(h_j(x_l), y_l) + \sum_{x_s \in \Omega_j} \mathcal{L}(h_j(x_s), y_s) \right), \quad (1)$$

where \mathbf{L}'_j is the j -th labeled subset and Ω_j is the j -th set of stable instances with high confidence. To optimize the training process, a pool of unlabeled instances is constructed using a similarity metric strategy. Each optimization iteration applied to the regressors is performed using the data in this pool.

3.2 Embedding space mapping

To address the first problem of co-training, an effective approach must employ an improved metric for instance selection. The typical approach involves randomly selecting a specific number of unlabeled instances and assigning them pseudo-labels. This approach overlooks the distribution of the given data, potentially resulting in the selection of noisy data that can contaminate the training set. Considering the smoothing assumption, the selected unlabeled data should conform to the same distribution as that of the labeled data. Next, a metric method is employed to find unlabeled data that are close to the labeled data. One potential approach is to use the K -NN algorithm, which enables the direct calculation of such metrics. However, constrained by the manifold assumption, we use f^d to replace the simple linear metric methods. For f^d , we utilize a triplet neural network. This

network consists of three identical subnetworks that receive an input pair and extract instance features. These features are then connected to form features that are associated with the input pair. The final layer of the network estimates the similarity with the input pairs.

A pairwise dataset $\{(x_i, \mathbf{P}_{x_i}, \mathbf{N}_{x_i})\}_{i=1}^{n_l}$ is generated from \mathbf{L} to train the Triplet neural network. Specifically, three objects derived from the pairs are the anchor x_i , positive instance set \mathbf{P}_{x_i} and negative instance set \mathbf{N}_{x_i} . A positive instance has high similarity to the anchor, and a negative instance has the opposite relationship. For the anchor $x_a \in \mathbf{L}$, the positive instance set \mathbf{P}_{x_a} and the negative instance set \mathbf{N}_{x_a} are constructed by d_{ij} . The initial similarity d_{ij} is calculated between instances x_i and x_j , while $i \neq j$ is calculated by the Euclidean distance measures. The numbers of positive and negative instances are set to 2.

After generating the pairwise dataset for all anchor instances, the Triplet neural network exploits $n_l(n_l - 1)$ pairs of instances. Then, the Ranked List Loss function \mathcal{L}_{RLL} [26] is employed as follows:

$$\begin{aligned}\mathcal{L}_P(x_a, \mathbf{P}_{x_a}) &= \sum_{x_j \in \mathbf{P}_{x_a}} \frac{\exp(\varepsilon(d_{aj} - (\alpha - M)))}{\sum_{x_j \in \mathbf{P}_{x_a}} \exp(\varepsilon(d_{aj} - (\alpha - M)))} (d_{aj} - (\alpha - M)), \\ \mathcal{L}_N(x_a, \mathbf{N}_{x_a}) &= \sum_{x_j \in \mathbf{N}_{x_a}} \frac{\exp(\varepsilon(\alpha - d_{aj}))}{\sum_{x_j \in \mathbf{N}_{x_a}} \exp(\varepsilon(\alpha - d_{aj}))} (\alpha - d_{aj}), \\ \mathcal{L}_{\text{RLL}} &= (1 - \gamma)\mathcal{L}_P(x_a, \mathbf{P}_{x_a}) + \gamma\mathcal{L}_N(x_a, \mathbf{N}_{x_a}),\end{aligned}\quad (2)$$

where α is the boundary parameter, ε is the temperature parameter, and M is the margin value, $\mathcal{L}_P(x_a, \mathbf{P}_{x_a})$ and $\mathcal{L}_N(x_a, \mathbf{N}_{x_a})$ represent the positive and negative instance losses, respectively. γ is the weight that controls the percentage of losses. In our experiments, we treat both targets equally and set γ to 0.5. In the feature space of the Triplet neural network, \mathcal{L}_{RLL} brings positive instances closer to each other than the negatives are by a margin of M .

The output of the Triplet neural network depends on the form of the input data. In the network, the last layer is des-

ignated as the similarity output layer, while the output of the penultimate hidden layer is employed as the embedding space. In the case of a single input, the network generates the deep embedding space of the input data. For pairs of inputs, the network outputs the similarity between pairs of inputs. In particular, when labeled or unlabeled data is employed as inputs, the network returns the data after the metric mapping. Furthermore, for pairs of inputs comprising labeled and unlabeled data, the network outputs the similarity between pairs of data. Consequently, the Triple neural network is capable of executing metric space mapping and similarity measuring operations.

Once the training of the Triple neural network has been completed, \mathbf{L} and \mathbf{U} are mapped to \mathbf{L}' and \mathbf{U}' through the network. For each instance in \mathbf{L}' , the similarity between it and any instance in \mathbf{U}' is computed by the network. This similarity is then compared to a threshold value, τ , and the two highest unlabeled instances with similarity below this threshold are selected. The similarity threshold τ is used to control the bounds of the similar instance selection procedure. These instances are used to construct an unlabeled data pool. In each iteration, the Triplet neural network selects a batch consisting of new data to update the pool according to τ . The threshold τ decreases with the number of iterations as follows:

$$\tau_{t+1} = \lambda\tau_t, \quad (3)$$

where t is the number of iterations and λ is the attenuation coefficient. The selected data are then removed from \mathbf{U}' . Figure 2 shows the unlabeled instances selection process.

3.3 Co-training with pseudo-label smearing

To address the second problem encountered during co-training, an effective training framework is proposed. k pairs of subsets $(\mathbf{L}'_i, \mathbf{C}'_i)$, $i \in \{1, \dots, k\}$ are randomly and repeatedly selected from \mathbf{L}' . \mathbf{L}'_i is used to initialize regressor h_i , and

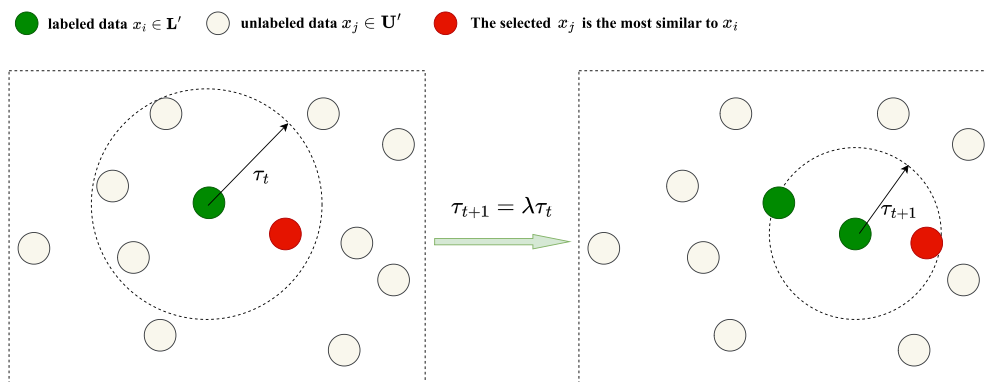


Fig. 2 The unlabeled instances selection process. For each embedded labeled instance $x_i \in \mathbf{L}'$, the Triplet neural network selects the most similar instances from the embedded set \mathbf{U}' according to the threshold τ . The threshold τ decreases with the number of iterations

C'_i is used as part of the confidence calculation. Specifically, C'_i is not included in the training data.

In each iteration, other regressors assign pseudo-labels to all instance data in the unlabeled data pool. Since the pseudo-labels assigned to the unlabeled data may be inaccurate, adding them to the training set of the current regressor will degrade the performance of this regressor. Utilizing these erroneous pseudo-labels is equivalent to introducing noise to the training set. To address this issue, we employ a smearing technique to assign stable pseudo-labels to the unlabeled data. The purpose of label smearing is to select stable instances for improving the robustness of the regressors by smoothing out the differences between different collaborative regressors. At its core, label smearing calculates the variance of the predicted values of the unlabeled data, treats the unlabeled data with smaller variances as stable instances, and smooths out the differences between the predictions of the different regressors by averaging them.

For regressor h_i , the other regressors $h_{j, j \neq i}$ compute the predicted values for all instances in the unlabeled data pool. Then, the stable instances are filtered from these instances through a stability assessment. Stability is expressed as the variance of the values predicted by these regressors for the instances, as follows:

$$STD_i(x_u) = \sqrt{\frac{1}{k-1} \sum_{j=1, j \neq i}^k (h_j(x_u) - \frac{1}{k-1} \sum_{j=1, j \neq i}^k h_j(x_u))^2}, \quad (4)$$

where x_u is an instance in the unlabeled data pool. In the context of the S2RMS algorithm, the smearing strategy involves calculating the stability of the unlabeled instances based on the variance of the predictions of the regressors. The instances with less variance are then selected as stable instances, and their predicted values are replaced with the mean prediction value, which serves as their corresponding pseudo-labels. Once stable instances are obtained and pseudo-labeling is completed, confidence calculation can be performed with h_i as follows:

$$\Delta_{x_u} = \frac{\delta - \delta'}{\delta + \delta'}, \quad (5)$$

where δ and δ' represent the validation errors induced on C'_i prior to and following the update of $\left(x_u, \frac{1}{k-1} \sum_{j=1, j \neq i}^k h_j(x_u)\right)$, respectively. In particular, the validation error in the algorithm is defined as the root mean square (RMSE). Instances with high confidence expand the training set of h_i , and the corresponding regressor is retrained.

