

Active Learning Using Hint Information

Chun-Liang Li, Chun-Sung Ferng, and Hsuan-Tien Lin

{b97018, r99922054, htlin}@csie.ntu.edu.tw

Department of Computer Science, National Taiwan University

Keywords: Active Learning, Support Vector Machine

Abstract

The abundance of real-world data and limited labeling budget calls for active learning, which is an important learning paradigm for reducing human labeling efforts. Many recently developed active learning algorithms consider both uncertainty and representativeness when making querying decisions. However, exploiting representativeness with uncertainty concurrently usually requires tackling sophisticated and challenging learning tasks, such as clustering. In this paper, we propose a new active learning framework, called *hinted sampling*, which takes both uncertainty and representativeness into account in a simpler way. We design a novel active learning algorithm within

the hinted sampling framework with an extended support vector machine. Experimental results validate that the novel active learning algorithm can result in a better and more stable performance than that achieved by state-of-the-art algorithms. We also show that the hinted sampling framework allows improving another active learning algorithm designed from the transductive support vector machine.

1 Introduction

Labeled data are the basic ingredients in training a good model in machine learning. It is common in real-world applications when one needs to cope with a large amount of data with costly labeling steps. For example, in the medical domain, a doctor may be required to distinguish (label) cancer patients from non-cancer patients according to their clinical records (data). In such applications, an important issue is to achieve high accuracy within a limited labeling budget. This issue demands active learning (Settles, 2009), which is a machine learning setup that allows iteratively querying the labeling oracle (doctor) in a strategic manner to label some selected instances (clinic records). By using a suitable query strategy, an active learning approach can achieve high accuracy within a few querying iterations i.e., only a few calls to the costly labeling oracle (Settles, 2009).

One intuitive approaches in active learning is called uncertainty sampling (Lewis and Gale, 1994). This approach maintains a classifier on hand, and queries the most uncertain instances, whose uncertainty is measured by the closeness to the decision boundary of the classifier, to fine-tune the boundary. However, the performance of un-

certainty sampling becomes restricted owing to the limited view of the classifier. In other words, uncertainty sampling can be hair-splitting on the local instances that confuse the classifier, but not considering the global distribution of instances. Therefore, queries may not represent the underlying data distribution well, leading to unsatisfactory performance of uncertainty sampling (Nguyen and Smeulders, 2004; Settles and Craven, 2008; Settles, 2009; Huang et al., 2010).

As suggested by Cohn et al. (1996) as well as Xu et al. (2003), active learning can be improved by considering the unlabeled instances in order to query the instance that is not only uncertain to the classifier on hand but also “representative” to the global data distribution. There are many existing algorithms that use unlabeled information to improve the performance of active learning, such as representative sampling (Xu et al., 2003).

Representative sampling makes querying decisions by not only the uncertainty of each instance, but also the representativeness, which is measured by determining whether the instances reside in a dense area. Typical representative sampling algorithms (Xu et al., 2003; Nguyen and Smeulders, 2004; Dasgupta and Hsu, 2008) estimate the underlying data distribution via clustering methods. However, the performance of the algorithms depends on the result of clustering, which is a sophisticated and non-trivial task, especially when the instances are within a high dimensional space. Another state-of-the-art algorithm (Huang et al., 2010) models the representativeness by estimating the potential label assignment of the unlabeled instances on the basis of the min-max view of active learning (Hoi et al., 2008). The performance of this algorithm depends on the results of estimating the label assignments, which is also a complicated

task. Yet another representative sampling algorithm makes potential label assignment of the unlabeled instances from the view of transductive learning, such as transductive SVM (TSVM; Wang et al., 2011).

In this work, we propose a novel framework of active learning, hinted sampling, which considers the unlabeled instances as hints (Abu-Mostafa, 1995) of the global data distribution, instead of directly clustering them or estimating their label assignments. This leads to a simpler active learning algorithm. Similar to representative sampling, hinted sampling also considers both uncertainty and representativeness. Somehow hinted sampling enjoys the advantage of simplicity by avoiding the clustering or label-assignment estimation steps. We demonstrate the effectiveness of hinted sampling by designing a novel algorithm with support vector machine (SVM; Vapnik, 1998). In the algorithm, we extend the usual SVM to a novel formulation, HintSVM, which is easier to solve than either clustering or label-assignment estimation. We then study a simple hint selection strategy to improve the efficiency and effectiveness of the proposed algorithm. Experimental results demonstrate that the simple HintSVM is comparable to the best of both uncertainty sampling and representative sampling algorithms, and results in better and more stable performance than other state-of-the-art active learning algorithms.

To demonstrate the generality of hinted sampling, we further extend the TSVM approach for active learning (Wang et al., 2011) to HintTSVM to show that the proposed framework can benefit not only uncertainty sampling but also representative sampling. Experimental results confirm the promising performance of HintTSVM as well as the usefulness of the proposed hinted sampling framework.

The rest of the paper is organized as follows. Section 2 introduces the formal problem definition and reviews related works. Section 3 describes our proposed hinted sampling framework as well as the HintSVM algorithms with the simple hint selection strategy, and reports experiment results and comparisons. Section 4 discusses TSVM and HintTSVM with experimental justifications. Section 5 compares hinted sampling with selective sampling (Cohn et al., 1994), another strong family of active learning algorithms. Finally, we conclude in Section 6.

A short version of the paper appeared in the 2012 Asian Conference on Machine Learning (Li et al., 2012). The paper was then enriched by discussing more related works in Section 2, refining the hint selection strategies along with broader experiments in Section 3, the novel extension of TSVM (Wang et al., 2011) to hinted sampling for active learning in Section 4, and the comparison to selective sampling in Section 5.

2 Problem Definition and Related Works

In this work, we focus on pool-based active learning for binary classification, which is one of the most common setups in active learning (Lewis and Gale, 1994). At the initial stage, the learning algorithm is presented with a labeled data pool and an unlabeled data pool. We denote the labeled data pool by $\mathcal{D}_l = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$ and the unlabeled data pool by $\mathcal{D}_u = \{\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2, \dots, \tilde{\mathbf{x}}_M\}$, where the input vectors $\mathbf{x}_i, \tilde{\mathbf{x}}_j \in \mathbf{R}^d$ and the labels $y_i \in \{-1, 1\}$. Usually, the labeled data pool \mathcal{D}_l is relatively small or even empty, whereas the unlabeled data pool \mathcal{D}_u is assumed to be large. Active learning is an iterative process that contains R iterations of querying and learning. That is, an

active learning algorithm can be split into two parts: the querying algorithm \mathcal{Q} and the learning algorithm \mathcal{L} .

Using the initial $\mathcal{D}_l \cup \mathcal{D}_u$, the learning algorithm \mathcal{L} is first called to learn a decision function $f^{(0)}: \mathbf{R}^d \rightarrow \mathbf{R}$ from some hypothesis set \mathcal{F} , where the function $\text{sign}(f^{(0)}(\mathbf{x}))$ is taken for predicting the label of any input vector \mathbf{x} . Then, in iteration r , where $r = 1, 2, \dots, R$, the querying algorithm \mathcal{Q} is allowed to select an instance $\tilde{\mathbf{x}}_s \in \mathcal{D}_u$ and query its label y_s from a labeling oracle. After querying, $(\tilde{\mathbf{x}}_s, y_s)$ is added to the labeled pool \mathcal{D}_l and $\tilde{\mathbf{x}}_s$ is removed from the unlabeled pool \mathcal{D}_u . The learning algorithm \mathcal{L} then learns a decision function $f^{(r)}$ from the updated $\mathcal{D}_l \cup \mathcal{D}_u$. The goal of active learning is to use the limited querying and learning opportunities properly to obtain a decent list of decision functions $[f^{(1)}, f^{(2)}, \dots, f^{(R)}]$ that can achieve low out-of-sample (test) error rates.

As discussed in a detailed survey (Settles, 2009), there are many active learning algorithms for binary classification. In this paper, we shall review some relevant ones. One of the most intuitive families of algorithms is called uncertainty sampling (Lewis and Gale, 1994). As the name suggests, the querying algorithm \mathcal{Q} of uncertainty sampling queries the most uncertain $\tilde{\mathbf{x}}_s \in \mathcal{D}_u$, where the uncertainty for each input vector $\tilde{\mathbf{x}}_j \in \mathcal{D}_u$ is usually computed by re-using the decision function $f^{(r-1)}$ returned from the learning algorithm \mathcal{L} . For instance, Tong and Koller (2000) take the support vector machine (SVM; Vapnik, 1998) as \mathcal{L} and measure the uncertainty of $\tilde{\mathbf{x}}_j$ by the distance between $\tilde{\mathbf{x}}_j$ and the boundary $f^{(r-1)} = 0$. In other words, the algorithm in Tong and Koller (2000) queries the $\tilde{\mathbf{x}}_s$ that is closest to the boundary.

Uncertainty sampling can be viewed as a greedy approach that queries instances

from the viewpoint only of the decision function $f^{(r-1)}$. When the decision function is not close enough to the ideal one, however, this limited viewpoint can hinder the performance of the active learning algorithm. Thus, Cohn et al. (1996) suggest that the viewpoint of the unlabeled pool \mathcal{D}_u should also be included. Their idea leads to another family of active learning algorithms, called representative sampling (Xu et al., 2003), or density-weighted sampling (Settles, 2009). Representative sampling takes both the uncertainty and the representativeness of each $\tilde{\mathbf{x}}_j \in \mathcal{D}_u$ into account concurrently in the querying algorithm \mathcal{Q} , where the representativeness of $\tilde{\mathbf{x}}_j$ with respect to \mathcal{D}_u is measured by the density of its neighborhood area. For instance, Xu et al. (2003) employ the SVM as the learning algorithm \mathcal{L} as do Tong and Koller (2000). They use a querying algorithm \mathcal{Q} that first clusters the unlabeled instances near the boundary of $f^{(r-1)}$ by a K -means algorithm, and then queries one of the centers of those clusters. In other words, the queried instance is not only uncertain for $f^{(r-1)}$ but also representative for \mathcal{D}_u . Some other works estimate the representativeness with a generative model. For instance, Nguyen and Smeulders (2004) propose a querying algorithm \mathcal{Q} that uses multiple Gaussian distributions to cluster all input vectors $\mathbf{x}_i \in \mathcal{D}_l$, $\tilde{\mathbf{x}}_j \in \mathcal{D}_u$ and estimate the prior probability $p(\mathbf{x})$; \mathcal{Q} then makes querying decisions based on the product of the prior probability and some uncertainty measurement. The idea of estimating the representativeness via clustering is a core element of many representative sampling algorithms (Xu et al., 2003; Nguyen and Smeulders, 2004; Dasgupta and Hsu, 2008). Nevertheless, clustering is a challenging task and it is not always easy to achieve satisfactory clustering performance. When the clustering performance is not satisfactory, it has been observed (Donmez et al., 2007; Huang et al., 2010) that repre-

sentative sampling algorithms could fail to achieve decent performance. In other words, the clustering step is usually the bottleneck of representative sampling.

