

A Survey of Active Learning Algorithms for Supervised Remote Sensing Image Classification

Devis Tuia, *Member, IEEE*, Michele Volpi, *Student Member, IEEE*, Loris Copa, Mikhail Kanevski, and Jordi Muñoz-Marí

Overview

- Defining an efficient **training set** → Fundamental phase for classification
- Active learning aims at building efficient training sets by **iteratively improving** the model performance through sampling.
- A user-defined heuristic ranks the unlabeled pixels according to a function of the **uncertainty**
- This paper reviews and tests the main families of active learning algorithms:
 1. committee,
 2. large margin,
 3. posterior probability-based

1. CONCEPTS AND DEFINITIONS

Algorithm 1: General active learning algorithm

Inputs

- Initial training set $X^\epsilon = \{\mathbf{x}_i, y_i\}_{i=1}^l$ ($X \in \mathcal{X}$, $\epsilon = 1$).
- Pool of candidates $U^\epsilon = \{\mathbf{x}_i\}_{i=l+1}^{l+u}$ ($U \in \mathcal{X}$, $\epsilon = 1$).
- Number of pixels q to add at each iteration (defining the batch of selected pixels S).

- 1: **repeat**
 - 2: Train a model with current training set X^ϵ .
 - 3: **for** each candidate in U^ϵ **do**
 - 4: Evaluate a user-defined *heuristic*
 - 5: **end for**
 - 6: Rank the candidates in U^ϵ according to the score of the heuristic
 - 7: Select the q most interesting pixels. $S^\epsilon = \{\mathbf{x}_k\}_{k=1}^q$
 - 8: The user assigns a label to the selected pixels.
 $S^\epsilon = \{\mathbf{x}_k, y_k\}_{k=1}^q$
 - 9: Add the batch to the training set $X^{\epsilon+1} = X^\epsilon \cup S^\epsilon$.
 - 10: Remove the batch from the pool of candidates
 $U^{\epsilon+1} = U^\epsilon \setminus S^\epsilon$
 - 11: $\epsilon = \epsilon + 1$
 - 12: **until** a stopping criterion is met.
-

2. COMMITTEE-BASED ACTIVE LEARNING

The first family of active learning methods quantifies the uncertainty of a pixel by considering a committee of learners.

1. *Normalized Entropy Query-by-Bagging*

K training sets built on a draw with replacement of the original data are defined. These draws account for a part of the available labeled pixels only. Then, each set is used to train a classifier and to predict the labels of the candidates.

$$\hat{\mathbf{x}}^{n\text{EQB}} = \arg \max_{\mathbf{x}_i \in U} \left\{ \frac{H^{\text{BAG}}(\mathbf{x}_i)}{\log(N_i)} \right\} \quad (1)$$

where

$$\begin{aligned} H^{\text{BAG}}(\mathbf{x}_i) \\ = - \sum_{\omega=1}^{N_i} p^{\text{BAG}}(y_i^* = \omega \mid \mathbf{x}_i) \log [p^{\text{BAG}}(y_i^* = \omega \mid \mathbf{x}_i)] \end{aligned} \quad (2)$$

where

$$p^{\text{BAG}}(y_i^* = \omega \mid \mathbf{x}_i) = \frac{\sum_{m=1}^k \delta(y_{i,m}^*, \omega)}{\sum_{m=1}^k \sum_{j=1}^{N_i} \delta(y_{i,m}^*, \omega_j)}.$$

2. COMMITTEE-BASED ACTIVE LEARNING

2. Adaptive Maximum Disagreement (AMD)

When confronted to high dimensional data, it may be relevant to construct the committee by splitting the feature space into a number of subsets, or *views*.