Algorithm 1 S2RMS algorithm.

input : Labeled dataset $\mathbf{L} = \{(x_i, y_i)\}_{i=1}^{n_l}$, unlabeled dataset $\mathbf{U} = \{x_j\}_{j=n_l+1}^{n_l+n_u}$, the initial value of τ , the maximum number of iterations T , the value of μ , the number of regressors k , the number of confident instances R .

output: Prediction regressor set $\mathbf{H} = \{h_j\}_{j=1}^k$.

```

1 Generate pairwise data from  $\mathbf{L}$ ;
2 Use the pairwise data to learn a Triplet neural network  $f^d$  according to (2);
3  $\mathbf{L}' \leftarrow f^d(\mathbf{L})$ ,  $\mathbf{U}' \leftarrow f^d(\mathbf{U})$ ;
4 for  $i \leftarrow 1$  to  $k$  do
5    $(\mathbf{L}'_i, \mathbf{C}'_i) \leftarrow \text{RandomSelect}(\mathbf{L}')$ ,  $h_i \leftarrow \text{Learn}(\mathbf{L}'_i)$ ;
6 for  $t \leftarrow 1$  to  $T$  do
7    $\mathbf{P} \leftarrow \emptyset$ ;
8   for each  $x_l \in \mathbf{L}'$  do
9     Calculate the similarity between  $x_l$  and each  $x_u \in \mathbf{U}'$  using  $f^d$ ;
10    Move the top-2 similar instances from  $\mathbf{U}'$  to  $\mathbf{P}$ ;
11  Decrease the value of  $\tau$  according to (3);
12  while  $\mathbf{P} \neq \emptyset$  do
13    for  $i \leftarrow 1$  to  $k$  do
14      Calculate the validation error  $\delta_i$  of  $h_i$  on  $\mathbf{C}'_i$ ;
15       $\Pi_i \leftarrow \emptyset$ ,  $\Omega_i \leftarrow \emptyset$ ;
16      for each  $x_u \in \mathbf{P}$  do
17        Calculate  $STD_i(x_u)$  according to (4);
18        if  $STD_i(x_u) \leq \mu$  then
19           $\tilde{y}_u \leftarrow \frac{1}{k-1} \sum_{j=1, j \neq i}^k h_j(x_u)$ ,
20           $\Pi_i \leftarrow \Pi_i \cup \{(x_u, \tilde{y}_u)\}$ ;
21           $h'_i \leftarrow \text{Learn}(\{(x_u, \tilde{y}_u)\} \cup \mathbf{L}'_i)$ ;
22          Calculate the validation error  $\delta'_i$  of  $h'_i$  on  $\mathbf{C}'_i$ ;
23           $\Delta_{x_u} \leftarrow \frac{\delta - \delta'}{\delta + \delta'}$ ;
24        for  $r \leftarrow 1$  to  $R$  do
25           $\Omega_i \leftarrow \Omega_i \cup \{\arg \max_{x_u \in \Pi_i} \Delta_{x_u} | \Delta_{x_u} > 0\}$ ;
26         $\mathbf{P} \leftarrow \mathbf{P} \setminus \Omega_i$ ,  $\mathbf{L}' \leftarrow \mathbf{L}' \cup \Omega_i$ ;
27      if  $\bigcup_{i=1}^k \Omega_i = \emptyset$  then
28        break; // No more confident instances.
29      for  $i \leftarrow 1$  to  $k$  do
30         $\mathbf{L}'_i \leftarrow \mathbf{L}'_i \cup \Omega_i$ ,  $h_i \leftarrow \text{Learn}(\mathbf{L}'_i)$ ;
31 return  $\mathbf{H} = \{h_j\}_{j=1}^k$ ;

```

3.4 Algorithm description

Algorithm 1 shows the complete training process of S2RMS. The inputs of S2RMS include the labeled dataset $\mathbf{L} = \{(x_i, y_i)\}_{i=1}^{n_l}$, the unlabeled dataset $\mathbf{U} = \{x_j\}_{j=n_l+1}^{n_l+n_u}$, the initial value of τ , the maximum number of iterations T , the value of μ , the number of regressors k , and the number of confident instances R .

A Triplet neural network f^d is trained to generate pairwise data from \mathbf{L} according to (2). In line 3, the labeled dataset \mathbf{L} and unlabeled dataset \mathbf{U} are mapped by f^d to \mathbf{L}' and \mathbf{U}' , respectively. The transformed data \mathbf{L}' are then randomly split into 2 subsets \mathbf{L}'_i and \mathbf{C}'_i , as shown in line 5. k different

subsets $\mathbf{L}'_i, i \in \{1, \dots, k\}$ are obtained after repeating the random selection procedure k times. These subsets are used to initialize the k regressors $\{h_j\}_{j=1}^k$.

The next step is to train these regressors. At the beginning of each iteration, in lines 7 to 10, an unlabeled data pool \mathbf{P} is created. This pool is generated by selecting unlabeled data using the metric method and the value of the threshold τ with the help of f^d . The threshold τ is then decreased according to (3). For each regressor h_i , the other regressors $h_j, j \neq i$ select the instances with the highest confidence value. First, in line 14, the validation error δ_i of h_i on \mathbf{C}'_i is calculated. Next, for each x_u in \mathbf{P} , its stability value is calculated according to (4). In lines 17 to 20, the instances x_u that satisfy the stability threshold μ receive pseudo-labels and are temporarily merged with the original training set \mathbf{L}'_i to train the regressor h'_i . Then, in line 21, the validation error δ_i of h'_i on \mathbf{C}'_i is calculated. In line 22, the confidence Δ_{x_u} of x_u is calculated according to (5). As shown in line 24, the instances x_u with the highest confidence values are selected and added to Ω_i . Eventually, at most R instances are contained in Ω_i . In line 29, the instances of Ω_i are integrated into the training set of h_i . Regressor h_i is subsequently retrained on the updated training set.

A set of regression models $\mathbf{H} = \{h_j\}_{j=1}^k$ is obtained for data prediction purposes at the end of the S2RMS training process. The predicted label of x_u is then given by

$$\hat{y}_{x_u} = \frac{1}{k} \sum_{j=1}^k h_j(x_u). \quad (6)$$

S2RMS controls the differences among the predictions of learners by using label smearing during the training process. This ensures that the prediction results of the learners are stable without the need for weighted averaging.

3.5 Complexity analysis

The co-training process is where the complexity of the entire algorithm lies. This stage involves processing each x_u contained in the unlabeled data pool individually. In the S2RMS algorithm, we use random forests as the underlying co-regressors. These random forest regressors are executed in parallel during the iterative process. Therefore, we multiply the dominant term of the corresponding time complexity by the number of iterations to form the final complexity representation.

The random forest model typically employs a classification and regression tree (CART) as its base model. The time complexity of constructing a single CART is denoted by $O(I_{\text{ins}} \times A_{\text{attr}} \times \log I_{\text{ins}})$, where I_{ins} denotes the number of instances and A_{attr} denotes the number of attributes. Therefore, the time complexity of constructing a random

forest model is $O(T_{\text{trees}} \times I_{\text{ins}} \times A_{\text{attr}} \times \log I_{\text{ins}})$, where T_{trees} represents the number of trees. As the computation process involves random feature selection, additional time is required to handle this process. Consequently, the average time complexity of constructing a random forest model is $O(T_{\text{trees}} \times I_{\text{ins}} \times A_{\text{attr}} \times \log I_{\text{ins}})$.

Co-training involves calculating confidence values and selecting instances based on these confidence levels. During the confidence calculation, additional random forest regressors are constructed for each stable instance x_u in the unlabeled data pool \mathbf{P} . When selecting confident instances, the time cost is only $O(R)$, where R is the number of confident instances. This cost can be ignored compared to the cost of the confidence calculation. Therefore, the total time complexity is

$$O(k \times N \times S \times T_{\text{trees}} \times I_{\text{ins}} \times A_{\text{attr}} \times \log I_{\text{ins}}),$$

where N is the maximum number of iterations and S is the number of stable instances. In the worst-case, all unlabeled instances are stable instances, hence the total time complexity is

$$O(k \times N \times |\mathbf{P}| \times T_{\text{trees}} \times I_{\text{ins}} \times A_{\text{attr}} \times \log I_{\text{ins}}),$$

where $|\mathbf{P}|$ is the cardinality of \mathbf{P} and $|\mathbf{P}| \geq S$.

4 Experiments

To validate the effectiveness of our approach, we employ S2RMS to conduct a diverse range of regression tasks involving inductive semi-supervised learning. First, we test it on a total of fourteen datasets obtained from machine learning repositories (e.g., UCI, Delve, and Statlib). We calculate the RMSE values of S2RMS and the comparison methods. The effectiveness of similarity-based instance selection strategy and pseudo-label smearing techniques are evaluated through ablation experiments. All the experiments are performed on PyTorch [50].

4.1 Experimental settings

The detailed configuration of the experimental setup is as follows.

Datasets Table 2 lists the fourteen real-world datasets used in our experiments. For each dataset, we randomly select 2000 instances as the training set, while the remaining instances are used for testing. The training set is divided into labeled and unlabeled portions based on a specific ratio. In our experiments, 0.5%, 2.5% and 5% of the training set are used as labeled data in the separate scenarios. In most cases (except Section 4.3.4), 20% of the training set is used for validation.

Table 2 Datasets information

Name	Instances	Features	Source
Abalone	4 177	8	UCI
Bank32nh	8 192	32	Delve
Cal_housing	20 460	8	StatLib
Cpu_act	8 192	21	Delve
Cpu_small	8 192	12	Delve
Elevators	9 517	7	UCI
Folds5x2_pp	9 565	4	UCI
Kin8nm	8 192	8	Delve
Parkinsons	5 875	21	UCI
Pollen	3 848	5	StatLib
Puma8NH	8 192	8	UCI
Space_ga	3 107	6	StatLib
Wind	6 574	14	StatLib
Wine_quality	6 497	11	UCI

As an illustration, let us consider the Abalone dataset. We randomly select 2000 instances as the training set, while the remaining 2177 instances are used for testing. With a label rate of 2.5%, we include 50 instances in the labeled dataset, whereas the remaining 1950 instances are included in the unlabeled dataset. Among them, 10 instances serve as the validation set.

