
Detecting Concept Drift with Support Vector Machines

Ralf Klinkenberg
Thorsten Joachims

KLINKENBERG@LS8.CS.UNI-DORTMUND.DE
JOACHIMS@LS8.CS.UNI-DORTMUND.DE

Informatik VIII, Universität Dortmund, Baroper Str. 301, 44221 Dortmund, Germany
<http://www-ai.cs.uni-dortmund.de/>

Abstract

For many learning tasks where data is collected over an extended period of time, its underlying distribution is likely to change. A typical example is information filtering, i.e. the adaptive classification of documents with respect to a particular user interest. Both the interest of the user and the document content change over time. A filtering system should be able to adapt to such concept changes. This paper proposes a new method to recognize and handle concept changes with support vector machines. The method maintains a window on the training data. The key idea is to automatically adjust the window size so that the estimated generalization error is minimized. The new approach is both theoretically well-founded as well as effective and efficient in practice. Since it does not require complicated parameterization, it is simpler to use and more robust than comparable heuristics. Experiments with simulated concept drift scenarios based on real-world text data compare the new method with other window management approaches. We show that it can effectively select an appropriate window size in a robust way.

1. Introduction

Machine learning methods are often applied to problems, where data is collected over an extended period of time. In many real-world applications this introduces the problem that the distribution underlying the data is likely to change over time. For example, companies collect an increasing amount of data like sales figures and customer data to find patterns in the customer behaviour and to predict future sales. As the customer behaviour tends to change over time, the model underlying successful predictions should be adapted accordingly.

The same problem occurs in information filtering, i.e. the adaptive classification of documents with respect to a particular user interest. Information filtering techniques are used, for example, to build personalized news filters, which learn about the news-reading preferences of a user or to filter e-mail. Both the interest of the user and the document content change over time. A filtering system should be able to adapt to such concept changes.

This paper proposes a new method for detecting and handling concept changes with support vector machines. The approach has a clear theoretical motivation and does not require complicated parameter tuning. After reviewing other work on adaptation to changing concepts and shortly describing support vector machines, this paper explains the new window adjustment approach and evaluates it in three simulated concept drift scenarios on real-world text data. The experiments show that the approach effectively selects an appropriate window size and results in a low predictive error rate.

2. Concept Drift

Throughout this paper, we study the problem of concept drift for the pattern recognition problem in the following framework. Each example $\vec{z} = (\vec{x}, y)$ consists of a feature vector $\vec{x} \in \mathbf{R}^N$ and a label $y \in \{-1, +1\}$ indicating its classification. Data arrives over time in batches. Without loss of generality these batches are assumed to be of equal size, each containing m examples.

$\vec{z}_{(1,1)}, \dots, \vec{z}_{(1,m)}, \vec{z}_{(2,1)}, \dots, \vec{z}_{(2,m)}, \dots, \vec{z}_{(t,1)}, \dots, \vec{z}_{(t,m)}, \vec{z}_{(t+1,1)}, \dots, \vec{z}_{(t+1,m)}$

$\vec{z}_{(ij)}$ denotes the j -th example of batch i . For each batch i the data is independently identically distributed with respect to a distribution $\text{Pr}_i(\vec{x}, y)$. Depending on the amount and type of concept drift, the example distribution $\text{Pr}_i(\vec{x}, y)$ and $\text{Pr}_{i+1}(\vec{x}, y)$ between batches will differ. The goal of the learner \mathcal{L} is to sequentially predict the labels of the next batch. For example, after

batch t the learner can use any subset of the training examples from batches 1 to t to predict the labels of batch $t + 1$. The learner aims to minimize the cumulated number of prediction errors.

In machine learning, changing concepts are often handled by time windows of fixed or adaptive size on the training data (Mitchell et al., 1994; Widmer & Kubat, 1996; Lanquillon, 1997; Klinkenberg & Renz, 1998) or by weighting data or parts of the hypothesis according to their age and/or utility for the classification task (Kunisch, 1996; Taylor et al., 1997). The latter approach of weighting examples has already been used for information filtering in the incremental relevance feedback approaches of Allan (1996) and Balabanovic (1997). In this paper, the earlier approach maintaining a window of adaptive size is explored. More detailed descriptions of the methods described above and further approaches can be found in Klinkenberg (1998).

For windows of fixed size, the choice of a “good” window size is a compromise between fast adaptivity (small window) and good generalization in phases without concept change (large window). The basic idea of *adaptive window management* is to adjust the window size to the current extent of concept drift.

The task of learning drifting or time-varying concepts has also been studied in computational learning theory. Learning a changing concept is infeasible, if no restrictions are imposed on the type of admissible concept changes,¹ but drifting concepts are provably efficiently learnable (at least for certain concept classes), if the rate or the extent of drift is limited in particular ways.