Huang et al. (2010) propose an improved algorithm that models representativeness without clustering. In the algorithm, the usefulness of each $\tilde{\mathbf{x}}_j$, which implicitly contains both uncertainty and representativeness, is estimated by using a technique in semi-supervised learning (Hoi et al., 2008) that checks approximately all possible label assignments for each unlabeled $\tilde{\mathbf{x}}_j \in \mathcal{D}_u$. The querying algorithm \mathcal{Q} proposed (Huang et al., 2010) is based on the usefulness of each $\tilde{\mathbf{x}}_j$; the learning algorithm \mathcal{L} is simply a stand-alone SVM. While the active learning algorithm (Huang et al., 2010) often achieves promising empirical results, its bottleneck is the label-estimation step, which is rather sophisticated and thus not always leading to satisfactory performance.

Another improvement of representative sampling is presented by Donmez et al. (2007), who report that representative sampling is less efficient than uncertainty sampling for later iterations, in which the decision function is closer to the ideal one. To combine the best properties of uncertainty sampling and representative sampling, Donmez et al. (2007) propose a mixed algorithm by extending representative sampling (Nguyen and Smeulders, 2004). The proposed query algorithm \mathcal{Q} (Donmez et al., 2007) is split into two stages. The first stage performs representative sampling (Nguyen and Smeulders, 2004) while estimating the expected error reduction. When the expected reduction is smaller than a given threshold, the querying algorithm \mathcal{Q} switches to uncertainty sampling for fine-tuning the decision boundary. The bottleneck of the algorithm (Donmez et al., 2007) is still the clustering step in the first stage.

One simple algorithm employs transductive SVM (TSVM; Joachims, 1999b) to re-

place SVM for uncertainty sampling (Wang et al., 2011). Note that TSVM aims to estimate the labels of unlabeled data to maximize its margin, which is similar to the algorithm proposed by Huang et al. (2010). Therefore, using TSVM to replace SVM in uncertainty sampling for querying can also be viewed a concrete instance of representative sampling. We will have more detailed discussions of the difference between the two algorithms in Section 4.

Finally, we connect our discussions to a strong family of active learning algorithms, called selective sampling (Cohn et al., 1994; Settles, 2009). Instead of query the most uncertain instance like uncertainty sampling, selective sampling defines a *region* of uncertainty, and makes a query (samples) from the region only. The original work of Cohn et al. (1994) maintains a hypothesis set \mathcal{F}_r of all the decision functions f that agree with all the labelled instances in iteration r , and samples within the region that the committee in \mathcal{F}_r still disagree on. The sampling-within-region step filters the underlying distribution that generated the instances. In the common case when the underlying distribution is unknown but a (large) pool from the distribution is available, as in the pool-based setting, one simple way is to keep sampling an instance uniformly from \mathcal{D}_u until one instance falls into the uncertainty region (Cohn et al., 1994).

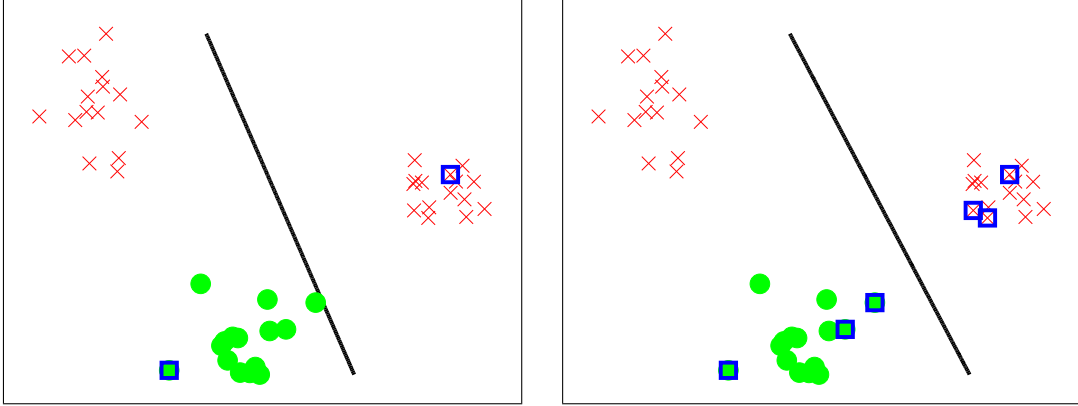
The idea of selective sampling can also be easily extended to the stream-based setup (Cohn et al., 1994). The setup assumes that instances $\tilde{\mathbf{x}}$ arrive as a data stream and the learner is allowed to decide whether to query the label of each instance on the fly. The stream-based setup and the original selective sampling algorithm form the basis of providing theoretical guarantees on the performance of many modern (stream-based) active learning algorithms (Balcan et al., 2006; Dasgupta et al., 2007; Hanneke, 2007;

Dasgupta, 2011). Under the stream-based setup, different selective sampling algorithms can then be viewed as different approaches for calculating the uncertainty of an instance to guide the query decisions. For instance, Beygelzimer et al. (2010) propose the importance-weighted active learning algorithm (IWAL) that considers a hypothesis set \mathcal{F}'_r that contains all the hypotheses that are close to the best hypothesis subject to the labelled instances, and takes the largest disagreement of the hypotheses in \mathcal{F}'_r for calculating the uncertainty of an instance.

As suggested implicitly in the algorithm steps of Hanneke (2007), selective sampling algorithms in the stream-based setup can be naturally applied to pool-based algorithm by simulating a data stream from \mathcal{D}_u . That is, in each iteration, we keep sampling an instance uniformly from \mathcal{D}_u to be $\tilde{\mathbf{x}}$ and feed $\tilde{\mathbf{x}}$ to the underlying selective sampling algorithm until the algorithm decides to query the label of the instance. The simulation allows the resulting algorithm to be representative because of the uniform sampling. That is, under the simulation, selective sampling can be viewed as yet another form of representative sampling that queries the denser part of the uncertainty region with higher probability. Later in Section 5, IWAL coupled with the simulation will be empirically compared with other uncertainty and representative sampling approaches, including our proposed one.

3 Hinted Sampling Framework

Instead of facing the challenges of either clustering or label-estimation, we propose to view the information in \mathcal{D}_u differently. In particular, the unlabeled instances $\tilde{\mathbf{x}}_j \in$

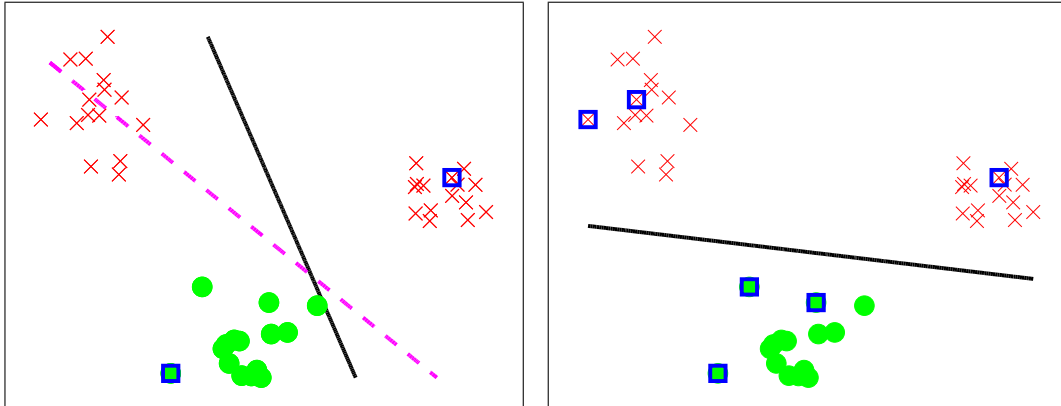


(a) the decision function (black) obtained from two labeled (blue) instances (b) when using the decision function in (a) for uncertainty sampling, the top-left cluster keeps being ignored

Figure 1: illustration of uncertainty sampling

\mathcal{D}_u are taken as hints (Abu-Mostafa, 1995) that guide the querying algorithm \mathcal{Q} . The idea of using hints leads to a simpler active learning algorithm with better empirical performance.

First, we illustrate the potential drawback of uncertainty sampling with a linear SVM classifier (Vapnik, 1998), which is applied to a two-dimensional artificial dataset. Figure 1 shows the artificial dataset, which consists of three clusters, each of which contains instances of a particular class. We denote one class by a red cross and the other by a filled green circle. The labeled instances in \mathcal{D}_l are marked with a blue square while other instances are in \mathcal{D}_u . In Figure 1(a), the initial two labeled instances reside in two of the clusters with different labels. The initial decision function $f^{(0)}$ trained on the labeled instances (from the two clusters) is not aware of the third cluster. The decision function $f^{(0)}$ then mis-classifies the instances in the third cluster, and causes the querying algorithm \mathcal{Q} (which is based on $f^{(0)}$) to query only from the instances near



(a) the hinted query boundary (dashed magenta) that is aware of the top-left cluster (b) when using the hinted query function in (a) for uncertainty sampling, all three clusters are explored

Figure 2: illustration of hinted sampling

the “wrong” boundary rather than exploring the third cluster. After several iterations, as shown in Figure 1(b), the uncertainty sampling algorithm still outputs an unsatisfactory decision function that mis-classifies the entire unqueried (third) cluster.