Parameters In this study, we choose a neural network with three hidden layers as Triplet metric model. In addition, we select specific regressors and parameters to suit our purposes. The three networks in the Triplet neural network share the same architecture, and each hidden layer uses the LeakyReLU [51] function to perform nonlinear transformations. The numbers of neurons in the three hidden layers are 16, 32 and 64. In particular, the dimensionality of the input layer is the number of features contained in the input data, while the output of the second hidden layer is taken as the embedding space. For the parameter settings of the Ranked List Loss function, α , ϵ and M are set to 40, -0.5 and 35, respectively.

To increase the diversity among the learners, three random forest regressors are initialized with different parameters. Specifically, the three regressors are based on 100 estimators, 120 estimators and 150 estimators. The minimum instance splits are set to 1, 2 and 3. In the experiments, the initial similarity threshold τ and the number of confident instances R are set to 0.017 and 2, respectively. The attenuation coefficient λ is set to 0.98. Instances with μ less than 0.02 are used for the selection of stable instances. The number of total experiments is set to 30.

Metrics The RMSE and R^2 metrics are chosen to evaluate the effectiveness of the tested algorithms. The RMSE is the root mean square error, which measures the square root of the average of the squared differences between the predicted and ground-truth values divided by the number of observations n_{test} . This metric calculates the difference between the predicted and ground-truth values and is particularly impacted by extreme values in the utilized dataset. The RMSE is calculated as follows:

$$\text{RMSE} = \sqrt{\frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} (y_i - \hat{y}_i)^2}, \quad (7)$$

where y_i represents the ground-truth value, \hat{y}_i represents the predicted value and n_{test} represents the number of instances.

The coefficient of determination R^2 is an indicator used to measure the degree of model fitness. It indicates the proportion of dependent variable variance that the model can explain. R^2 is calculated as follows:

$$R^2 = 1 - \frac{\sum_{i=1}^{n_{\text{test}}} (y_i - \hat{y}_i)^2 / n_{\text{test}}}{\sum_{i=1}^{n_{\text{test}}} (y_i - \bar{y})^2 / n_{\text{test}}}, \quad (8)$$

where \bar{y} represents the predicted average value. R^2 measures the extent to which the model explains the total variance, with a normal range of values between 0 and 1. A value closer to 1 indicates a better level of fit, while a value closer to 0 indicates a worse level of fit. A value less than 0 indicates that the model does not fit well and cannot make predictions.

4.2 Comparison with semi-supervised methods

S2RMS is compared with the following methods to further evaluate its performance.

- **COREG**: COREG uses a fixed K value and different distance metric algorithms to obtain the best experimental results. Two 3-NN regressors based on the Euclidean and Mahalanobis distance metrics are utilized. A new method using the mean squared error (MSE) is employed to calculate the confidence values.
- **SAFER** [17]: The SAFER semi-supervised regression method transforms the task of interest into a geometric projection problem. This approach utilizes multiple regressors to learn security predictions. The problem that unlabeled data may degrade the performance of the utilized algorithm is effectively circumvented.

- **DML-S2R** [52]: DML-S2R aims to learn the metric space of similar instances by efficiently using both unlabeled and labeled data. This method consists of two main steps: 1) pairwise similarity modeling using labeled data and 2) metric learning based on a Triplet network with a large amount of unlabeled data.

Table 3 presents the comparison results produced by S2RMS and other algorithms under different label rates.

The best performance entries in Table 3 are identified with bullets. First, S2RMS is compared with COREG and SAFER, which are based on divergence. At all label rates, S2RMS achieves the lowest average RMSE, while COREG and SAFER exhibit significantly greater average RMSEs than does S2RMS. S2RMS outperforms both COREG and SAFER on average on more than 10 datasets under different label rates. With a 2.5% label rate, S2RMS outperforms COREG on 12 datasets and SAFER on 11 datasets. It is evident that S2RMS outperforms both COREG and SAFER by significant margins.

Second, S2RMS is compared with DML-S2R using the metric method. Table 3 shows that S2RMS outperforms DML-S2R on all datasets under all label rates. The average performance improvement is greater than 13%. With a lower label rate, S2RMS is more effective. The above results show that S2RMS outperforms DML-S2R. In addition, S2RMS performs better in cases with extremely small rates.

Regardless of whether the label rate is 2.5% or 5%, S2RMS has poor performance on the Folds5x2_pp and Bank32NH datasets. The reasons for this result are as follows. The first reason is the inaccurate embedding space mapping process. Due to the number of features and the structural characteristics of the datasets, a fixed embedding space dimensionality may lead to data misalignment. This issue affects the training of the learner. The second reason is the special distribution of the embedding space. Because of the first reason, algorithms such as K -NN perform well in this space, but the performance of random forests decreases significantly. In contrast to COREG and SAFER, S2RMS employs a decision tree for regression purposes. As a result, its use of random forests leads to a decrease in performance compared to that of the other two methods using K -NN.

As S2RMS and COREG share co-training similarities, we adapt the co-learner of S2RMS to adopt the same K -NN regressor as that of COREG. This enables a clearer and more precise comparison between the performance results of the algorithms, as they all use the same regression model. Figure 3 shows the results of the experimental comparisons.

The experimental results indicate that when the base regressor is replaced with the K -NN algorithm, the prediction performance of S2RMS tends to decrease on some datasets relative to that attained with random forests. How-

ever, S2RMS outperforms COREG on most datasets. On the other hand, S2RMS falls behind COREG in all aspects on the Folds5x2_pp and Puma8NH datasets. The primary reason for this outcome is that the Euclidean distance measure, which is used by K -NN, is not suitable for the metric embedding spaces of these two datasets. Additionally, higher RMSEs are observed for the Cal_housing, Space_ga, and Wine_quality datasets, where the labeled data rate is set to 5%. S2RMS has an overall advantage over COREG. The experimental results indicate that the performance of the proposed algorithm is affected by the differences among the base regressors. The random forest-based co-learner outperforms K -NN under the influence of the metric mapping strategy proposed in this paper.

4.3 Ablation study

4.3.1 Similarity-based instance selection strategy

First, the regression model is initialized by labeled data. According to the smoothing assumption in semi-supervised learning, two instances that are close in terms of distance in high-density regions should have similar outputs. Therefore, for an unlabeled instance that is close to a labeled instance, the regression model can predict its target value more accurately. Thus, we construct an unlabeled data pool by selecting instances that are similar or close to the labeled data from the unlabeled data. This unlabeled data pool not only accelerates the iteration speed of the algorithm but also provides available unlabeled data for the subsequent co-training process. In conclusion, we choose a deep metric network to simultaneously complete metric space mapping and similarity-based instance selection. In the metric space, similar instances are close to each other, while instances with large differences are far from each other. Based on this characteristic of the metric space, the selection of similar instances can be better realized, and the effectiveness of the unlabeled data pool is also guaranteed.

We compile a list of the performance improvements (in percentages) resulting from the similarity-based instance selection strategy. This improvement is calculated by subtracting the RMSE of the similarity-based instance selection strategy from the RMSE obtained without using the similarity-based instance selection strategy. Then, the result is divided by the RMSE obtained without similarity-based instance selection. Table 4 shows that the performances achieved by the model on most datasets improve.

In the absence of a similarity-based instance selection strategy, an unlabeled data pool is constructed in accordance with the pool construction strategy outlined in COREG. In particular, the size of the pool is fixed and instances are randomly selected from the unlabeled dataset to be added to the pool until it is filled. After one round of iteration, the unlabeled

Table 3 The values of RMSE produced by the comparison algorithms under different label rates