Helmbold, & Long (1994) assume a possibly permanent but slow concept drift and define the *extent of drift* as the probability that two subsequent concepts disagree on a randomly drawn example. Their results include an upper bound for the extend of drift maximally tolerable by any learner and algorithms that can learn concepts that do not drift more than a certain constant extent of drift. Furthermore they show that it is sufficient for a learner to see a fixed number of the most recent examples. Hence a window of a certain minimal fixed size allows to learn concepts for which the extent of drift is appropriately limited.

¹E.g. a function randomly jumping between the values one and zero cannot be predicted by any learner with more than 50% accuracy.

While Helmbold and Long restrict the extend of drift, Kuh et al. (1991) determine a maximal *rate of drift* that is acceptable by any learner, i. e. a maximally acceptable frequency of concept changes, which implies a lower bound for the size of a fixed window for a time-varying concept to be learnable, which is similar to the lower bound of Helmbold and Long.

In practice, however, it usually cannot be guaranteed that the application at hand obeys these restrictions, e.g. a reader of electronic news may change his interests (almost) arbitrarily often and radically. Furthermore the large time window sizes, for which the theoretical results hold, would be impractical. Hence more application oriented approaches rely on far smaller windows of fixed size or on window adjustment heuristics that allow far smaller window sizes and usually perform better than fixed and/or larger windows (Widmer & Kubat, 1996; Lanquillon, 1997; Klinkenberg & Renz, 1998). While these heuristics are intuitive and work well in their particular application domain, they usually require tuning their parameters, are often not transferable to other domains, and lack a proper theoretical foundation.

Syed et al. (1999) describe an approach to incrementally learning support vector machines that handles *virtual* concept drift implied by incrementally learning from several subsamples of a large training set, but they do not address the problem of (*real*) concept drift addressed here.

3. Support Vector Machines

The window adjustment approach described in this paper uses support vector machines (Vapnik, 1998) as their core learning algorithm. Support vector machines are based on the *structural risk minimization* principle (Vapnik, 1998) from statistical learning theory. In their basic form, SVMs learn linear decision rules

$$h(\vec{x}) = \text{sign}\{\vec{w} \cdot \vec{x} + b\} = \begin{cases} +1, & \text{if } \vec{w} \cdot \vec{x} + b > 0 \\ -1, & \text{else} \end{cases} \quad (1)$$

described by a weight vector \vec{w} and a threshold b . The idea of structural risk minimization is to find a hypothesis h for which one can guarantee the lowest probability of error. For SVMs, Vapnik (1998) shows that this goal can be translated into finding the hyperplane with maximum soft-margin.² Computing this hyperplane is equivalent to solving the following optimization problem.

²See Burges (1998) for an introduction to SVMs.

Optimization Problem 1 (SVM (primal))

$$\text{minimize: } V(\vec{w}, b, \vec{\xi}) = \frac{1}{2} \vec{w} \cdot \vec{w} + C \sum_{i=1}^n \xi_i \quad (2)$$

$$\text{subject to: } \forall_{i=1}^n : y_i [\vec{w} \cdot \vec{x}_i + b] \geq 1 - \xi_i \quad (3)$$

$$\forall_{i=1}^n : \xi_i > 0 \quad (4)$$

In this optimization problem, the Euclidean length $\|\vec{w}\|$ of the weight vector is inversely proportional to the soft-margin of the decision rule. The constraints (3) require that all training examples are classified correctly up to some slack ξ_i . If a training example lies on the “wrong” side of the hyperplane, the corresponding ξ_i is greater or equal to 1. Therefore $\sum_{i=1}^n \xi_i$ is an upper bound on the number of training errors. The factor C in (2) is a parameter that allows trading-off training error vs. model complexity.

For computational reasons it is useful to solve the Wolfe dual (Fletcher, 1987) of optimization problem 1 instead of solving optimization problem 1 directly (Vapnik, 1998).

Optimization Problem 2 (SVM (dual))

$$\text{minimize: } W(\vec{\alpha}) = -\sum_{i=1}^n \alpha_i + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j (\vec{x}_i \cdot \vec{x}_j) \quad (5)$$

$$\text{subject to: } \sum_{i=1}^n y_i \alpha_i = 0 \quad (6)$$

$$\forall_{i=1}^n : 0 \leq \alpha_i \leq C \quad (7)$$

In this paper, *SVM^{light}* (Joachims, 1999) is used for computing the solution of this optimization problem.³ Support vectors are those training examples \vec{x}_i with $\alpha_i > 0$ at the solution. From the solution of optimization problem 2 the decision rule can be computed as

$$\vec{w} \cdot \vec{x} = \sum_{i=1}^n \alpha_i y_i (\vec{x}_i \cdot \vec{x}) \quad \text{and} \quad b = y_{usv} - \vec{w} \cdot \vec{x}_{usv} \quad (8)$$

The training example (\vec{x}_{usv}, y_{usv}) for calculating b must be a support vector with $\alpha_{usv} < C$. Finally, the training losses ξ_i can be computed as $\xi_i = \max(1 - y_i [\vec{w} \cdot \vec{x}_i + b], 0)$.