The unsatisfactory performance of uncertainty sampling originates in its lack of awareness of candidate unlabeled instances that should be queried. When trained on only a few labeled instances, the resulting (linear) decision function is overly confident about the unlabeled instances that are far from the boundary. Intuitively, uncertainty sampling could be improved if the querying algorithm \mathcal{Q} were aware of and less confident about the unqueried regions. Both clustering (Nguyen and Smeulders, 2004) and label-estimation (Huang et al., 2010) are based on this intuition, but they explore the unlabeled regions in a rather sophisticated way.

We propose a simpler alternative as follows. Note that the uncertainty sampling

algorithm measures the uncertainty by the distance between instances and the boundary. In order to make \mathcal{Q} less confident about the unlabeled instances, we seek a “query boundary” that not only classifies the labeled instances correctly but also passes through the unqueried regions, denoted by the dashed magenta line in Figure 2(a). Then, in the later iterations, the query algorithm \mathcal{Q} , using the query boundary, would be less confident about the unqueried regions, and thus be able to explore them. The instances in the unqueried regions give hints as to where the query boundary should pass. Using these hints about the unqueried regions, the uncertainty sampling algorithm can take both uncertainty and the underlying distribution into account concurrently, and achieve better performance, as shown in Figure 2(b).

Based on this idea, we propose a novel active learning framework, *hinted sampling*. The learning algorithm \mathcal{L} in hinted sampling is similar to that in uncertainty sampling, but the querying algorithm is different. In particular, the querying algorithm \mathcal{Q} is provided with some unlabeled instances, called the hint pool $\mathcal{D}_h \subseteq \mathcal{D}_u$. When properly using the information in the hint pool \mathcal{D}_h , both uncertainty and representativeness can be considered concurrently to obtain a query boundary that assists \mathcal{Q} in making query decisions. We sketch the framework of *Active Learning under Hinted Sampling* (ALHS) in Algorithm 1.

Next, we design a concrete active learning algorithm of hinted sampling based on SVM, which is also used as the core of many state-of-the-art algorithms (Tong and Koller, 2000; Xu et al., 2003; Huang et al., 2010), as both \mathcal{L} and \mathcal{Q} . Before illustrating the complete algorithm, we show how SVM can be appropriately extended to use the information in \mathcal{D}_h for \mathcal{Q} .

Algorithm 1 Active Learning under Hinted Sampling (ALHS) framework

Input: the number of rounds R ; a labeled pool \mathcal{D}_l ; an unlabeled pool \mathcal{D}_h ; parameters $\theta_{\mathcal{Q}}$ for querying algorithm and $\theta_{\mathcal{L}}$ for learning algorithm

Output: decision functions $f^{(1)}, \dots, f^{(R)}$

For $r \leftarrow 1$ **to** R **do**

Select \mathcal{D}_h from \mathcal{D}_u

$g \leftarrow \mathcal{Q}(\theta_{\mathcal{Q}}, \mathcal{D}_h \cup \mathcal{D}_l)$

$(\tilde{\mathbf{x}}_s, y_s) \leftarrow \text{Query}(g, \mathcal{D}_u)$

$\mathcal{D}_u \leftarrow \mathcal{D}_u \setminus \tilde{\mathbf{x}}_s; \mathcal{D}_l \leftarrow \mathcal{D}_l \cup (\tilde{\mathbf{x}}_s, y_s)$

$f^{(r)} \leftarrow \mathcal{L}(\theta_{\mathcal{L}}, \mathcal{D}_l)$

End

3.1 HintSVM

The extended SVM is called HintSVM, which takes hints into account. The goal of HintSVM is to locate a query boundary which does well on two objectives: (1) classifying labeled instances in \mathcal{D}_l , and (2) being close to the unlabeled instances in hint pool \mathcal{D}_h . Note that the two objectives are different from the usual semi-supervised SVM (Bennett and Demiriz, 1998) such as transductive SVM (Joachims, 1999b), which pushes the unlabeled instances away from the decision boundary.

The first objective matches an ordinary support vector classification (SVC) problem. To deal with the second objective, we consider ϵ -support vector regression (ϵ -SVR) and set regression targets to 0 for all instances in \mathcal{D}_h , which means that instances in \mathcal{D}_h should be close to the query boundary. By combining the objective functions of SVC and ϵ -SVR together, HintSVM solves the following convex optimization problem,

which simultaneously achieves the two objectives.

$$\begin{aligned}
\min_{\mathbf{w}, b, \xi, \tilde{\xi}, \tilde{\xi}^*} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C_l \sum_{i=1}^{|\mathcal{D}_l|} \xi_i + C_h \sum_{j=1}^{|\mathcal{D}_h|} (\tilde{\xi}_j + \tilde{\xi}_j^*) \\
\text{subject to} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \text{for } (\mathbf{x}_i, y_i) \in \mathcal{D}_l, \\
& \mathbf{w}^T \tilde{\mathbf{x}}_j + b \leq \epsilon + \tilde{\xi}_j \quad \text{for } \mathbf{x}_j \in \mathcal{D}_h, \\
& -(\mathbf{w}^T \tilde{\mathbf{x}}_j + b) \leq \epsilon + \tilde{\xi}_j^* \quad \text{for } \mathbf{x}_j \in \mathcal{D}_h, \\
& \xi_i \geq 0 \quad \text{for } (\mathbf{x}_i, y_i) \in \mathcal{D}_l, \\
& \tilde{\xi}_j, \tilde{\xi}_j^* \geq 0 \quad \text{for } \mathbf{x}_j \in \mathcal{D}_h.
\end{aligned} \tag{1}$$

Here ϵ is the margin of tolerance for being close to the boundary, and C_l, C_h are the weights of the classification errors (on \mathcal{D}_l) and hint errors (on \mathcal{D}_h), respectively. Similar to the usual SVC and ϵ -SVR, the convex optimization problem can be transformed to the dual form to allow using the kernel trick. Define $\hat{\mathbf{x}}_i = \mathbf{x}_i$, $\hat{\mathbf{x}}_{|\mathcal{D}_l|+j} = \hat{\mathbf{x}}_{|\mathcal{D}_l|+|\mathcal{D}_h|+j} = \tilde{\mathbf{x}}_j$, $\hat{y}_i = y_i$, $\hat{y}_{|\mathcal{D}_l|+j} = 1$, and $\hat{y}_{|\mathcal{D}_l|+|\mathcal{D}_h|+j} = -1$ for $1 \leq i \leq |\mathcal{D}_l|$ and $1 \leq j \leq |\mathcal{D}_h|$. The dual problem of (1) can be written as follows:

$$\begin{aligned}
\min_{\alpha} \quad & \frac{1}{2} \alpha^T Q \alpha + \mathbf{p}^T \alpha \\
\text{subject to} \quad & \hat{\mathbf{y}}^T \alpha = 0, \\
& 0 \leq \alpha_i \leq C_l \quad \text{for } i = 1, 2, \dots, |\mathcal{D}_l|, \\
& 0 \leq \alpha_j \leq C_h \quad \text{for } j = |\mathcal{D}_l| + 1, \dots, |\mathcal{D}_l| + 2|\mathcal{D}_h|,
\end{aligned}$$

where $p_i = -1$, $p_j = \epsilon$, and $Q_{ab} = \hat{y}_a \hat{y}_b \hat{\mathbf{x}}_a^T \hat{\mathbf{x}}_b$. The derived dual form can be easily solved by any state-of-the-art quadratic programming solver, such as the one implemented in LIBSVM (Chang and Lin, 2011).

3.2 Hint Selection Strategy

A naïve strategy for selecting a proper hint pool $\mathcal{D}_h \subseteq \mathcal{D}_u$ is to directly let $\mathcal{D}_h = \mathcal{D}_u$, which retains all the information about the unlabeled data. However, given that the size of \mathcal{D}_u is usually much larger than the size of \mathcal{D}_l , this strategy may cause the hints to overwhelm HintSVM, which leads to performance and computational concerns. In our earlier version of this work (Li et al., 2012), a specifically designed selection strategy that drops hints with radial functions have been studied. We have conducted some broader studies that suggest the sufficiency of using a simpler uniform sampling strategy. Next, we will use the strategy to demonstrate the essence, validity and usefulness of hint information.

3.3 Hint Influence Control

Tong and Koller (2000) show that uncertainty sampling can outperform uniform sampling when enough examples have been queried. Thus, after querying more examples, there can be advantages by changing the focus of the active learning approach to “refine” the boundary by uncertainty sampling. (Donmez et al., 2007) tries to dynamically balance the representative sampling and uncertainty sampling. Our earlier work (Li et al., 2012) exploits two strategies, *hint dropping* and *hint termination*, to achieve this goal. Similar ideas have also been widely used in the bandit problem (Langford and Zhang, 2007) to balance exploration and exploitation. Here we take a simple alternative by multiplying the parameter C_h by a ratio δ each iteration, where $0 < \delta < 1$, to gradually change our focus to uncertainty sampling. That is, in iteration r , the cost parameter $C_h^{(r)}$ of hint instances is $C_h^{(1)} \times \delta^{r-1}$. After enough many iterations, $C_h^{(r)}$

will be closed to 0, which essentially transforms hinted sampling to typical uncertainty sampling.

3.4 Hinted Sampling with HintSVM

Next, we incorporate the proposed ALHS with the derived HintSVM formulation to make a novel active learning algorithm, ALHS-SVM.

The querying algorithm \mathcal{Q} of ALHS-SVM selects unlabeled instances from the unlabeled pool \mathcal{D}_u as the hint pool \mathcal{D}_h and trains HintSVM from \mathcal{D}_l and \mathcal{D}_h to obtain the query boundary for uncertainty sampling. The use of both \mathcal{D}_l and \mathcal{D}_h combines uncertainty and representativeness. The learning algorithm \mathcal{L} of ALHS-SVM, on the other hand, trains a stand-alone SVM from \mathcal{D}_l to get a decision function $f^{(r)}$, just like \mathcal{L} in uncertainty sampling (Tong and Koller, 2000). The full ALHS-SVM algorithm is listed in Algorithm 2.