Ratio	Dataset	COREG	SAFER	DML-S2R	S2RMS
0.5%	Abalone	0.1090 \pm 0.0106	0.1158 \pm 0.0226	0.1153 \pm 0.0146	● 0.1072 \pm 0.0089
	Bank32nh	0.1598 \pm 0.0221	0.1579 \pm 0.0124	0.1642 \pm 0.0086	● 0.1545 \pm 0.0021
	Cal_housing	0.2345 \pm 0.0369	0.2496 \pm 0.0408	0.2674 \pm 0.0286	● 0.1987 \pm 0.0010
	Cpu_act	0.1813 \pm 0.0175	0.1625 \pm 0.0230	0.2033 \pm 0.0067	● 0.1357 \pm 0.0009
	Cpu_small	0.1767 \pm 0.0095	0.1598 \pm 0.0205	0.1789 \pm 0.0067	● 0.1243 \pm 0.0008
	Elevators	● 0.0732 \pm 0.0070	0.0746 \pm 0.0077	0.0755 \pm 0.0067	0.0749 \pm 0.0048
	Folds5x2_pp	0.1224 \pm 0.0129	● 0.1127 \pm 0.0069	0.1489 \pm 0.0039	0.1683 \pm 0.0196
	Kin8nm	0.1862 \pm 0.0116	0.1857 \pm 0.0142	0.1855 \pm 0.0152	● 0.1641 \pm 0.0023
	Parkinsons	0.1269 \pm 0.0038	● 0.1265 \pm 0.0048	0.1374 \pm 0.0027	0.1344 \pm 0.0020
	Pollen	0.3227 \pm 0.0078	0.3287 \pm 0.0078	0.3334 \pm 0.0077	● 0.3179 \pm 0.0096
	Puma8NH	0.2373 \pm 0.0034	● 0.2351 \pm 0.0078	0.2741 \pm 0.0012	0.2419 \pm 0.0192
	Space_ga	0.1985 \pm 0.0139	0.2210 \pm 0.0111	0.1999 \pm 0.0121	● 0.1696 \pm 0.0008
	Wind	0.1427 \pm 0.0094	0.1328 \pm 0.0091	0.1462 \pm 0.0117	● 0.1187 \pm 0.0024
	Wine_quality	0.1463 \pm 0.0141	0.1473 \pm 0.0182	0.1721 \pm 0.0044	● 0.1358 \pm 0.0024
2.5%	Abalone	0.0920 \pm 0.0075	0.0964 \pm 0.0052	0.1066 \pm 0.0052	● 0.0854 \pm 0.0011
	Bank32nh	● 0.1503 \pm 0.0053	0.1527 \pm 0.0054	0.1609 \pm 0.0039	0.1575 \pm 0.0055
	Cal_housing	0.1844 \pm 0.0219	0.1975 \pm 0.0218	0.2245 \pm 0.0071	● 0.1837 \pm 0.0029
	Cpu_act	0.1317 \pm 0.0217	● 0.1079 \pm 0.0231	0.1188 \pm 0.0054	0.1202 \pm 0.0030
	Cpu_small	0.1218 \pm 0.0293	0.0991 \pm 0.0242	0.1749 \pm 0.0055	● 0.0952 \pm 0.0096
	Elevators	0.0645 \pm 0.0014	0.0643 \pm 0.0016	0.0689 \pm 0.0016	● 0.0590 \pm 0.0056
	Folds5x2_pp	0.0835 \pm 0.0044	● 0.0795 \pm 0.0024	0.1109 \pm 0.0017	0.0848 \pm 0.0029
	Kin8nm	0.1612 \pm 0.0035	0.1602 \pm 0.0053	0.1621 \pm 0.0023	● 0.1530 \pm 0.0082
	Parkinsons	0.1047 \pm 0.0040	0.1079 \pm 0.0072	0.1195 \pm 0.0036	● 0.0999 \pm 0.0048
	Pollen	0.3186 \pm 0.0076	0.3207 \pm 0.0075	0.3369 \pm 0.0076	● 0.3026 \pm 0.0054
	Puma8NH	0.2092 \pm 0.0060	0.2081 \pm 0.0042	0.2470 \pm 0.0017	● 0.1946 \pm 0.0082
	Space_ga	0.1725 \pm 0.0117	0.1952 \pm 0.0089	0.1844 \pm 0.0098	● 0.1682 \pm 0.0065
	Wind	0.1185 \pm 0.0175	0.1130 \pm 0.0049	0.1203 \pm 0.0046	● 0.0984 \pm 0.0030
	Wine_quality	0.1381 \pm 0.0121	0.1374 \pm 0.0083	0.1669 \pm 0.0028	● 0.1332 \pm 0.0016
5%	Abalone	0.0877 \pm 0.0037	0.0949 \pm 0.0047	0.1016 \pm 0.0032	● 0.0874 \pm 0.0071
	Bank32nh	● 0.1494 \pm 0.0024	0.1518 \pm 0.0036	0.1591 \pm 0.0030	0.1567 \pm 0.0014
	Cal_housing	● 0.1713 \pm 0.0128	0.1817 \pm 0.0069	0.2093 \pm 0.0075	0.1846 \pm 0.0047
	Cpu_act	0.1065 \pm 0.0153	0.0863 \pm 0.0281	0.0879 \pm 0.0054	● 0.0856 \pm 0.0113
	Cpu_small	0.0911 \pm 0.0203	● 0.0664 \pm 0.0277	0.0902 \pm 0.0054	0.0917 \pm 0.0177
	Elevators	0.0634 \pm 0.0014	0.0628 \pm 0.0012	0.0645 \pm 0.0017	● 0.0599 \pm 0.0021
	Folds5x2_pp	0.0754 \pm 0.0020	● 0.0738 \pm 0.0013	0.0848 \pm 0.0022	0.0834 \pm 0.0032
	Kin8nm	0.1477 \pm 0.0021	0.1510 \pm 0.0030	0.1623 \pm 0.0028	● 0.1474 \pm 0.0016
	Parkinsons	0.0962 \pm 0.0049	0.0995 \pm 0.0072	0.0956 \pm 0.0034	● 0.0928 \pm 0.0012
	Pollen	0.3177 \pm 0.0071	0.3185 \pm 0.0073	0.3293 \pm 0.0074	● 0.3045 \pm 0.0028
	Puma8NH	0.2026 \pm 0.0029	0.2019 \pm 0.0015	0.2029 \pm 0.0018	● 0.1927 \pm 0.0012
	Space_ga	0.1659 \pm 0.0166	0.1852 \pm 0.0071	0.1782 \pm 0.0097	● 0.1658 \pm 0.0221
	Wind	0.1111 \pm 0.0323	0.1079 \pm 0.0051	0.1001 \pm 0.0060	● 0.0996 \pm 0.0031
	Wine_quality	0.1353 \pm 0.0086	0.1347 \pm 0.0071	0.1685 \pm 0.0038	● 0.1330 \pm 0.0005

The best results obtained on each dataset are denoted with bullets

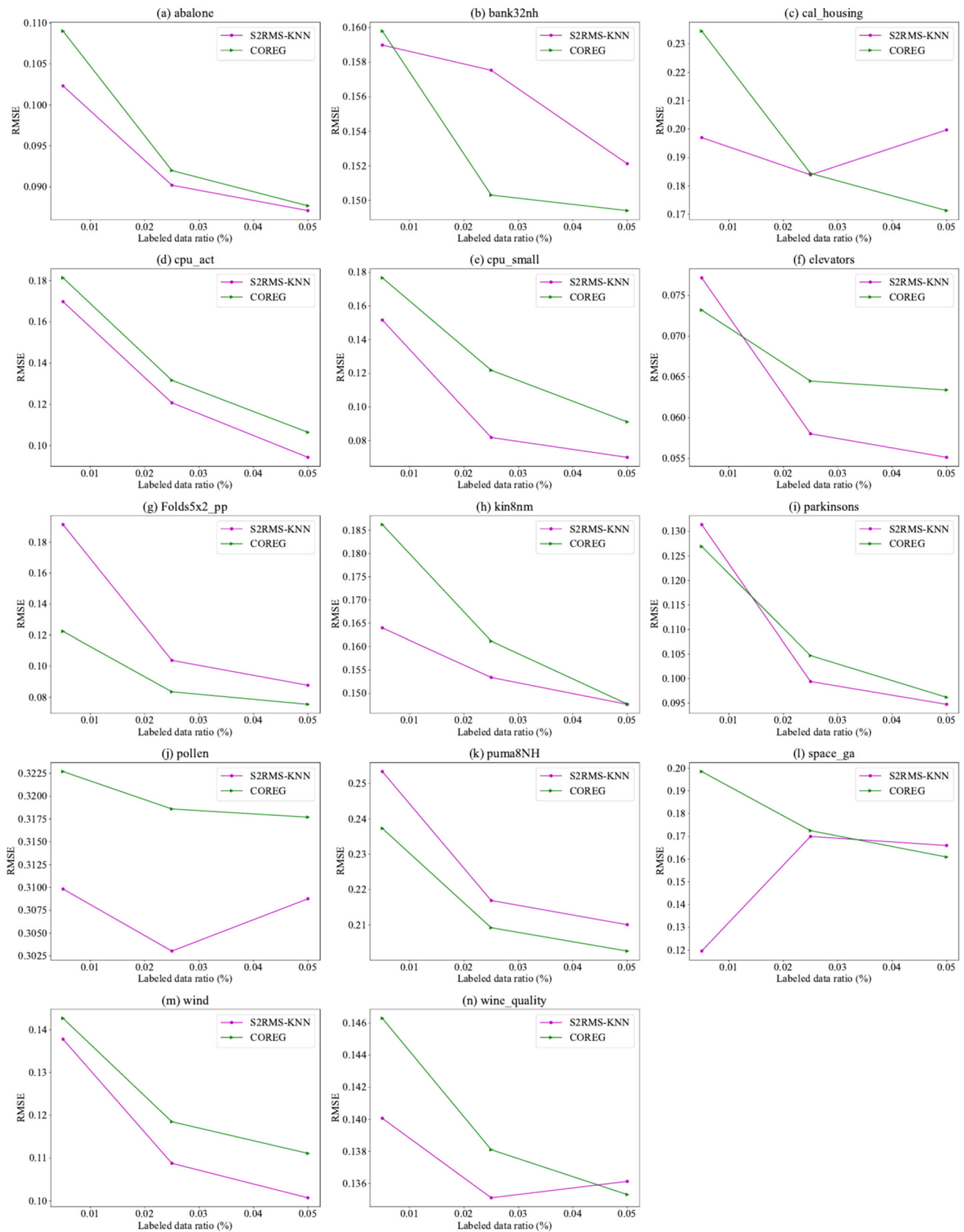


Fig. 3 RMSE comparisons between S2RMS-KNN and COREG conducted on fourteen datasets under different label rates. S2RMS-KNN denotes the S2RMS algorithm constructed after applying K -NN as the base regressor

Table 4 The RMSE values of the two methods were compared, one employing similarity-based instance selection and the other random selection

Dataset	Randomization-based strategy	Similarity-based strategy	Improvements
Abalone	0.0901 \pm 0.0062	● 0.0854 \pm 0.0011	5.22%
Bank32nh	0.1642 \pm 0.0084	● 0.1575 \pm 0.0055	4.14%
Cal_housing	0.1894 \pm 0.0191	● 0.1837 \pm 0.0029	3.01%
Cpu_act	0.1409 \pm 0.0181	● 0.1202 \pm 0.0030	14.69%
Cpu_small	0.0958 \pm 0.0102	● 0.0952 \pm 0.0096	0.62%
Elevators	0.0615 \pm 0.0035	● 0.0590 \pm 0.0056	4.07%
Folds5x2_pp	● 0.0842 \pm 0.0033	0.0848 \pm 0.0029	−0.71%
Kin8nm	0.1601 \pm 0.0118	● 0.1530 \pm 0.0082	4.43%
Parkinsons	0.1009 \pm 0.0024	● 0.0999 \pm 0.0048	0.99%
Pollen	0.3061 \pm 0.0003	● 0.3026 \pm 0.0054	1.14%
Puma8NH	0.2020 \pm 0.0041	● 0.1946 \pm 0.0082	3.66%
Space_ga	0.1750 \pm 0.0162	● 0.1682 \pm 0.0065	3.89%
Wind	0.0986 \pm 0.0041	● 0.0984 \pm 0.0030	0.20%
Wine_quality	0.1335 \pm 0.0064	● 0.1332 \pm 0.0016	0.22%

The label rate of the experiment is set to 2.5%

beled data pool is replenished by randomly selecting again. In this way, the construction process of the unlabeled data pool, as outlined in lines 7 to 11 of Algorithm 1, is replaced by the aforementioned method.