For both solving optimization problem 2 as well as applying the learned decision rule, it is sufficient to be able to calculate inner products between feature vectors. Exploiting this property, Boser et al. introduced the use of kernels $K(\vec{x}_1, \vec{x}_2)$ for learning non-linear decision rules. Depending on the type of kernel function, SVMs learn polynomial classifiers, radial basis

³*SVM^{Light}* is available at http://www-ai.informatik.uni-dortmund.de/svm_light

function (RBF) classifiers, or two layer sigmoid neural nets. Such kernels calculate an inner-product in some feature space and replace the inner-product in the formulas above.

4. Window Adjustment by Optimizing Performance

Our approach to handling drift in the distribution of examples uses a window on the training data. This window should include only those example which are sufficiently “close” to the current target concept. Assuming the amount of drift increases with time, the window includes the last n training examples. Previous approaches used similar windowing strategies. Their shortcomings are that they either fix the window size (Mitchell et al., 1994) or involve complicated heuristics (Widmer & Kubat, 1996; Lanquillon, 1997; Klinkenberg & Renz, 1998). A fixed window size makes strong assumptions about how quickly the concept changes. While heuristics can adapt to different speed and amount of drift, they involve many parameters that are difficult to tune. Here, we present an approach to selecting an appropriate window size that does not involve complicated parameterization. Their key idea is to select the window size so that the estimated generalization error on new examples is minimized. To get an estimate of the generalization error we use a special form of $\xi\alpha$ -estimates (Joachims, 2000). $\xi\alpha$ -estimates are a particularly efficient method for estimating the performance of a SVM.

4.1 $\xi\alpha$ -Estimators

$\xi\alpha$ -estimators are based on the idea of leave-one-out estimation (Lunts & Brailovskiy, 1967). The leave-one-out estimator of the error rate proceeds as follows. From the training sample $S = ((\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n))$ the first example (\vec{x}_1, y_1) is removed. The resulting sample $S^{\setminus 1} = ((\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n))$ is used for training, leading to a classification rule $h_{\mathcal{L}}^{\setminus 1}$. This classification rule is tested on the held out example (\vec{x}_1, y_1) . If the example is classified incorrectly it is said to produce a leave-one-out error. This process is repeated for all training examples. The number of leave-one-out errors divided by n is the leave-one-out estimate of the generalization error.

While the leave-one-out estimate is usually very accurate, it is very expensive to compute. With a training sample of size n , one must run the learner n times. $\xi\alpha$ -estimators overcome this problem using an upper bound on the number of leave-one-out errors instead of calculating them brute force. They owe their name

to the two arguments they are computed from. $\vec{\xi}$ is the vector of training losses at the solution of the primal SVM training problem. $\vec{\alpha}$ is the solution of the dual SVM training problem. Based on these two vectors — both are available after training the SVM at no extra cost — the $\xi\alpha$ -estimators are defined using the following two counts. With R_{Δ}^2 being the maximum difference of any two elements of the Hessian (i.e. $R_{\Delta}^2 \geq \max_{\vec{x}, \vec{x}'} (\mathcal{K}(\vec{x}, \vec{x}) - \mathcal{K}(\vec{x}, \vec{x}'))$),

$$d = |\{i : (\alpha_i R_{\Delta}^2 + \xi_i) \geq 1\}| \quad (9)$$

counts the number of training examples, for which the quantity $\alpha_i R_{\Delta}^2 + \xi_i$ exceeds one. Since the document vectors are normalized to unit length in the experiments described in this paper, here $R_{\Delta}^2 = 1$. It is proven in Joachims (2000) that d is an approximate upper bound on the number of leave-one-out errors in the training set. With n as the total number of training examples, the $\xi\alpha$ -estimators of the error rate is

$$Err_{\xi\alpha}^n(h_{\mathcal{L}}) = \frac{|\{i : (\alpha_i R_{\Delta}^2 + \xi_i) \geq 1\}|}{n} \quad (10)$$

The theoretical properties of this $\xi\alpha$ -estimator are discussed in Joachims (2000). It can be shown that the estimator is pessimistically biased, overestimating the true error rate on average. Experiments show that the bias is acceptably small for text classification problems and that the variance of the $\xi\alpha$ -estimator is essentially as low as that of a holdout estimate using twice as much data. It is also possible to design similar estimators for precision and recall, as well as combined measures like $F1$ (Joachims, 2000).