Algorithm 2 The ALHS-SVM algorithm

Input: the number of rounds R ; a labeled pool \mathcal{D}_l ; an unlabeled pool \mathcal{D}_u ; parameters for HintSVM and SVM; ratio δ

Output: decision functions $f^{(1)}, \dots, f^{(R)}$

For $r \leftarrow 1$ **to** R **do**

Uniformly select \mathcal{D}_h from \mathcal{D}_u

$g \leftarrow \text{Train HintSVM}(C_h, C_l, \epsilon, \mathcal{D}_h, \mathcal{D}_l)$

$(\tilde{\mathbf{x}}_s, y_s) \leftarrow \text{Query}(g, \mathcal{D}_u)$

$\mathcal{D}_u \leftarrow \mathcal{D}_u \setminus \tilde{\mathbf{x}}_s; \mathcal{D}_l \leftarrow \mathcal{D}_l \cup (\tilde{\mathbf{x}}_s, y_s)$

$f^{(r)} \leftarrow \text{Train SVM}(C, \mathcal{D}_l)$

$C_h \leftarrow C_h \times \delta$

End

Uncertainty sampling with SVM is a special case of ALHS-SVM when always set-

ting $C_h = 0$. In other words, ALHS can be viewed as a generalization of uncertainty sampling that considers representativeness through the hints. The simple use of hints avoids the challenges in clustering or label-estimation steps.

3.5 Experimental Studies of ALHS-SVM

Next, we compared the proposed ALHS-SVM algorithm with the following active learning algorithms: (1) UNCERTAIN (Tong and Koller, 2000): uncertainty sampling with SVM, (2) REPRESENT (Xu et al., 2003): representative sampling with SVM and clustering, (3) DUAL (Donmez et al., 2007): mixture of uncertainty and representative sampling, (4) QUIRE (Huang et al., 2010): representative sampling with label estimation based on the min-max view. We also list the results of another related active learning algorithm, TSVM-SVM (Wang et al., 2011), which conducts representative sampling with transductive SVM, and will compare it with ALHS-SVM in detail in Section 4.

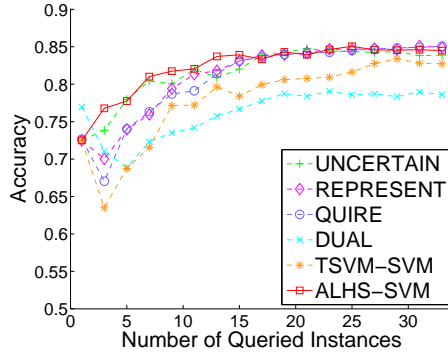
We conducted experiments on eight UCI benchmarks (Frank and Asuncion, 2010), which are *australian*, *diabetes*, *german*, *splice*, *wdbc*, *letetrMvsN*, *letterVvsY* (Donmez et al., 2007; Huang et al., 2010) and *segment-binary* (Ratsch et al., 2001; Donmez et al., 2007) as taken by other related works. For each dataset, we randomly divided it into two parts with equal size. One part was treated as the unlabeled pool \mathcal{D}_u for active learning algorithms. The other part was reserved as the test set. Before querying, we randomly select one positive instance and one negative instance to form the labeled pool \mathcal{D}_l . For each dataset, we ran the algorithms 20 times with different random splits.

There are two usage scenarios of active learning algorithms. For one, there are

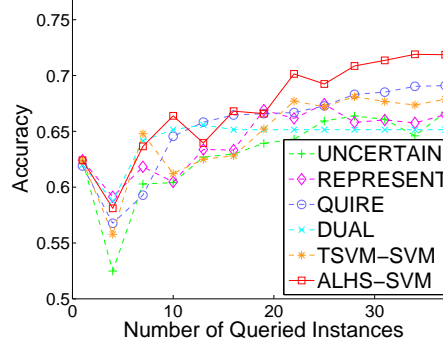
some (related) labelled data on hand for tuning reasonable parameters, and some works like (Zha et al., 2012; Wang and Ye, 2013; Quanquan Gu, 2014) are based on this setting. For the other, active learning is conducted from scratch and there is little information about reasonable parameters. In that case, the algorithms can only rely on fixed, pre-defined parameters, and some other works (Donmez et al., 2007; Huang et al., 2010; Cao et al., 2014; Chen et al., 2014) are based on this setting. For instance, in the micro-processor design space exploration, the simulation will take hundreds of hours, which makes tuning parameters expensive and infeasible (Chen et al., 2014). In this paper, we provide results for *both* settings for completeness. We will first show the results on the fixed setting and then the tuned setting.

Comparison under Fixed Parameters. In the experiments, We adapt SVM-light (Joachims, 1999a) for TSVM and LIBSVM (Chang and Lin, 2011) for other SVM-based algorithms with the RBF kernel and the default parameters, except for $C = 5$. Correspondingly, the parameter λ in the works of Donmez et al. (2007) and Huang et al. (2010) was set to $\lambda = \frac{1}{C}$. These parameters ensure that all four algorithms behave in a stable manner. For ALHS-SVM, we fixed $\delta = 0.5$ and uniformly sample 10% data from \mathcal{D}_u as \mathcal{D}_h without any further tuning for each dataset. For other algorithms, we take the parameters in the original papers.

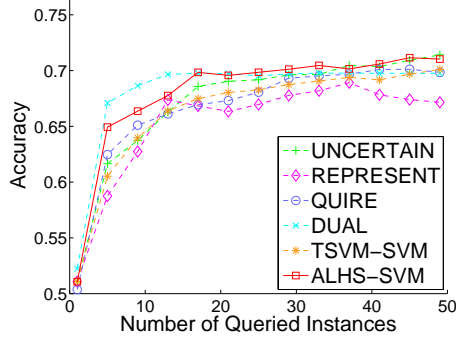
Figure 3 presents the accuracy of different active learning algorithms along with the number of rounds R , which equals the number of queried instances. Tables 1 and 2 list the mean and standard error of accuracy when $R = |\mathcal{D}_u| \times 5\%$ and $R = |\mathcal{D}_u| \times 10\%$, respectively. The highest mean accuracy is shown in boldface for each dataset. We



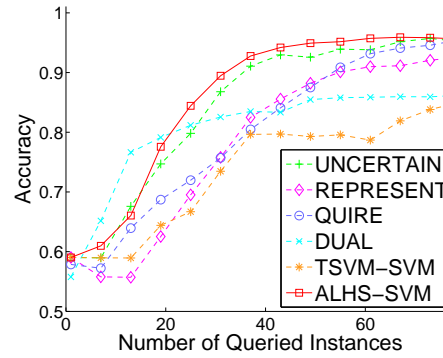
(a) *australian*



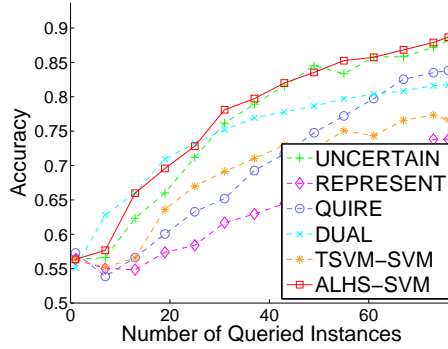
(b) *diabetes*



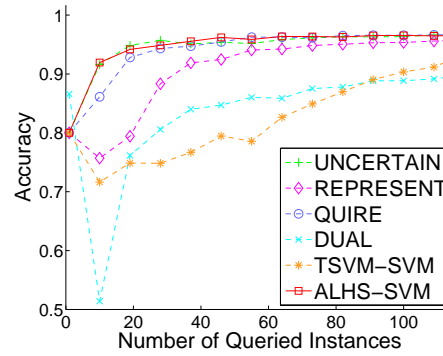
(c) *german*



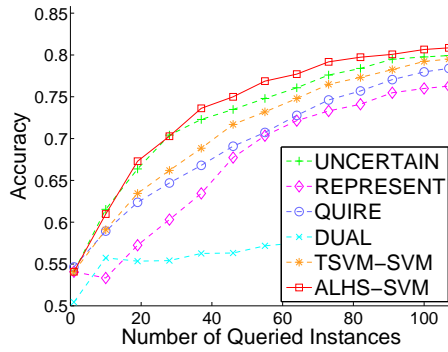
(d) *leterMvsN*



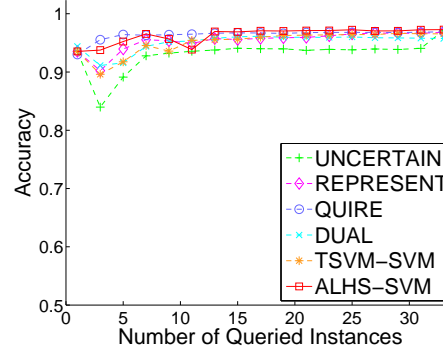
(e) *letterVvsY*



(f) *segment*



(g) *splice*



(h) *wdbc*

Figure 3: Comparison on different datasets

Table 1: Comparison on accuracy (mean \pm se) after querying 5% of unlabeled pool

data	Algorithms (%), the highest accuracy for each dataset is in boldface					
	UNCERTAIN	REPRESENT	QUIRE	DUAL	TSVM-SVM	ALHS-SVM
<i>australian</i>	82.188 \pm 1.571	83.739 \pm 0.548	82.319 \pm 1.126	81.304 \pm 0.647	78.116 \pm 1.490	83.362 \pm 0.439
<i>diabetes</i>	63.229 \pm 2.767	66.758 \pm 0.505	66.771 \pm 0.960	65.143 \pm 0.381	67.148 \pm 0.570	68.438 \pm 0.709
<i>german</i>	69.060 \pm 0.497	67.240 \pm 1.099	68.750 \pm 0.605	69.620 \pm 1.323	68.160 \pm 0.382	69.330 \pm 0.373
<i>letterMvsN</i>	89.632 \pm 1.103	83.463 \pm 1.348	81.372 \pm 1.693	83.437 \pm 1.211	80.051 \pm 1.165	91.112 \pm 0.444
<i>letterVvsY</i>	79.245 \pm 1.176	63.523 \pm 2.335	68.516 \pm 2.132	76.213 \pm 1.549	71.258 \pm 0.973	79.300 \pm 0.695
<i>segment</i>	95.437 \pm 0.367	94.390 \pm 0.482	96.074 \pm 0.224	86.078 \pm 2.834	79.160 \pm 0.716	96.095 \pm 0.204
<i>splice</i>	74.430 \pm 0.606	69.117 \pm 1.452	70.340 \pm 0.942	56.969 \pm 0.576	72.847 \pm 0.332	76.334 \pm 0.315
<i>wdbc</i>	93.842 \pm 3.137	95.616 \pm 0.711	96.613 \pm 0.230	96.056 \pm 0.250	95.777 \pm 0.264	97.020 \pm 0.123