Experiments are conducted under a 2.5% label rate. Each dataset is used for 30 tests implemented under different label rates. All the data from the training process are mapped to the metric space. Table 4 presents the experimental RMSE results obtained before and after using similar selection. The pool size of the random selection process is set to 50. The results show that the similarity-based instance selection strategy always performs better.

The data pool of S2RMS is constructed by selecting unlabeled data that are similar to the labeled data. In the embedding space, similar instances are close to one another. More appropriate unlabeled data can be selected based on the distribution of the labeled data. According to the smoothing and manifold assumptions, these unlabeled data can be effectively used to expand the training set. Thus, the proposed method significantly improves upon the performance of random selection.

However, for the Folds5x2_pp dataset, the performance degrades. One possible reason for this finding is that this dataset has fewer features. After mapping the data to a higher-dimensional embedding space, some ambiguous features cannot correctly represent the data. Therefore, it is not possible to effectively select instances. This, in turn, affects the training process of the regressor.

4.3.2 Embedding space mapping

The COREG and SAFER comparison algorithms use regression and pseudo-labeling without embedding learning. In this

case, we compare the S2RMS algorithm with the COREG and SAFER algorithms after applying the embedding space mapping strategy. Specifically, we additionally apply steps 1 to 3 of Algorithm 1 to the compared algorithms.

We then compare the performances achieved by the COREG and SAFER algorithms before and after adopting the mapping strategy. This approach aims to accurately evaluate the impact of embedding space mapping on the resulting algorithmic performance.

For an example with a 2.5% label rate, Table 5 shows the results obtained in the comparison experiments by the tested methods when using the mapping strategy. The experimental results indicate that S2RMS outperforms the other methods on the 11 datasets, even though COREG and SAFER also implement the mapping strategy.

Table 6 displays the results of comparison between the R^2 values produced before and after applying embedding mapping in the COREG and SAFER algorithms.

Table 6 shows that the two algorithms produce divergent results after applying the mapping strategy. Compared with that of COREG, the performance of COREG-ESP is significantly improved on only 6 datasets and is basically equal on 3 datasets. However, performance degradations are observed on 5 datasets: Bank32nh, Folds5x2_pp, Pollen, Space_ga, and Wine_quality. As discussed in the previous subsections, the metrics used in COREG may not be suitable for metric embedding mappings, resulting in negative shifts on some datasets. In contrast, the embedding mapping strategy works relatively well in the SAFER algorithm. The performance of SAFER-ESP is effectively improved on 12 datasets, with only slight decreases on the Bank32nh and Pollen datasets. Overall, the metric mapping strategy can improve the performance of these algorithms on most datasets.

Table 5 The values of R^2 produced by the comparison algorithms at a label rate of 2.5%

Dataset	COREG-ESP	SAFER-ESP	S2RMS
Abalone	0.2738 \pm 0.0458	0.1132 \pm 0.0734	● 0.3886 \pm 0.0460
Bank32nh	-0.2234 \pm 0.0872	-0.2447 \pm 0.1077	● -0.1723 \pm 0.0847
Cal_housing	0.1080 \pm 0.0802	0.3183 \pm 0.1036	● 0.3762 \pm 0.1359
Cpu_act	0.3686 \pm 0.1095	0.4731 \pm 0.1096	● 0.6866 \pm 0.1941
Cpu_small	0.4483 \pm 0.1435	0.4726 \pm 0.1821	● 0.8685 \pm 0.1979
Elevators	0.4871 \pm 0.0334	● 0.4960 \pm 0.0246	0.4843 \pm 0.0256
Folds5x2_pp	0.8616 \pm 0.0217	● 0.8811 \pm 0.0700	0.8459 \pm 0.0101
Kin8nm	● 0.3678 \pm 0.0318	0.3406 \pm 0.0107	0.3117 \pm 0.0429
Parkinsons	0.2642 \pm 0.0530	0.2628 \pm 0.0748	● 0.4887 \pm 0.0912
Pollen	-2.9128 \pm 0.1539	-2.9131 \pm 0.1504	● -0.2070 \pm 0.1513
Puma8NH	0.1533 \pm 0.0551	0.1538 \pm 0.0480	● 0.3411 \pm 0.0349
Space_ga	-1.3920 \pm 0.2160	-1.3540 \pm 0.2347	● 0.4155 \pm 0.2040
Wind	0.0953 \pm 0.0695	0.2804 \pm 0.0602	● 0.3989 \pm 0.0551
Wine_quality	-0.3879 \pm 0.1020	-0.2922 \pm 0.1155	● 0.1419 \pm 0.1555

The best results obtained on each dataset are denoted with bullets. COREG-ESP and SAFER-ESP denote the COREG and SAFER algorithms constructed after applying embedding space mapping (ESP), respectively

4.3.3 Pseudo-label smearing

In the pseudo-label calculation process, S2RMS uses the smearing technique to smooth out the predictions of the regressors. To evaluate the effectiveness of the proposed method, we compare the RMSE values obtained with and without pseudo-label smearing. Lines 16 to 22 of Algorithm 1 represent the process of label smearing. The core of this process is the calculation of stability, which is performed using (4) to filter the stable instances. Once the stable instances have been identified, the subsequent confidence assessment operation is performed. In the baseline experiment where label

smearing is not employed, the step in line 17 of Algorithm 1 is ignored. Line 19 and subsequent operations are conducted directly without any form of judgement. The objective of this experiment is to assess the impact of stable instance filtering on the performance of the algorithm in the context of label smearing.

Figure 4 shows the results of the experimental comparisons.

Under the same configuration, the pseudo-label smearing technique can effectively improve the performance of the algorithm. At a label rate of 0.5%, the average RMSE is notably greater when employing the smearing technique than

Table 6 The values of R^2 produced by COREG and SAFER at a label rate of 2.5%

Dataset	COREG	COREG-ESP	SAFER	SAFER-ESP
	COREG		SAFER	
Abalone	0.2661 \pm 0.0538	● 0.2738 \pm 0.0458	0.1011 \pm 0.0147	● 0.1132 \pm 0.0734
Bank32nh	● -0.2224 \pm 0.0905	-0.2234 \pm 0.0872	● -0.2042 \pm 0.0903	-0.2447 \pm 0.1077
Cal_housing	0.0851 \pm 0.1737	● 0.1080 \pm 0.0802	0.2739 \pm 0.1503	● 0.3183 \pm 0.1036
Cpu_act	0.2552 \pm 0.1778	● 0.3686 \pm 0.1095	0.3381 \pm 0.1898	● 0.4731 \pm 0.1096
Cpu_small	0.2685 \pm 0.2345	● 0.4483 \pm 0.1435	0.3822 \pm 0.1977	● 0.4726 \pm 0.1821
Elevators	0.3500 \pm 0.0225	● 0.4871 \pm 0.0334	0.4475 \pm 0.0246	● 0.4960 \pm 0.0246
Folds5x2_pp	● 0.7135 \pm 0.0159	0.8616 \pm 0.0217	0.8155 \pm 0.0084	● 0.8811 \pm 0.0700
Kin8nm	0.2288 \pm 0.0239	● 0.3678 \pm 0.0318	0.2424 \pm 0.0347	● 0.3406 \pm 0.0107
Parkinsons	0.1019 \pm 0.0483	● 0.2642 \pm 0.0530	0.2210 \pm 0.0900	● 0.2628 \pm 0.0748
Pollen	● -2.9091 \pm 0.1492	-2.9128 \pm 0.1539	● -2.9090 \pm 0.1511	-2.9131 \pm 0.1504
Puma8NH	0.1292 \pm 0.0485	● 0.1533 \pm 0.0551	0.1386 \pm 0.0405	● 0.1538 \pm 0.0480
Space_ga	● -1.3672 \pm 0.2581	-1.3920 \pm 0.2160	-1.4029 \pm 0.2048	● -1.3540 \pm 0.2347
Wind	0.0152 \pm 0.2447	● 0.0953 \pm 0.0695	-0.1860 \pm 0.0536	● 0.2804 \pm 0.0602
Wine_quality	-0.1888 \pm 0.2463	● -0.3879 \pm 0.1020	-0.3210 \pm 0.1593	● -0.2922 \pm 0.1155

The best results obtained by each dataset are denoted with bullets. COREG-ESP and SAFER-ESP denote the COREG and SAFER algorithms constructed after applying embedding space mapping (ESP), respectively