$$\hat{\mathbf{x}}^{\text{AMD}} = \arg \max_{\mathbf{x}_i \in U} H^{\text{MV}}(\mathbf{x}_i) \quad (3)$$

where the multiview entropy H^{MV} is assessed over the predictions of classifiers using a specific view v :

$$H^{\text{MV}}(\mathbf{x}_i) = - \sum_{\omega=1}^{N_i} p^{\text{MV}}(y_{i,v}^* = \omega \mid \mathbf{x}_i^v) \times \log [p^{\text{MV}}(y_{i,v}^* = \omega \mid \mathbf{x}_i^v)] \quad (4)$$

where

$$p^{\text{MV}}(y_i^* = \omega \mid \mathbf{x}_i^v) = \frac{\sum_{v=1}^V W^{\epsilon-1}(v, \omega) \delta(y_{i,v}^*, \omega)}{\sum_{v=1}^V \sum_{j=1}^{N_i} W^{\epsilon-1}(v, \omega)}$$

3. LARGE-MARGIN-BASED ACTIVE LEARNING

The second family of methods is specific to margin-based classifiers (SVM)

The distance of a sample \mathbf{x}_i from the SVM hyperplane is given by

$$f(\mathbf{x}_i) = \sum_{j=1}^n \alpha_j y_j K(\mathbf{x}_j, \mathbf{x}_i) + b$$

1. *Margin Sampling (MS)*

$$\hat{\mathbf{x}}^{\text{MS}} = \arg \min_{\mathbf{x}_i \in U} \left\{ \min_{\omega} |f(\mathbf{x}_i, \omega)| \right\}$$

2. *Multiclass Level Uncertainty (MCLU)*

$$\hat{\mathbf{x}}^{\text{MCLU}} = \arg \min_{\mathbf{x}_i \in U} \{ f(\mathbf{x}_i)^{\text{MC}} \} \quad (8)$$

where

$$f(\mathbf{x}_i)^{\text{MC}} = \max_{\omega \in N} f(\mathbf{x}_i, \omega) - \max_{\omega \in N \setminus \omega^+} f(\mathbf{x}_i, \omega) \quad (9)$$

3. LARGE-MARGIN-BASED ACTIVE LEARNING

3. *Significance Space Construction (SSC)*

The support vector coefficients are used to convert the multiclass classification problem into a binary support vector detection problem. This second classifier predicts which pixels are likely to become support vectors:

$$\hat{\mathbf{x}}^{\text{SSC}} = \arg_{\mathbf{x}_i \in U} f^{\text{SSC}}(\mathbf{x}_i) > 0. \quad (10)$$

Once the candidates more likely to become support vectors have been highlighted, a random selection among them is done.

3. LARGE-MARGIN-BASED ACTIVE LEARNING

4. *On the Need for a Diversity Criterion*

- The heuristic, called “most ambiguous and orthogonal” (MAO) is iterative: starting from the samples selected by MS, , this heuristic iteratively chooses the samples minimizing the highest values between the candidates list and the samples already included in the batch .

$$\hat{\mathbf{x}}^{\text{MAO}} = \arg \min_{\mathbf{x}_i \in U^{\text{MS}}} \left\{ \max_{\mathbf{x}_j \in S} K(\mathbf{x}_i, \mathbf{x}_j) \right\} .$$

- the MAO criterion is combined with the MCLU uncertainty estimation in the “multiclass level uncertainty—angle based diversity” (MCLU-ABD) heuristic.

$$\hat{\mathbf{x}}^{\text{MCLU-ABD}} = \arg \min_{\mathbf{x}_i \in U^{\text{MCLU}}} \left\{ \lambda f(\mathbf{x}_i)^{\text{MC}} + (1 - \lambda) \max_{\mathbf{x}_j \in S} \frac{K(\mathbf{x}_i, \mathbf{x}_j)}{\sqrt{K(\mathbf{x}_i, \mathbf{x}_i)K(\mathbf{x}_j, \mathbf{x}_j)}} \right\} \quad (12)$$

where $f(\mathbf{x}_i)^{\text{MC}}$ is the multiclass uncertainty function defined by (9).