Table 2: Comparison on accuracy (mean \pm se) after querying 10% of unlabeled pool

data	Algorithms (%), the highest accuracy for each dataset is in boldface					
	UNCERTAIN	REPRESENT	QUIRE	DUAL	TSVM-SVM	ALHS
<i>australian</i>	83.884 \pm 0.460	84.884 \pm 0.367	84.870 \pm 0.455	81.174 \pm 0.798	82.725 \pm 0.628	84.4391 \pm 0.259
<i>diabetes</i>	66.706 \pm 2.632	66.484 \pm 1.223	67.500 \pm 1.337	65.143 \pm 0.381	67.878 \pm 0.609	71.758 \pm 0.618
<i>german</i>	71.410 \pm 0.488	67.150 \pm 0.773	70.250 \pm 0.560	69.760 \pm 0.299	70.110 \pm 0.342	71.030 \pm 0.257
<i>letterMvsN</i>	95.369 \pm 0.315	92.433 \pm 0.777	95.114 \pm 0.486	86.893 \pm 0.870	84.625 \pm 0.879	95.578 \pm 0.214
<i>letterVvsY</i>	88.213 \pm 0.635	73.806 \pm 1.551	84.723 \pm 0.891	80.123 \pm 1.359	76.645 \pm 0.677	88.561 \pm 0.431
<i>segment</i>	96.528 \pm 0.143	95.684 \pm 0.155	96.658 \pm 0.110	89.519 \pm 1.760	92.420 \pm 0.217	96.445 \pm 0.100
<i>splice</i>	79.931 \pm 0.274	76.274 \pm 0.895	78.560 \pm 0.648	58.947 \pm 0.853	79.503 \pm 0.197	80.750 \pm 0.143
<i>wdbc</i>	97.155 \pm 0.141	96.818 \pm 0.191	96.862 \pm 0.206	95.748 \pm 0.247	96.554 \pm 0.219	97.170 \pm 0.127

also conducted the t -test at 95% significance level (Melville and Mooney, 2004; Guo and Greiner, 2007; Donmez et al., 2007). The t -test results are given in Table 3, which summarizes the number of datasets in which ALHS-SVM performs significantly better (or worse) than the other algorithms.

Comparison between ALHS-SVM and Uncertainty Sampling For some datasets, such as *wdbc* and *diabetes* in Figure 3(h) and 3(b), the result for UNCERTAIN is unsatisfactory. This unsatisfactory performance is possibly caused by the lack of awareness of unlabeled instances, which echoes our illustration in Figure 1. Note that we have considered some more aggressive querying criteria (Tong and Koller, 2000) than UNCERTAIN on the side, and have observed that those criteria are designed for hard-margin SVM and hence can be worse than UNCERTAIN with soft-margin SVM in our

Table 3: ALHS-SVM versus the other algorithm based on t -test at 95% significance level

Percentage of queries	Algorithms (win/tie/loss)				
	UNCERTAIN	REPRESENT	QUIRE	DUAL	TSVM-SVM
5%	1/7/0	6/2/0	5/3/0	4/4/0	6/2/0
10%	3/5/0	7/1/0	6/2/0	5/3/0	8/0/0

experiments. Thus, we excluded them from the tables. In these two cases, ALHS-SVM improves on UNCERTAIN with much lower standard error by using the hint information to avoid the worse local optimal. The results demonstrate the validity of the proposed ALHS framework. For other datasets which UNCERTAIN performs well on them, ALHS-SVM is still competitive and can sometimes reach even better performance. For instance, in *splice*, ALHS-SVM results in significantly higher accuracy after 30 queries. The observation further justifies that the hint information can be useful in boosting the performance of UNCERTAIN.

Comparison between ALHS-SVM and Representative Sampling We will leave the detailed comparison to TSVM-SVM in the next section. For DUAL, CLUSTER and QUIRE, we see that ALHS-SVM is only worse than DUAL on *german* and *letter M vs N* when ALHS-SVM has not queried enough many instances. For all other datasets and situations, ALHS-SVM generally results in better performance than all the other three algorithms. For instance, in Figure 3(g), since *splice* is a larger and higher dimensional dataset, representative sampling algorithms that perform clustering (REPRESENT, DUAL) or label estimation (QUIRE) fail to reach a decent performance. We attribute the results to the fact that it is usually non-trivial to perform distribution esti-

mation, clustering or label estimations in a high dimensional space. On the other hand, ALHS uses the hint information without aggressive assumptions, and can thus result in better and stabler performance.

In summary, Figure 3 shows that ALHS-SVM can achieve comparable results to those of the best representative sampling and uncertainty sampling algorithms. As shown in Tables 1 and 2, after querying 5% of the unlabeled instances (Table 1), ALHS-SVM achieves the highest mean accuracy in 6 out of 8 datasets; after querying 10% of unlabeled instances (Table 2), ALHS-SVM achieves the highest mean accuracy in 5 out of 8 datasets. Table 3 further confirms that ALHS-SVM usually outperforms each of the other algorithms at the 95% significance level.

Comparison under Tuned Parameters. Besides using fixed parameters with default values for the learning algorithm, we also study the performance of different active learning approaches with tuned parameters. Next, we compare ALHS-SVM with UNCERTAIN, REPRESENT and QUIRE. We leave out DUAL because the learning algorithm of DUAL is tightly coupled with its query algorithm, and it is difficult to tune the parameters of the learning algorithm fairly in advance.

For each dataset, we randomly divide it into three parts. The first part (fully labelled) is used for choosing parameters (C, γ) of SVM by cross-validation from $\{5 \times 10^{-2}, 5 \times 10^{-1}, \dots, 5 \times 10^2\} \times \{5 \times 10^{-2}, 5 \times 10^{-1}, \dots, 5 \times 10^2\}$. The second part is treated as the unlabeled pool for active learning, and the third part was reserved as the test set. The ratio of three parts is 1 : 2 : 2. The other settings are the same as the descriptions in Section 3.5.

Table 4: Comparison on accuracy (mean \pm se) after querying 5% of unlabeled pool with tuned parameters

Algorithms (%), the highest accuracy for each dataset is in boldface				
data	UNCERTAIN	REPRESENT	QUIRE	ALHS
<i>australian</i>	81.123 \pm 1.708	76.866 \pm 2.411	79.656 \pm 1.860	83.533 \pm 0.586
<i>diabetes</i>	65.649 \pm 2.040	67.597 \pm 1.476	64.659 \pm 2.000	70.308 \pm 0.924
<i>german</i>	69.050 \pm 0.885	66.112 \pm 1.074	69.300 \pm 0.613	68.112 \pm 1.007
<i>letterMvsN</i>	91.198 \pm 1.005	91.175 \pm 0.921	90.214 \pm 1.053	92.532 \pm 0.755
<i>letterVvsY</i>	87.871 \pm 0.649	82.774 \pm 0.812	87.347 \pm 0.533	89.831 \pm 0.569
<i>segment</i>	96.645 \pm 0.137	94.167 \pm 0.374	92.267 \pm 1.568	96.672 \pm 0.149
<i>splice</i>	72.195 \pm 0.394	75.011 \pm 0.429	71.891 \pm 0.619	74.379 \pm 0.414
<i>wdbc</i>	94.927 \pm 1.206	95.604 \pm 0.424	82.930 \pm 6.492	92.985 \pm 2.192

Table 5: Comparison on accuracy (mean \pm se) after querying 10% of unlabeled pool with tuned parameters

Algorithms (%), the highest accuracy for each dataset is in boldface				
data	UNCERTAIN	REPRESENT	QUIRE	ALHS
<i>australian</i>	83.478 \pm 0.609	80.127 \pm 1.332	83.351 \pm 0.441	84.457 \pm 0.359
<i>diabetes</i>	70.097 \pm 0.845	71.981 \pm 0.465	66.769 \pm 1.619	72.468 \pm 0.572
<i>german</i>	69.950 \pm 0.475	68.237 \pm 0.417	69.963 \pm 0.613	69.250 \pm 0.409
<i>letterMvsN</i>	97.246 \pm 0.271	94.230 \pm 0.272	96.897 \pm 0.396	97.635 \pm 0.430
<i>letterVvsY</i>	96.734 \pm 0.215	92.121 \pm 0.298	95.815 \pm 0.231	96.863 \pm 0.168
<i>segment</i>	97.543 \pm 0.070	95.952 \pm 0.130	96.326 \pm 0.343	97.657 \pm 0.069
<i>splice</i>	79.207 \pm 0.195	78.374 \pm 0.285	79.017 \pm 0.258	79.672 \pm 0.277
<i>wdbc</i>	96.502 \pm 0.146	95.714 \pm 0.322	92.454 \pm 1.460	96.593 \pm 0.139

The results of $R = |\mathcal{D}_u| \times 5\%$ and $R = |\mathcal{D}_u| \times 10\%$ are listed in Table 4 and 5, respectively. The highest mean accuracy is shown in boldface for each dataset. The t -test results are given in Table 6.