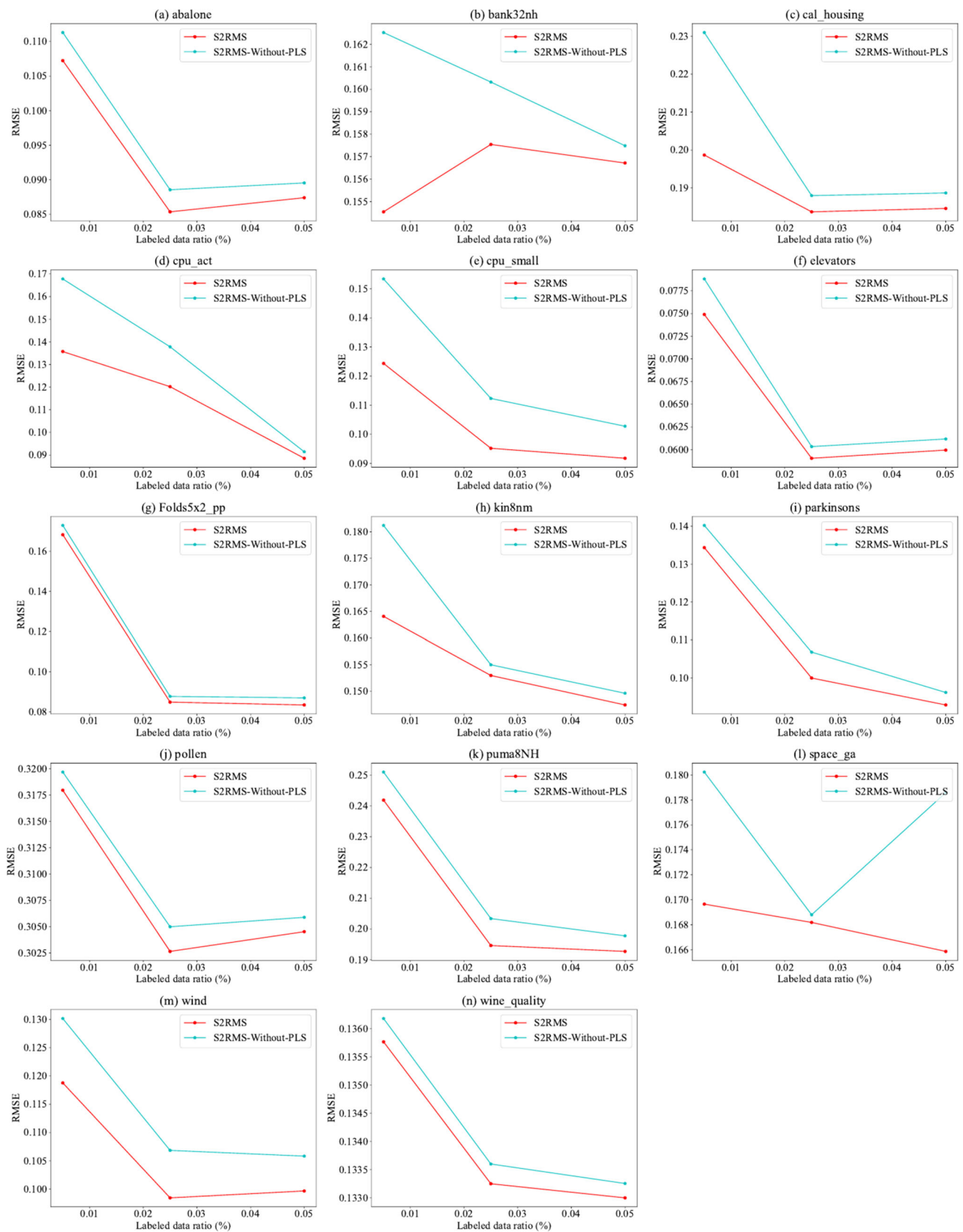


Fig. 4 Average RMSE values produced over thirty tests on fourteen datasets. S2RMS-Without-PLS denotes the S2RMS algorithm without the pseudo-label smearing (PLS) technique

that attained when not using it. With the help of smearing, the RMSE decreases smoothly on most of the datasets as the label rate increases. Especially on the Space_ga dataset, the RMSE is more volatile without the smearing technique. This indicates that the smearing technique can effectively leverage stable instances to optimize the regression performance of the algorithm.

4.3.4 The influence of validation set

Considering the relatively small size of the labeled dataset, the validation set C'_i is also expected to be small. To investigate potential biases in the confidence calculation executed during the pseudo-label smearing process caused by the size of the validation set, we conduct comparative experiments with different C'_i sizes.

The experiment of this subsection aims to study the impact of the validation set size. Therefore, we add some additional labeled data to increase the validation set. In this way, the size can be 30 or 50. We also consider the extreme case where the size of the validation set is 0. In this case, it is impossible to perform a confidence calculation, and the algorithm degenerates into a self-training mode in its ordinary form. Consequently, lines 14, 21, 22, and 24 of Algorithm 1 are disabled, and the strategy that replaces them is to randomly select R instances from Π_i .

Table 7 shows the experimental results obtained on the validation set C'_i under different size settings. The performance of the algorithm gradually improves as the size of the validation set increases. A larger validation set leads to more accurate calculation results for (5). Inevitably, the labeled dataset size is positively correlated with the size of the vali-

dation set. With a small labeled dataset, a small validation set can cause bias in the confidence computation. Nevertheless, the algorithm also performs better with a small validation set than it does without it.

4.3.5 Parameter sensitivity

The parameters τ and μ are key parameters that affect the performance of the S2RMS algorithm. One of the parameters, τ , is the attenuation coefficient, which is applied to the process of creating a pool of unlabeled instances. The value of τ reflects the size of the instance pool. The number of instances in the pool directly affects the performance of the co-regressors in subsequent iterations. The use of too few unlabeled instances is meaningless for the learner. An instance pool that is excessively large makes the learning process less efficient. To obtain a suitable value for τ , we observe the performance achieved by the algorithm with different parameter settings. When testing parameter τ , it is important to note that we only focus on the process related to parameter τ to clearly observe the experimental results. The co-training process removes the label smearing strategy and eliminates the experimental bias caused by the stable instances controlled by parameter μ .

Table 8 shows the experimental results produced by S2RMS with different τ values. A higher τ value implies that the conditions for the similarity threshold are relaxed to a greater extent. When τ is set to 0.008, this lower value prevents the selection of more unlabeled instances that do not satisfy the imposed similarity requirement. This results in a final unlabeled data pool that is too small or even empty after just a few iterations due to the decaying value of τ . The

Table 7 The values of R^2 obtained on the validation set C'_i under sizes of 0, 30 and 50 at a label rate of 2.5%

Dataset	$ C'_i = 0$	$ C'_i = 30$	$ C'_i = 50$
Abalone	0.3217 \pm 0.0046	0.3617 \pm 0.0510	● 0.3880 \pm 0.0264
Bank32nh	-0.6241 \pm 0.1990	-0.2864 \pm 0.2465	● -0.2122 \pm 0.1837
Cal_housing	0.3879 \pm 0.0074	0.3722 \pm 0.0537	● 0.3973 \pm 0.0278
Cpu_act	0.6230 \pm 0.0009	0.6673 \pm 0.0115	● 0.6926 \pm 0.0364
Cpu_small	0.6854 \pm 0.1011	0.7104 \pm 0.0238	● 0.7426 \pm 0.0120
Elevators	0.4582 \pm 0.0006	0.4824 \pm 0.0101	● 0.4874 \pm 0.0152
Folds5x2_pp	0.8376 \pm 0.0041	0.8318 \pm 0.0023	● 0.8406 \pm 0.0070
Kin8nm	0.1743 \pm 0.0330	● 0.2872 \pm 0.0394	0.2685 \pm 0.0128
Parkinsons	0.3172 \pm 0.1101	0.4298 \pm 0.0748	● 0.4354 \pm 0.1142
Pollen	-0.1654 \pm 0.0779	● -0.1137 \pm 0.0233	-0.1289 \pm 0.0407
Puma8NH	0.1669 \pm 0.1201	● 0.3237 \pm 0.1863	0.2115 \pm 0.1337
Space_ga	● 0.4721 \pm 0.1367	0.4222 \pm 0.0131	0.4458 \pm 0.0271
Wind	0.3763 \pm 0.0361	0.4154 \pm 0.0707	● 0.4289 \pm 0.0890
Wine_quality	0.1647 \pm 0.1103	0.2153 \pm 0.1291	● 0.2283 \pm 0.1321

The best results obtained on each dataset version are denoted with bullets. $|C'_i|$ denotes the cardinality of the validation set C'_i

Table 8 The values of R^2 produced under various settings of parameter τ at a label rate of 2.5%

Dataset	$\tau = 0.008$	$\tau = 0.011$	$\tau = 0.014$	$\tau = 0.017$	$\tau = 0.020$
Abalone	0.3298 \pm 0.1258	0.3049 \pm 0.1094	● 0.3378 \pm 0.0827	0.3339 \pm 0.1371	0.3245 \pm 0.1104
Bank32nh	−0.1567 \pm 0.1844	● −0.1370 \pm 0.1652	−0.1380 \pm 0.1285	−0.1421 \pm 0.1771	−0.1500 \pm 0.1801
Cal_housing	0.3176 \pm 0.0120	0.3721 \pm 0.0240	0.3254 \pm 0.0553	● 0.3975 \pm 0.0683	0.3345 \pm 0.0334
Cpu_act	0.6507 \pm 0.0889	0.6122 \pm 0.0548	0.6487 \pm 0.0856	● 0.6859 \pm 0.0400	0.6666 \pm 0.0798
Cpu_small	0.6625 \pm 0.0832	0.6808 \pm 0.0639	0.6575 \pm 0.1057	● 0.7198 \pm 0.0752	0.7108 \pm 0.0560
Elevators	0.4594 \pm 0.0266	0.4889 \pm 0.0222	0.4662 \pm 0.0315	● 0.4908 \pm 0.0260	0.4532 \pm 0.0301
Folds5x2_pp	0.8373 \pm 0.0141	0.8322 \pm 0.0064	0.8435 \pm 0.0189	● 0.8436 \pm 0.0127	0.8401 \pm 0.0191
Kin8nm	0.2237 \pm 0.0283	0.2580 \pm 0.0094	0.2510 \pm 0.0083	0.2459 \pm 0.0242	● 0.2680 \pm 0.0101
Parkinsons	● 0.3645 \pm 0.0825	0.3576 \pm 0.0938	0.3366 \pm 0.0992	0.3450 \pm 0.0960	0.3332 \pm 0.0887
Pollen	−0.1418 \pm 0.0401	−0.1454 \pm 0.0572	−0.1388 \pm 0.0317	● −0.1193 \pm 0.0636	−0.1297 \pm 0.0443
Puma8NH	0.2916 \pm 0.1027	0.3046 \pm 0.1188	● 0.2983 \pm 0.1418	0.2698 \pm 0.0975	0.2721 \pm 0.1770
Space_ga	0.4022 \pm 0.0843	0.3789 \pm 0.0460	0.4518 \pm 0.0489	0.4153 \pm 0.0380	● 0.4522 \pm 0.0599
Wind	0.4121 \pm 0.0102	0.4120 \pm 0.0110	0.4117 \pm 0.0027	0.4160 \pm 0.0229	● 0.4183 \pm 0.0149
Wine_quality	0.2000 \pm 0.0204	0.1945 \pm 0.0334	0.1830 \pm 0.0205	● 0.2070 \pm 0.0285	0.1934 \pm 0.0289

The best results obtained on each dataset are denoted with bullets

performance of the subsequent regressors remains relatively low due to the insufficient availability of unlabeled data.