3. LARGE-MARGIN-BASED ACTIVE LEARNING

4. On the Need for a Diversity Criterion

- *Constraining the MS solution to pixels associated to different closest support*

$$\hat{\mathbf{x}}^{\text{cSV}} = \arg \min_{\mathbf{x}_i \in U^{\text{MS}}} \{ |f(\mathbf{x}_i, \omega)| \mid \text{cSV}_i \notin \text{cSV}_\theta \} \quad (13)$$

where $\theta = [1, \dots, q - 1]$ are the indices of the already selected candidates and cSV is the set of selected closest support vectors.

- *Finally, diversity can be ensured using clustering in the feature space*

$$\begin{aligned} \hat{\mathbf{x}}^{\text{MCLU-ECBD}} &= \arg \min_{\mathbf{x}_i \in c_m} \{ f(\mathbf{x}_i)^{\text{MC}} \}, \\ m &= [1, \dots, q], \quad \mathbf{x}_i \in U^{\text{MCLU}} \end{aligned} \quad (14)$$

where c_m is one among the q clusters defined with kernel k -means.

4. POSTERIOR PROBABILITY BASED ACTIVE LEARNING

□ The third class of methods uses the estimation of posterior probabilities of class membership (i.e., $P(y/x)$) to rank the candidates.

- *KL-Max*

The first idea is to sample the pixels whose inclusion in the training set

$$\hat{\mathbf{x}}^{\text{KL-max}} = \arg \max_{\mathbf{x}_i \in U} \left\{ \sum_{\omega \in N} \frac{1}{(u-1)} \times \text{KL} \left(p^+(\omega | \mathbf{x}) \parallel p(\omega | \mathbf{x}) \right) p(y_i^* = \omega | \mathbf{x}_i) \right\} \quad (16)$$

where the condition $n_{c_m}^{bSV} = 0$ ensures that the cluster queried does not contain any bounded support vector sampled at the previous iteration

4. POSTERIOR PROBABILITY BASED ACTIVE LEARNING

- *B. Breaking Ties (BT)*

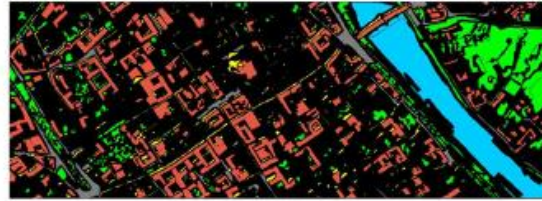
Another strategy, closer to the idea of EQB presented in Section III-A, consists of building a heuristic exploiting the conditional probability of predicting a given label for each candidate .

In this case, the per-class posterior probability is assessed fitting a sigmoid function to the SVM decision function [50]:

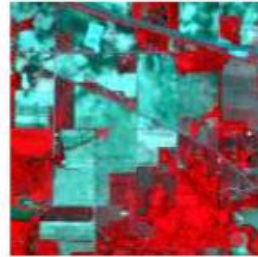
$$p(y_i^* = \omega \mid \mathbf{x}_i) = \frac{1}{1 + e^{(Af(\mathbf{x}_i, \omega) + B)}} \quad (18)$$

$$\hat{\mathbf{x}}^{\text{BT}} = \arg \min_{\mathbf{x}_i \in U} \left\{ \max_{\omega \in N} \{p(y_i^* = \omega \mid \mathbf{x}_i)\} - \max_{\omega \in N \setminus \omega^+} \{p(y_i^* = \omega \mid \mathbf{x}_i)\} \right\} \quad (19)$$

4. DATASETS



ROSIS Pavia



AVIRIS Indian Pines



QuickBird Zurich

Fig. 2. Images considered in the experiments: (top) ROSIS image of the city of Pavia, Italy (bands [56 – 31 – 6] and corresponding ground survey); (middle) AVIRIS Indian Pines hyperspectral data (bands [40 – 30 – 20] and corresponding ground survey); (bottom) QuickBird multispectral image of a suburb of the city of Zurich, Switzerland (bands [3 – 2 – 1] and corresponding ground survey).