The results of tuned-parameters settings are similar to the results of the above fixed-parameter settings. Compared with UNCERTAIN, ALHS-SVM outperforms UNCERTAIN in 6 out of 8 datasets with regards to the mean accuracy. Moreover, in Table 6, ALHS-SVM reaches significantly better performance than UNCERTAIN in 3 out of 8 datasets when we only query 5% of \mathcal{D}_u and ties with UNCERTAIN in the remaining cases. Compared with representative sampling algorithms, ALHS-SVM also has better performance than REPRESENT and QUIRE with regards to both mean accuracy and results of t -test. The results demonstrate that ALHS-SVM is also effective under the

Table 6: ALHS-SVM versus the other algorithm with tuned parameters based on t -test at 95% significance level

Percentage of queries	Algorithms (win/tie/loss)		
	UNCERTAIN	REPRESENT	QUIRE
5%	3/5/0	3/5/0	6/2/0
10%	1/7/0	7/1/0	6/2/0

tuned parameters in the early stage of active learning. Furthermore, after querying more instances, ALHS-SVM keeps having competitive or better performance than existing uncertainty sampling and representative sampling algorithms when $R = |\mathcal{D}_u| \times 10\%$.

From the study of both fixed and tuned-parameter settings, ALHS-SVM yield better and more stable than other algorithms. The results demonstrate the validity and generality of the proposed ALHS-SVM.

4 Transductive SVM versus HintSVM

As discussed, transductive SVM (TSVM; Joachims, 1999b) has been considered for active learning (Wang et al., 2011) as a representative sampling approach. TSVM arises from semi-supervised learning and has been demonstrated useful for many applications, such text mining (Joachims, 1999b). TSVM aims to maximize its margin between the labeled data and the unlabeled data by assigning suitable labels to the unlabeled data.

The formulation is as follows.

$$\begin{aligned}
\min_{\mathbf{w}, b, \xi, \bar{\xi}, \bar{y}} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C_l \sum_{i=1}^{|\mathcal{D}_l|} \xi_i + C_u \sum_{j=1}^{|\mathcal{D}_u|} \bar{\xi}_j \\
\text{subject to} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \text{for } (\mathbf{x}_i, y_i) \in \mathcal{D}_l, \\
& \bar{y}_j(\mathbf{w}^T \mathbf{x}_j + b) \geq 1 - \bar{\xi}_j \quad \text{for } \mathbf{x}_j \in \mathcal{D}_u, \\
& \xi_i \geq 0 \quad \text{for } (\mathbf{x}_i, y_i) \in \mathcal{D}_l, \\
& \bar{\xi}_j \geq 0 \quad \text{for } \mathbf{x}_j \in \mathcal{D}_u, \\
& \bar{y}_j \in \{+1, -1\} \quad \text{for } \mathbf{x}_j \in \mathcal{D}_u.
\end{aligned} \tag{2}$$

Existing approach (Wang et al., 2011) uses TSVM for querying, which can be viewed as a form of representative sampling that locates the querying boundary by estimating the labels of unlabeled instances.

Comparing formulation (2) of TSVM and formulation (1) of HintSVM as described in Section 3.1, we see that the formulations share some similarities, but focus on very different objective functions. In this section, we study the validity and effectiveness of these two formulations for active learning as querying and/or learning algorithms.

4.1 Comparison between HintSVM and TSVM

We study the case when using HintSVM and TSVM as the querying algorithm while taking the stand-alone SVM as the learning algorithm. The two algorithms are denoted as ALHS-SVM and TSVM-SVM, respectively. For other experimental settings, we follow the same setup as described in Section 3.5.

Figure 3 has presented the accuracy of TSVM-SVM and ALHS-SVM as well as the baseline UNCERTAIN algorithm with SVM. The mean and standard error of accuracy at different rounds of querying are also readily listed in Tables 1 and 2. Clearly, ALHS-

SVM performs generally better than TSVM-SVM across datasets. The results again justify the usefulness of the proposed ALHS framework. We discuss the performance difference as follows.

In formulation (2), TSVM-SVM aims to estimate the possible labels on the unlabeled data, which is similar to QUIRE (Huang et al., 2010). Nevertheless, in QUIRE, the estimation is used for exploring a better query, but TSVM-SVM takes the estimation with a goal of a better classifier. Thus, TSVM-SVM pushes the unlabeled data away from the boundary for better classification ability. Note that ALHS-SVM, on the other hand, aims at a boundary close to parts of unlabeled data (hints) to explore like QUIRE. In the earlier iterations of active learning, exploration rather than pushing the unlabeled data away (as if they are certain) can be important. Thus, TSVM-SVM can be inferior to ALHS-SVM.

On the other hand, in the latter iterations of active learning, ALHS-SVM is similar to the baseline UNCERTAIN approach, which is known to perform decently when the learning algorithm is SVM. Nevertheless, because the boundaries obtained from TSVM and SVM can be quite different, the instances queried by TSVM may not be “uncertain” for learning with SVM. This explanation matches one interesting observation that TSVM-SVM is also worse than the baseline UNCERTAIN algorithm in many datasets such as *australian* and *letter I vs J*. Thus, ALHS-SVM also holds an advantage over TSVM-SVM in the latter iterations.

4.2 Hint Transductive SVM

The discussion in Section 4.1 shows that the discrepancy between TSVM and SVM may be part of the reason that TSVM-SVM is inferior for active learning. What if we take

TSVM for learning instead? Next, we couple TSVM for learning with two querying approaches: HintSVM or TSVM. The two algorithms are named HintSVM-TSVM and TSVM-TSVM, and the results will be shown in Figure 4.

In our study, we observe that HintSVM-TSVM may not result in satisfactory performance in some datasets. We address this problem to the same reason that TSVM-SVM is inferior: HintSVM, which is based on a stand-alone SVM, is very different from TSVM. Thus, uncertain instances queried by HintSVM may not be uncertain to the learning TSVM. Thus, HintSVM-TSVM is yet another combination where the discrepancy between the querying and the learning parts results in unsatisfactory performance. Thus, HintSVM-TSVM results in inferior performance.

The results in Section 4.1 show that considering the learning algorithm is an important issue when designing the querying algorithm to avoid discrepancy. Similar idea has also been studied in other active learning works (Donmez et al., 2007). One interesting question is then whether ALHS can be used when employing TSVM as the learning algorithm. Next, we demonstrate one such possibility. We combine formulation (2) of TSVM and formulation (1) to an extension of TSVM with hint information, which is called Hint Transductive SVM (HintTSVM). HintTSVM can then be used for hinted sampling with TSVM, and expected to be a better match of ALHS when employing

TSVM as the learning algorithm. The formulation of HintTSVM is as follows:

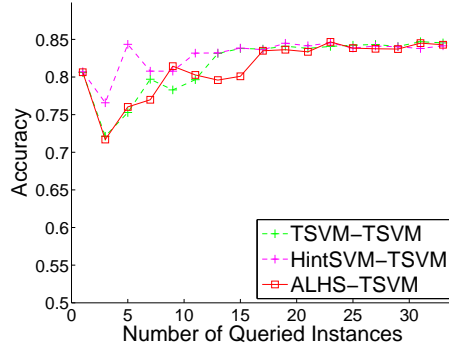
$$\begin{aligned}
& \min_{\mathbf{w}, b, \xi, \bar{\xi}, \tilde{\xi}, \bar{y}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C_l \sum_{i=1}^{|\mathcal{D}_l|} \xi_i + C_u \sum_{j=1}^{|\mathcal{D}'_u|} \bar{\xi}_j + C_h \sum_{j=1}^{|\mathcal{D}_h|} (\tilde{\xi}_j + \tilde{\xi}_j^*) \\
& \text{subject to} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \text{for } (\mathbf{x}_i, y_i) \in \mathcal{D}_l, \\
& \quad \bar{y}_j(\mathbf{w}^T \mathbf{x}_j + b) \geq 1 - \bar{\xi}_j \quad \text{for } \mathbf{x}_j \in \mathcal{D}'_u, \\
& \quad \mathbf{w}^T \tilde{\mathbf{x}}_j + b \leq \epsilon + \tilde{\xi}_j \quad \text{for } \mathbf{x}_j \in \mathcal{D}_h, \\
& \quad -(\mathbf{w}^T \tilde{\mathbf{x}}_j + b) \leq \epsilon + \tilde{\xi}_j^* \quad \text{for } \mathbf{x}_j \in \mathcal{D}_h, \\
& \quad \xi_i \geq 0 \quad \text{for } (\mathbf{x}_i, y_i) \in \mathcal{D}_l, \\
& \quad \bar{\xi}_j \geq 0 \quad \text{for } \mathbf{x}_j \in \mathcal{D}'_u, \\
& \quad \tilde{\xi}_j, \tilde{\xi}_j^* \geq 0 \quad \text{for } \mathbf{x}_j \in \mathcal{D}_h, \\
& \quad \bar{y}_j \in \{+1, -1\} \quad \text{for } \mathbf{x}_j \in \mathcal{D}'_u.
\end{aligned} \tag{3}$$

where $\mathcal{D}_u = \mathcal{D}'_u \cup \mathcal{D}_h$.

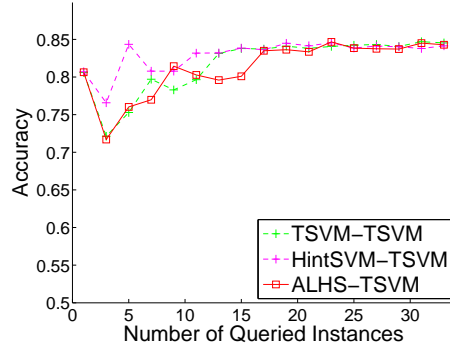
Different from HintSVM, it is difficult to solve HintTSVM efficiently since the TSVM part of the formulation is NP-hard to solve. Therefore, we consider a simple approximation by splitting the training of HintTSVM into two stages. In the first stage, we only consider training the TSVM part on \mathcal{D}_l and \mathcal{D}'_u by any existing algorithm (Joachims, 1999b). In the second stage, we use the inferred labels and cost parameters C_h from the first stage to train a HintSVM as described in Section 3.1. We call the variant of ALHS with HintTSVM as ALHS-TSVM, and list the details in Algorithm 3.