As the value of τ increases to 0.017, more unlabeled instances are able to satisfy the larger threshold. Additionally, sufficient training rounds and data are available for improving the performance of the regressor until the threshold τ decays to the point where the instance pool can no longer be constructed. As the value of τ increases again, some unlabeled data that do not satisfy the imposed conditions are selected, affecting the performance of the algorithm. The experimental results indicate that a parameter τ setting of approximately 0.017 is appropriate.

The parameter μ denotes the stability coefficient, which is used to screen for stable instances that satisfy the preset conditions. The stability of a instance is indicated by the variance of the values predicted by multiple regressors. For each instance in the pool of unlabeled instances, when its stability value is lower than the parameter μ , it is regarded as a stable instance and marked with a pseudo-label. In summary, parameter μ controls the prediction variance among the co-learners by setting a stability threshold to ensure the robustness of the algorithm. Similar to when testing parameter τ , we focus solely on the process related to parameter μ and exclude the similarity-based instance selection process that is controlled by parameter τ . The pool of unlabeled

Table 9 The values of R^2 produced with various settings of μ at a label rate of 2.5%

Dataset	$\mu = 0.01$	$\mu = 0.02$	$\mu = 0.03$	$\mu = 0.04$
Abalone	0.3316 \pm 0.1323	● 0.3594 \pm 0.1372	0.3378 \pm 0.1140	0.3129 \pm 0.1322
Bank32nh	−0.1580 \pm 0.1305	● −0.1434 \pm 0.1250	−0.1499 \pm 0.1228	−0.1522 \pm 0.1893
Cal_housing	0.3859 \pm 0.0032	● 0.3947 \pm 0.0083	0.3505 \pm 0.0290	0.3434 \pm 0.0360
Cpu_act	0.6418 \pm 0.0385	0.6153 \pm 0.0107	● 0.6912 \pm 0.0435	0.6286 \pm 0.0209
Cpu_small	0.6854 \pm 0.0162	● 0.7662 \pm 0.0392	0.6338 \pm 0.0167	0.6578 \pm 0.0085
Elevators	0.4682 \pm 0.0308	● 0.4967 \pm 0.0111	0.4645 \pm 0.0376	0.4466 \pm 0.0094
Folds5x2_pp	0.8373 \pm 0.0075	0.8357 \pm 0.0050	● 0.8444 \pm 0.0012	0.8411 \pm 0.0111
Kin8nm	0.2396 \pm 0.0119	● 0.2796 \pm 0.0231	0.2716 \pm 0.0017	0.2221 \pm 0.0056
Parkinsons	0.3645 \pm 0.1136	0.3274 \pm 0.1214	0.3570 \pm 0.1133	● 0.3805 \pm 0.0933
Pollen	−0.1424 \pm 0.0059	● −0.1269 \pm 0.0233	−0.1454 \pm 0.0604	−0.1373 \pm 0.0030
Puma8NH	● 0.3033 \pm 0.1594	0.2970 \pm 0.1821	0.2826 \pm 0.1301	0.2729 \pm 0.1929
Space_ga	0.4520 \pm 0.0166	● 0.4484 \pm 0.0165	0.4048 \pm 0.0124	0.3926 \pm 0.0333
Wind	● 0.4117 \pm 0.0209	0.4050 \pm 0.0028	0.3971 \pm 0.0033	0.4067 \pm 0.0099
Wine_quality	0.2114 \pm 0.0204	0.2163 \pm 0.0460	● 0.2180 \pm 0.0010	0.1955 \pm 0.0325

The best results obtained on each dataset are denoted with bullets

instances is then constructed by randomly selecting instances with a fixed size.

Table 9 shows the experimental results yielded by S2RMS under different μ values. The parameter μ indicates the degree of the predictive differences among multiple regressors. When μ is set to 0.01, only unlabeled instances with prediction variances of 0.009 or less are considered stable instances. Due to variations in the regression model settings, producing identical outputs for unlabeled data is essentially impossible. Lower parameter settings may not provide sufficient unlabeled data for sufficiently augmenting the regressors. Conversely, higher settings, such as 0.04, may result in the selection of the vast majority of the instances in the instance pool, which would interfere with the performance of the regressor due to the presence of widely varying labels. The experiments show that setting the unlabeled instance rate to 0.02 provides enough, but not all, instances to satisfy the stability constraint. These unlabeled instances can effectively enhance the training set of the corresponding regressor by labeling them with pseudo-labels.

5 Conclusions and further works

This paper proposes a semi-supervised regression approach via embedding space mapping and pseudo-label smearing. The method involves constructing a deep metric network based on pairwise data, mapping the raw training data to the metric embedding space through this network, and then training co-regressors on the transformed mapping data. The parameters of the regression model are iteratively updated using the label smearing technique. Tests conducted on 14 benchmark datasets show that the proposed algorithm outperforms 3 other excellent semi-supervised regression methods. An ablation study indicates that both the metric approach and the label smearing technique are effective for improving the performance of the algorithm.

Our future research will focus on designing more suitable mapping strategies and developing new confidence calculation methods. It is important to consider both the original spatial distribution of the given training data and the feature representations during the metric mapping process. Additionally, we will investigate improved confidence evaluation methods. Inspired by the evaluation metrics used in supervised learning research, novel reliability estimation strategies will be applied to increase the accuracy of confidence calculation [53]. Since the current study focuses on single-target regression, we will explore the potential to refine the proposed methods for applications involving multi-target regression tasks in our future research.

Acknowledgements This study is partially supported by the Nanchong Municipal Government-Universities Scientific Cooperation Project (Nos. 23XNSYSX0062, 23XNSYSX0013, and SXHZ051).

Author Contributions **Liyan Liu:** Methodology, Validation, Supervision, Writing - Review & Editing. **Jin Zhang:** Conceptualization, Methodology, Software, Visualization, Writing - Original Draft. **Kun Qian:** Supervision, Writing - Review & Editing. **Fan Min:** Methodology, Supervision, Writing - Review & Editing.

Declarations

Conflicts of interest The authors have no relevant financial or non-financial interests to disclose.

Ethical and informed consent for data used The data used in this study were obtained from publicly available datasets as described in Table 2. The original data collection was approved by the appropriate institutional ethical review board, and informed consent was obtained from all participants.

References

1. Chapelle O, Schölkopf B, Zien A (2009) Semi-supervised learning. *IEEE Trans Neural Netw* 20(3):542–542. <https://doi.org/10.1109/TNN.2009.2015974>
2. Goldman S, Zhou Y (2000) Enhancing supervised learning with unlabeled data. In: *ICML*, pp 327–334
3. Blum A, Mitchell T (1998) Combining labeled and unlabeled data with co-training. In: *COLT*, pp 92–100. <https://doi.org/10.1145/279943.279962>
4. Gong M, Zhou H, Qin AK, Liu W, Zhao Z (2023) Self-paced co-training of graph neural networks for semi-supervised node classification. *IEEE Trans Neural Netw Learn Syst* 34(11):9234–9247. <https://doi.org/10.1109/TNNLS.2022.3157688>
5. Wei C, Sohn K, Mellina C, Yuille A, Yang F (2021) Crest: A class-rebalancing self-training framework for imbalanced semi-supervised learning. In: *CVPR*, pp 10852–10861. <https://doi.org/10.1109/CVPR46437.2021.01071>
6. Li J, Zhu Q (2020) A boosting self-training framework based on instance generation with natural neighbors for k nearest neighbor. *Appl Intell* 50(11):3535–3553. <https://doi.org/10.1007/s10489-020-01732-1>
7. Miyato T, Maeda S-I, Koyama M, Ishii S (2019) Virtual adversarial training: A regularization method for supervised and semi-supervised learning. *IEEE Trans Pattern Anal Mach Intell* 41(8):1979–1993. <https://doi.org/10.1109/TPAMI.2018.2858821>
8. Verma V, Kawaguchi K, Lamb A, Kannala J, Solin A, Bengio Y, Lopez-Paz D (2022) Interpolation consistency training for semi-supervised learning. *Neural Netw* 145:90–106. <https://doi.org/10.1016/j.neunet.2021.10.008>
9. Jiang Y, Li X, Che Y, He Y, Xu Q, Yang Z, Cao X, Huang Q (2023) Maxmatch: Semi-supervised learning with worst-case consistency. *IEEE Trans Pattern Anal Mach Intell* 45(5):5970–5987. <https://doi.org/10.1109/TPAMI.2022.3208419>
10. Duan Y, Zhao Z, Qi L, Wang L, Zhou L, Shi Y, Gao Y (2022) Mutexmatch: Semi-supervised learning with mutex-based consistency regularization. *IEEE Trans Neural Netw Learn Syst* 1–15. <https://doi.org/10.1109/TNNLS.2022.3228380>
11. Kostopoulos G, Karlos S, Kotsiantis S, Ragos O, Tiwari S, Trivedi M, Kohle ML (2018) Semi-supervised regression: A recent