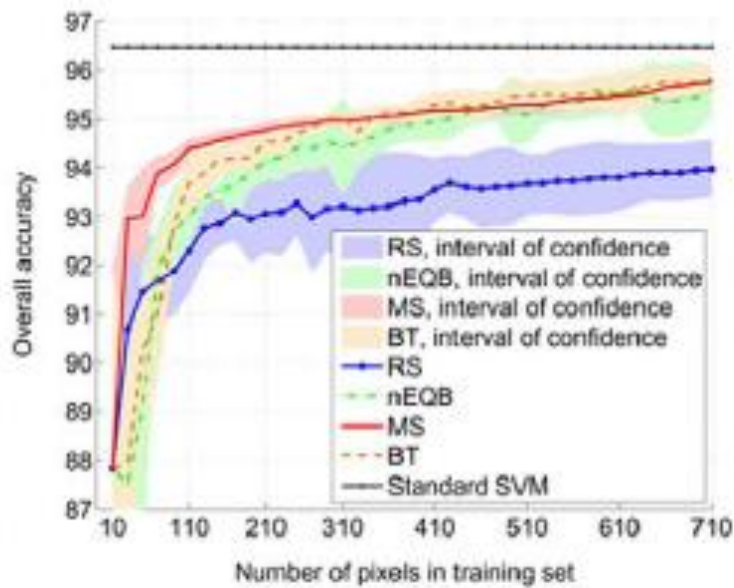
5. EXPERIMENTAL SETUP

- ❑ In the experiments, **SVM classifiers with RBF kernel** and LDA classifiers have been considered for the experiments.
- ❑ When using SVM, free parameters have been optimized by **five-fold cross validation** optimizing an accuracy criterion.
- ❑ The active learning algorithms have been run in two settings, adding **N+5** and **N+20** pixels per iteration.

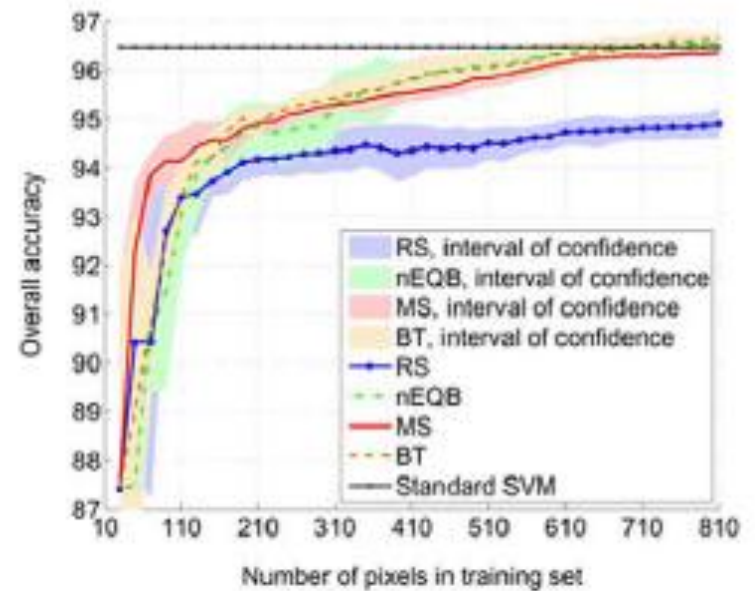
6. NUMERICAL RESULTS

Pavia ROSIS

$N + 5$



$N + 20$



5. CONCLUSION

- ❑ A series of heuristics have been classified by their characteristics into three families.
- ❑ Active learning has a strong potential for remote sensing data processing.
- ❑ Some recent examples can be found in the active selection of unlabeled pixels for semi-supervised classification.