4.3 Experimental Studies of ALHS-TSVM

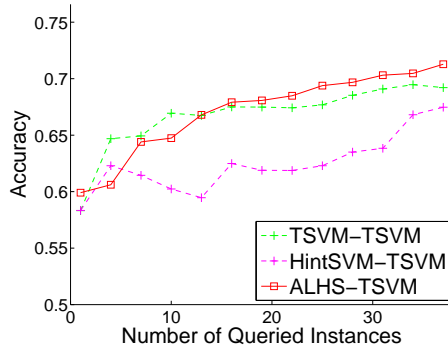
We follow the same experiment setup as previous experiments (fixed parameter setting) and report the results in Figure 4 and Table 7. According to Figure 4, When using



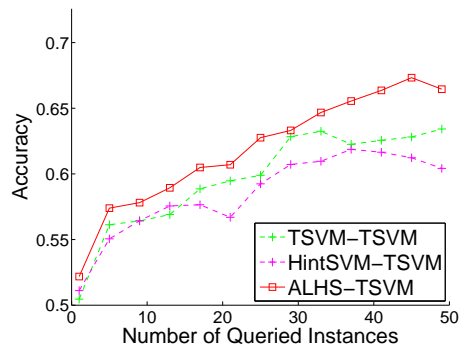
(a) *australian*



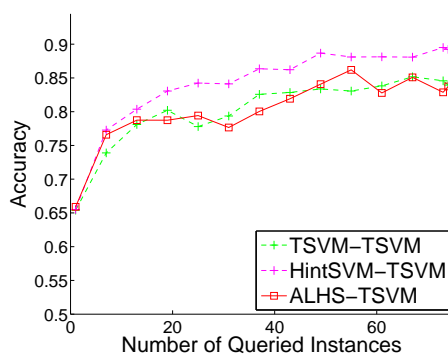
(b) *diabetes*



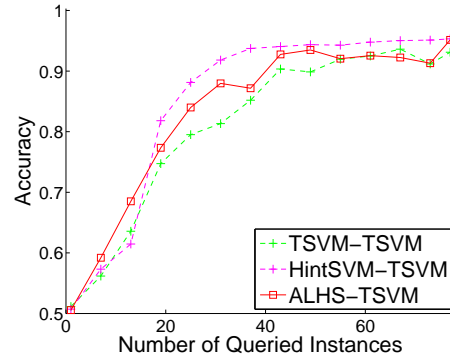
(c) *german*



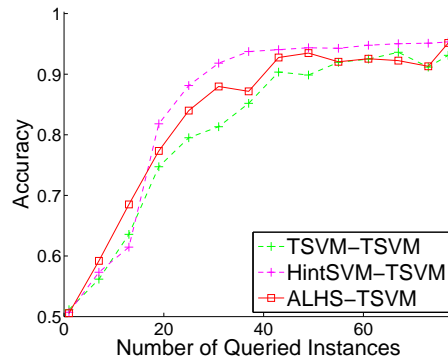
(d) *letterIvsJ*



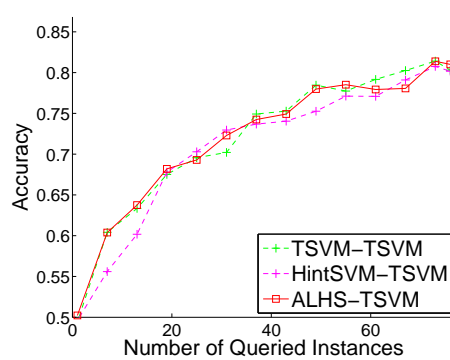
(e) *letterMvsN*



(f) *letterVvsY*



(g) *splice*



(h) *wdbc*

Figure 4: Comparison between TSVM and HintSVM for querying by using TSVM for learning on different datasets

Algorithm 3 The ALHS-TSVM algorithm

Input: the number of rounds R ; a labeled pool \mathcal{D}_l ; an unlabeled pool \mathcal{D}_h ; parameters for HintTSVM and TSVM; ratio δ

Output: decision functions $f^{(1)}, \dots, f^{(R)}$

For $r \leftarrow 1$ **to** R **do**

 Uniformly select \mathcal{D}_h from \mathcal{D}_u

$\mathcal{D}'_u \leftarrow \mathcal{D}_u \setminus \mathcal{D}_h$

$g \leftarrow \text{Train HintTSVM}(C_h, C_u, C_l, \epsilon, \mathcal{D}_h, \mathcal{D}'_u, \mathcal{D}_l)$

$(\tilde{\mathbf{x}}_s, y_s) \leftarrow \text{Query}(g, \mathcal{D}_u)$

$\mathcal{D}_u \leftarrow \mathcal{D}_u \setminus \tilde{\mathbf{x}}_s$; $\mathcal{D}_l \leftarrow \mathcal{D}_l \cup (\tilde{\mathbf{x}}_s, y_s)$

$f^{(r)} \leftarrow \text{Train TSVM}(C_l, C_u, \mathcal{D}_l, \mathcal{D}_u)$

$C_h \leftarrow C_h \times \delta$

End

Table 7: ALHS-TSVM versus the other algorithm based on t -test at 95% significance level

Percentage of queries	Algorithms (win/tie/loss)	
	TSVM-TSVM	HintSVM-TSVM
5%	2/6/0	4/1/3
10%	3/5/0	5/1/3

TSVM (TSVM-TSVM) as the learning algorithm, HintSVM still benefits the querying stage for exploration and results in significantly better performance than TSVM-TSVM in some datasets, such as *letter M vs N* and *splice*. The exploration makes the learning boundary converge to better local optimal one within a few queries. Nevertheless, we observe that HintSVM-TSVM may not result in satisfactory performance in other datasets, such as *diabete*, *german* and *wdbc*. In such cases, ALHS-TSVM results in better or competitive performance than both of them. The results demonstrate the validity

of employing HintTSVM in ALHS-TSVM to explore unknown region of the data to TSVM and resolve the potential drawback of HintSVM-SVM as our previous discussion. It also justifies the advantage of using hint for different learning algorithms.

From the results, we see that the only two datasets that ALHS-TSVM do not perform the strongest are *letter I vs J* and *letter M vs N*. On those datasets, HintSVM-TSVM reaches the best performance. We address this to the difficulty of properly training HintTSVM in ALHS-TSVM, and the simple HintSVM results in more stable performance.

5 Selective Sampling versus Hinted Sampling

As introduced in Section 2, selective sampling (Balcan et al., 2006; Dasgupta et al., 2007; Beygelzimer et al., 2010) along with the pool-to-stream simulation also allows taking both uncertainty and representativeness into account. Next, we compare ALHS-SVM with a state-of-the-art selective sampling algorithm, importance-weighted active learning (IWAL; Beygelzimer et al., 2010) to evaluate the effectiveness of hinted sampling. Upon receiving $\tilde{\mathbf{x}}$ in a simulated data stream in iteration r , IWAL calculates

$$p_r(\tilde{\mathbf{x}}) = \min\left\{1, s(f^{(r)}, \tilde{\mathbf{x}}) \times \frac{\log r}{r}\right\},$$

where s is a scoring function that depends on the goodness of the current hypothesis $f^{(r)}$ and the instance $\tilde{\mathbf{x}}$ in the simulated data stream. In the early iterations, $p_r(\tilde{\mathbf{x}})$ is generally close to 1, and IWAL acts more like representative sampling based on the data distribution in the pool. After several iterations, $p_r(\tilde{\mathbf{x}})$ will be reduced from 1 and IWAL acts more like uncertainty sampling with uncertainty calculated from s . From this perspective, IWAL shares one similarity to ALHS. They both gradually change the

focus from representative sampling to uncertainty sampling.

5.1 Experimental Studies of Selective Sampling

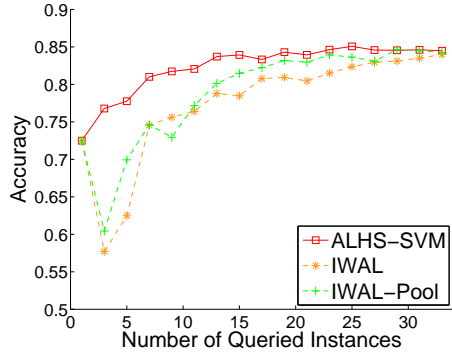
We follow the same setup in Section 3.5. For IWAL, we also use SVM as the learning algorithm and tune the parameter C_0 to be 2 to make IWAL perform better than the original setting with $C_0 = 8$ used by Beygelzimer et al. (2010).

The results are shown in Figure 5. Clearly, IWAL generally performs worse than ALHS-SVM under the pool-based setup. One possible reason for the big difference is that IWAL can only get the information of one instance in each iteration through the pool-to-stream simulation rather than the whole unlabeled pool. The partial information makes it difficult to query useful instances.

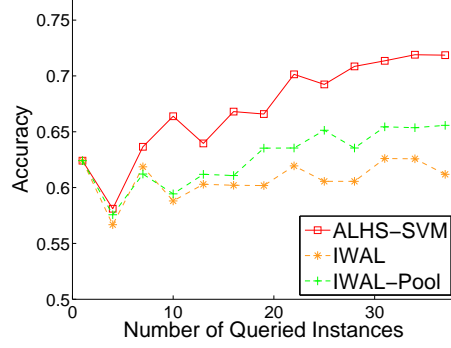
To make a fair comparison to IWAL, we further design another learning algorithm IWAL-Pool. The algorithm estimates $p_r(\tilde{\mathbf{x}})$ of all instances in the unlabeled pool and query the one with the highest probability. The results are also presented in Figure 5. We see that IWAL-Pool is competitive to IWAL and can significantly outperform IWAL in *diabetes* and *letter V vs Y*. Nevertheless, ALHS-SVM, with its simple yet direct use of the unlabelled pool, is generally still better than IWAL-Pool. The results demonstrate the importance of acquiring direct information of unlabeled pool in the pool-based setup. They also highlight the difference of pool-based setup and stream-based setup in active learning.