- review. *J Intell Fuzzy Syst* 35(2):1483–1500. <https://doi.org/10.3233/JIFS-169689>
12. Jing W, Lu L, Liu Q (2023) Graph regularized discriminative non-negative tucker decomposition for tensor data representation. *Appl Intell* 53(20):23864–23882. <https://doi.org/10.1007/s10489-023-04738-7>
13. Dong H, Yang L, Wang X (2021) Robust semi-supervised support vector machines with laplace kernel-induced correntropy loss functions. *Appl Intell* 51(2):819–833. <https://doi.org/10.1007/s10489-020-01865-3>
14. Balcan M-F, Blum A (2010) A discriminative model for semi-supervised learning. *J ACM* 57(3):1–46. <https://doi.org/10.1145/1706591.1706599>
15. Yang T, Priebe CE (2011) The effect of model misspecification on semi-supervised classification. *IEEE Trans Pattern Anal Mach Intell* 33(10):2093–2103. <https://doi.org/10.1109/TPAMI.2011.45>
16. Li Y-F, Tsang IW, Kwok JT, Zhou Z-H (2013) Convex and scalable weakly labeled SVMs. *J Mach Learn Res* 14(1):2151–2188
17. Li Y-F, Zha H-W, Zhou Z-H (2017) Learning safe prediction for semi-supervised regression. In: AAAI, pp 2217–2223
18. Ma F, Meng D, Xie Q, Li Z, Dong X (2017) Self-paced co-training. In: ICML, pp 2275–2284
19. Breiman L (2001) Random forests. *Mach Learn* 45:5–32. <https://doi.org/10.1023/A:1010933404324>
20. Song K, Han J, Cheng G, Lu J, Nie F (2022) Adaptive neighborhood metric learning. *IEEE Trans Pattern Anal Mach Intell* 44(9):4591–4604. <https://doi.org/10.1109/TPAMI.2021.3073587>
21. Milbich T, Roth K, Bratoli B, Ommer B (2022) Sharing matters for generalization in deep metric learning. *IEEE Trans Pattern Anal Mach Intell* 44(1):416–427. <https://doi.org/10.1109/TPAMI.2020.3009620>
22. Bromley J, Guyon I, LeCun Y, Säckinger E, Shah R (1993) Signature verification using a “siamese” time delay neural network. *Int J Pattern Recognit Artif Intell* 7(4):669–688. <https://doi.org/10.1142/S0218001493000339>
23. Hoffer E, Ailon N (2015) Deep metric learning using triplet network. In: SIMBAD, pp 84–92. https://doi.org/10.1007/978-3-319-24261-3_7
24. Ge W, Huang W, Dong D, Scott MR (2018) Deep metric learning with hierarchical triplet loss. In: ECCV, pp 272–288. https://doi.org/10.1007/978-3-030-01231-1_17
25. Deng J, Guo J, Xue N, Zafeiriou S (2019) Arcface: Additive angular margin loss for deep face recognition. In: CVPR, pp 4685–4694. <https://doi.org/10.1109/CVPR.2019.00482>
26. Wang X, Hua Y, Kodirov E, Robertson NM (2022) Ranked list loss for deep metric learning. *IEEE Trans Pattern Anal Mach Intell* 44(9):5414–5429. <https://doi.org/10.1109/TPAMI.2021.3068449>
27. Yao X, She D, Zhang H, Yang J, Cheng M-M, Wang L (2021) Adaptive deep metric learning for affective image retrieval and classification. *IEEE Trans Multimed* 23:1640–1653. <https://doi.org/10.1109/TMM.2020.3001527>
28. Duan C, Feng Y, Zhou M, Xiong X, Wang Y, Qiang B, Jia W (2023) Multilevel similarity-aware deep metric learning for fine-grained image retrieval. *IEEE Trans Ind Inform* 19(8):9173–9182. <https://doi.org/10.1109/TII.2022.3227721>
29. Chang X, Ma Z, Wei X, Hong X, Gong Y (2020) Transductive semi-supervised metric learning for person re-identification. *Pattern Recognit* 108:107569. <https://doi.org/10.1016/j.patcog.2020.107569>
30. Uzun B, Cevikalp H, Saribas H (2023) Deep discriminative feature models (ddfms) for set based face recognition and distance metric learning. *IEEE Trans Pattern Anal Mach Intell* 45(5):5594–5608. <https://doi.org/10.1109/TPAMI.2022.3205939>
31. Zhou Z-H, Li M (2005) Tri-training: exploiting unlabeled data using three classifiers. *IEEE Trans Knowl Data Eng* 17(11):1529–1541. <https://doi.org/10.1109/TKDE.2005.186>
32. Gong Y, Wu Q (2023) SIVLC: improving the performance of co-training by sufficient-irrelevant views and label consistency. *Appl Intell* 53:20710–20729. <https://doi.org/10.1007/s10489-023-04611-7>
33. Qiao S, Shen W, Zhang Z, Wang B, Alan Y (2018) Deep co-training for semi-supervised image recognition. In: ECCV, pp 142–159. https://doi.org/10.1007/978-3-030-01267-0_9
34. Chen D-D, Wang W, Gao W, Zhou Z-H (2018) Tri-net for semi-supervised deep learning. In: IJCAI, pp 2014–2020. <https://doi.org/10.24963/ijcai.2018/278>
35. Zhou Z-H, Li M (2007) Semisupervised regression with cotraining-style algorithms. *IEEE Trans Knowl Data Eng* 19(11):1479–1493. <https://doi.org/10.1109/TKDE.2007.190644>
36. Cover T, Hart PE (1967) Nearest neighbor pattern classification. *IEEE Trans Inf Theory* 13(1):21–27. <https://doi.org/10.1109/TIT.1967.1053964>
37. Hady MFA, Schwenker F, Palm G (2009) Semi-supervised learning for regression with co-training by committee. In: ICANN, pp 121–130. https://doi.org/10.1007/978-3-642-04274-4_13
38. Moody J, Darken CJ (1989) Fast learning in networks of locally-tuned processing units. *Neural Comput* 1(2):281–294. <https://doi.org/10.1162/neco.1989.1.2.281>
39. Min F, Li Y, Liu L (2022) Self-paced safe co-training for regression. In: PAKDD, pp 71–82. https://doi.org/10.1007/978-3-031-05936-0_6
40. Lu S, Shi X, Li M, Jiao J, Feng L, Wang G (2021) Semi-supervised random forest regression model based on co-training and grouping with information entropy for evaluation of depression symptoms severity. *Math Biosci Eng* 18(4):4586–4602. <https://doi.org/10.3934/mbe.2021233>
41. Li H, Wang S, Liu B, Fang M, Cao R, He B, Liu S, Hu C, Dong D, Wang X, Wang H, Tian J (2023) A multi-view co-training network for semi-supervised medical image-based prognostic prediction. *Neural Netw* 164:455–463. <https://doi.org/10.1016/j.neunet.2023.04.030>
42. Shen Z, Cao P, Yang H, Liu X, Yang J, Zaiane OR (2023) Co-training with high-confidence pseudo labels for semi-supervised medical image segmentation. In: IJCAI, pp 4199–4207. <https://doi.org/10.24963/ijcai.2023/467>
43. Xie H, Fu C, Zheng X, Zheng Y, Sham C-W, Wang X (2023) Adversarial co-training for semantic segmentation over medical images. *Comput Biol Med* 157:106736. <https://doi.org/10.1016/j.compbiomed.2023.106736>
44. Li Q, Chen Y, He X, Huang L (2024) Co-training transformer for remote sensing image classification, segmentation, and detection. *IEEE Trans Geosci Remote Sens* 62:1–18. <https://doi.org/10.1109/TGRS.2024.3354783>
45. Gong Y, Wu Q, Zhou M, Wen J (2023) Self-paced multi-label co-training. *Inf Sci* 622:269–281. <https://doi.org/10.1016/j.ins.2022.11.153>
46. Liu Z, Ma Q, Ma P, Wang L (2023) Temporal-frequency co-training for time series semi-supervised learning. In: AAAI, pp 8923–8931. <https://doi.org/10.1609/aaai.v37i7.26072>
47. Jiang Z, Zhao L, Zhan Y (2023) A boosted co-training method for class-imbalanced learning. *Expert Syst* 40(9):13377. <https://doi.org/10.1111/exsy.13377>
48. Lu J, Gong Y (2021) A co-training method based on entropy and multi-criteria. *Appl Intell* 51(6):3212–3225. <https://doi.org/10.1007/s10489-020-02014-6>
49. Li Y, Xiong H, Wang Q, Kong L, Liu H, Li H, Bian J, Wang S, Chen G, Dou D, Yin D (2023) COLTR: Semi-supervised learning to rank with co-training and over-parameterization for web search. *IEEE Trans Knowl Data Eng* 35(12):12542–12555. <https://doi.org/10.1109/TKDE.2023.3270750>
50. Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, Killeen T, Lin Z, Gimelshein N, Antiga L, Desmaison A, Köpf A,

- Yang E, DeVito Z, Raison M, Tejani A, Chilamkurthy S, Steiner B, Fang L, Bai J, Chintala S (2019) Pytorch: An imperative style, high-performance deep learning library. In: NeurIPS, pp 8026–8037
51. Khalid M, Baber J, Kasi MK, Bakhtyar M, Devi V, Sheikh N (2020) Empirical evaluation of activation functions in deep convolution neural network for facial expression recognition. In: TSP, pp 204–207. <https://doi.org/10.1109/TSP49548.2020.9163446>
 52. Zell A, Sumbul G, Demir B (2022) Deep metric learning-based semi-supervised regression with alternate learning. In: ICIP, pp 2411–2415. <https://doi.org/10.1109/ICIP46576.2022.9897939>
 53. Bosnić Z, Kononenko I (2008) Comparison of approaches for estimating reliability of individual regression predictions. *Data Knowl Eng* 67(3):504–516. <https://doi.org/10.1016/j.datak.2008.08.001>

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.