6 Conclusion

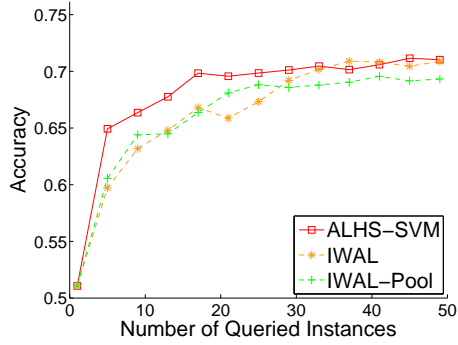
We propose a new framework of active learning, *hinted sampling*, which exploits the unlabeled instances as hints. Hinted sampling can take both uncertainty and represen-



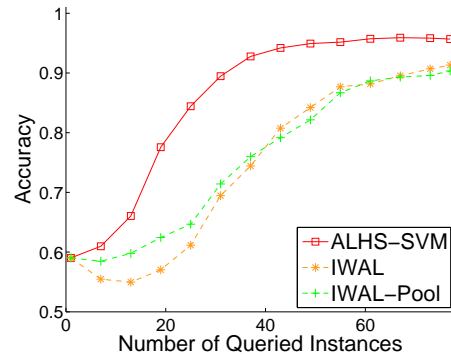
(a) *australian*



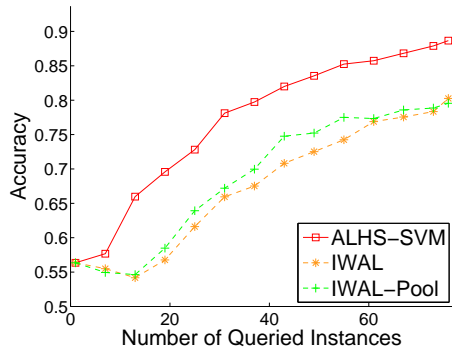
(b) *diabetes*



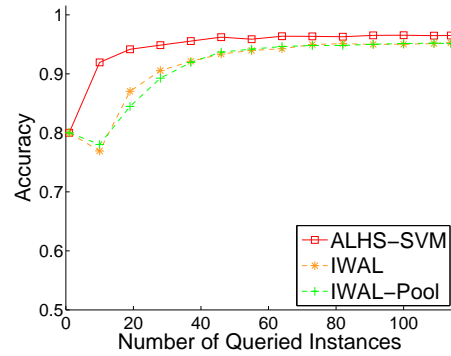
(c) *german*



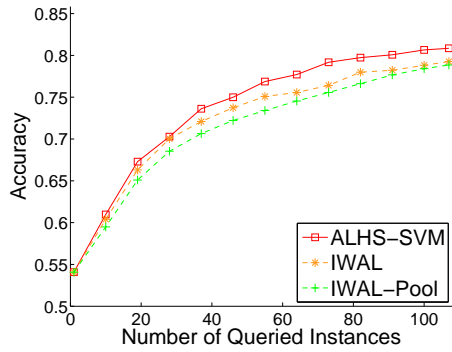
(d) *leterMvsN*



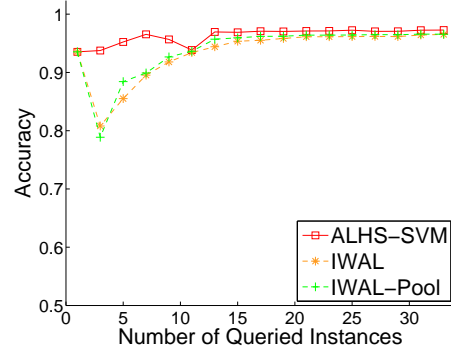
(e) *letterVvsY*



(f) *segment*



(g) *splice*



(h) *wdbc*

Figure 5: Comparison between IWAL and ALHS-SVM on different datasets

tativeness into account concurrently in a more natural and simpler way. We design a novel active learning algorithm ALHS within the framework, and couple the algorithm with a promising hint selection strategy. Because ALHS models the representativeness by hints, it avoids the potential problems of other more sophisticated approaches that are employed by other representative sampling algorithms. Hence, ALHS results in a significantly better and more stable performance than other state-of-the-art algorithms, and can be used to immediately improve SVM-based uncertainty sampling and TSVM-based representative sampling. On the other hand, compared with selective sampling algorithms taking both uncertainty and representativeness into account, the proposed ALHS also has better performance. It not only justifies the effectiveness of ALHS again but also demonstrates the importance of considering the whole unlabeled pool in the querying stage for pool-based active learning problems.

Because of the simplicity and effectiveness of *hinted sampling*, it is worth studying more about this framework. An intensive research direction is to couple hinted sampling with other classification algorithms, and investigate deeper on the hint selection strategies. While we use SVM in ALHS, this framework could be generalized to other classification algorithms. In the future, we plan to investigate more general hint selection strategies and extend hinted sampling from binary classification to other classification problem.

Acknowledgments

A preliminary version of this paper appeared in the Asian Conference on Machine Learning 2012. We thank the reviewers of the conference as well as reviewers for

all versions of this paper for their many useful suggestions. This research has been supported by the National Science Council of Taiwan via NSC 101-2628-E-002-029-MY2..

References

- Abu-Mostafa, Y. S. (1995). Hints. *Neural Computation*, 4:639–671.
- Balcan, M.-F., Beygelzimer, A., and Langford, J. (2006). Agnostic active learning. In *Proceedings of the 23rd International Conference on Machine Learning*.
- Bennett, K. P. and Demiriz, A. (1998). Semi-supervised support vector machines. In *Advances in Neural Information Processing Systems 11*, pages 368–374.
- Beygelzimer, A., Hsu, D., Langford, J., and Zhang, T. (2010). Agnostic active learning without constraints. In *Advances in Neural Information Processing Systems 23*.
- Cao, H., Bao, C., Li, X., and Woon, D. Y.-K. (2014). Class augmented active learning. In *Proceedings of the SIAM International Conference on Data Mining*.
- Chang, C.-C. and Lin, C.-J. (2011). LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, pages 27:1–27:27.
- Chen, T., Chen, Y., Guo, Q., Zhou, Z.-H., Li, L., and Xu, Z. (2014). Effective and efficient microprocessor design space exploration using unlabeled design configurations. *ACM Transactions on Intelligent Systems and Technology*.
- Cohn, D., Atlas, L., and Ladner, R. (1994). Improving generalization with active learning. *Machine Learning*.

- Cohn, D. A., Ghahramani, Z., and Jordan, M. I. (1996). Active learning with statistical models. *Journal of Artificial Intelligence Research*, 4:129–145.
- Dasgupta, S. (2011). Two faces of active learning. *Theoretical Computer Science*.
- Dasgupta, S. and Hsu, D. (2008). Hierarchical sampling for active learning. In *Proceedings of the 25th International Conference on Machine Learning*, pages 208–215.
- Dasgupta, S., Hsu, D., and Monteleoni, C. (2007). A general agnostic active learning algorithm. In *Advances in Neural Information Processing Systems 20*.
- Donmez, P., Carbonell, J. G., and Bennett, P. N. (2007). Dual strategy active learning. In *Proceedings of the 18th European Conference on Machine Learning*, pages 116–127.
- Frank, A. and Asuncion, A. (2010). UCI machine learning repository.
- Guo, Y. and Greiner, R. (2007). Optimistic active learning using mutual information. In *Proceedings of the 20th International Joint Conference on Artificial Intelligence*, pages 823–829.
- Hanneke, S. (2007). Teaching dimension and the complexity of active learning. In *Proceedings of the Annual Conference on Learning Theory*.
- Hoi, S. C. H., Jin, R., Zhu, J., and Lyu, M. R. (2008). Semi-supervised SVM batch mode active learning for image retrieval. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 1–7.
- Huang, S.-J., Jin, R., and Zhou, Z.-H. (2010). Active learning by querying informative and representative examples. In *Advances in Neural Information Processing Systems 23*, pages 892–900.

- Joachims, T. (1999a). Advances in kernel methods. chapter Making large-scale support vector machine learning practical.
- Joachims, T. (1999b). Transductive inference for text classification using support vector machines. In *Proceedings of the Sixteenth International Conference on Machine Learning*.
- Langford, J. and Zhang, T. (2007). The epoch-greedy algorithm for contextual multi-armed bandits. In *Advances in Neural Information Processing Systems 20*.
- Lewis, D. D. and Gale, W. A. (1994). A sequential algorithm for training text classifiers. In *Proceedings of the 17th ACM International Conference on Research and Development in Information Retrieval*, pages 3–12.
- Li, C.-L., Ferng, C.-S., and Lin, H.-T. (2012). Active learning with hinted support vector machine. In *Proceedings of the forth Asian Conference on Machine Learning*.
- Melville, P. and Mooney, R. J. (2004). Diverse ensembles for active learning. In *Proceedings of the 21st International Conference on Machine Learning*, pages 584–591.
- Nguyen, H. T. and Smeulders, A. (2004). Active learning using pre-clustering. In *Proceedings of the 21st International Conference on Machine Learning*, pages 623–630.
- Quanquan Gu, Tong Zhang, J. H. (2014). Batch-mode active learning via error bound minimization. In *Proceedings of the 30th Conference Annual Conference on Uncertainty in Artificial Intelligence*.

- Ratsch, G., Onoda, T., and Müller, K. R. (2001). Soft margins for AdaBoost. *Machine Learning*, 2:27:1–27:27.
- Settles, B. (2009). Active learning literature survey. Technical report, University of Wisconsin–Madison.
- Settles, B. and Craven, M. (2008). An analysis of active learning strategies for sequence labeling tasks. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- Tong, S. and Koller, D. (2000). Support vector machine active learning with applications to text classification. In *Proceedings of the 17th International Conference on Machine Learning*, pages 999–1006.
- Vapnik, V. (1998). *Statistical learning theory*. Wiley.
- Wang, Z., Yan, S., and Zhang, C. (2011). Active learning with adaptive regularization. *Pattern Recognition*.
- Wang, Z. and Ye, J. (2013). Querying discriminative and representative samples for batch mode active learning. In *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
- Xu, Z., Yu, K., Tresp, V., Xu, X., and Wang, J. (2003). Representative sampling for text classification using support vector machines. In *Proceedings of the 25th European Conference on Information Retrieval Research*, pages 393–407.
- Zha, Z.-J., Wang, M., Zheng, Y.-T., Yang, Y., Hong, R., and Chua, T.-S. (2012). Interactive video indexing with statistical active learning. *IEEE Transactions on Multimedia*.