4.2 Window Adjustment Algorithm

A window adjustment algorithm has to solve the following trade-off. A large window provides the learner with much training data, allowing it to generalize well given that the concept did not change. On the other hand, a large window can contain old data that is no longer relevant (or even confusing) for the current target concept. Finding the right size means trading-off the quality against the number of training examples.

To answer this question the window adjustment algorithm proposed in the following uses $\xi\alpha$ -estimates in a particular way. At batch t , it essentially tries various window sizes, training a SVM for each resulting training set.

$$\vec{z}_{(t,1)}, \dots, \vec{z}_{(t,m)} \quad (11)$$

$$\vec{z}_{(t-1,1)}, \dots, \vec{z}_{(t-1,m)}, \vec{z}_{(t,1)}, \dots, \vec{z}_{(t,m)} \quad (12)$$

$$\vec{z}_{(t-2,1)}, \dots, \vec{z}_{(t-2,m)}, \vec{z}_{(t-1,1)}, \dots, \vec{z}_{(t-1,m)}, \vec{z}_{(t,1)}, \dots, \vec{z}_{(t,m)} \quad (13)$$

⋮

For each window size it computes a $\xi\alpha$ -estimate based on the result of training. In contrast to the previous section, the $\xi\alpha$ -estimator used here considers only the last batch, that is the m most recent training examples $\vec{z}_{(t,1)}, \dots, \vec{z}_{(t,m)}$.

$$Err_{\xi\alpha}^m(h_{\mathcal{L}}) = \frac{|\{i : 1 \leq i \leq m \wedge (\alpha_{(t,i)} R_{\Delta}^2 + \xi_{(t,i)}) \geq 1\}|}{m} \quad (14)$$

This reflects the assumption that the most recent examples are most similar to the new examples in batch $t + 1$. The window size minimizing the $\xi\alpha$ -estimate of the error rate is selected by the algorithm.

The algorithm can be summarized as follows:

- **input:** S training sample consisting of t batches containing m examples each
- **for** $h \in \{0, \dots, t - 1\}$
 - **train SVM** on examples $\vec{z}_{(t-h,1)}, \dots, \vec{z}_{(t-h,m)}$
 - **compute $\xi\alpha$ -estimate** on examples $\vec{z}_{(t,1)}, \dots, \vec{z}_{(t,m)}$
- **output:** window size which minimizes $\xi\alpha$ -estimate

5. Experiments

5.1 Experimental Setup

Each of the following *data management approaches* is evaluated in combination with the SVM:

- *“Full Memory”*: The learner generates its classification model from all previously seen examples, i.e. it cannot “forget” old examples.
- *“No Memory”*: The learner always induces its hypothesis only from the most recent batch. This corresponds to using a window of the fixed size of one batch.
- Window of *“Fixed Size”*: A window of the fixed size of three batches is used.
- Window of *“Adaptive Size”*: The window adjustment algorithm proposed in the previous section adapts the window size to the current concept drift situation.

The experiments are performed in an information filtering domain, a typical application area for learning drifting concept. Text documents are represented as attribute-value vectors (*bag of words* model), where each distinct word corresponds to a feature whose

Table 1. Relevance of the categories in the concept change scenarios A, B, and C.

Scenario	Category	Probability of being relevant for a document of the specified category at the specified time step (batch)																			
		0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
A	1	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
B	1	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.8	0.6	0.4	0.2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.2	0.4	0.6	0.8	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
C	1	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.0	0.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
	3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

value is the “l₁”-TF/IDF-weight (Salton & Buckley, 1988) of that word in that document. Words occurring less than three times in the training data or occurring in a given list of stop words are not considered. Each document feature vector is normalized to unit length to abstract from different document lengths.

The performance of a classifier is measured by the three metrics prediction error, recall, and precision. *Recall* is the probability, that the classifier recognizes a relevant document as relevant. *Precision* is the probability, that a document classified as relevant actually is relevant. All reported results are estimates averaged over ten runs.

The experiments use a subset of 2608 documents of the data set of the *Text REtrieval Conference (TREC)* consisting of English business news texts. Each text is assigned to one or several categories. The categories considered here are 1 (Antitrust Cases Pending), 3 (Joint Ventures), 4 (Debt Rescheduling), 5 (Dumping Charges), and 6 (Third World Debt Relief). For the experiments, three concept change scenarios are simulated. The texts are randomly split into 20 batches of equal size containing 130 documents each.⁴ The texts of each category are distributed as equally as possible over the 20 batches.

Table 1 describes the relevance of the categories in the three concept change scenarios A, B, and C. For each time step (batch), the probability of being relevant (interesting to the user) is specified for documents of categories 1 and 3, respectively. Documents of the classes 4, 5, and 6 are never relevant in any of these scenarios. In the first scenario (*scenario A*), first documents of category 1 are considered relevant for the user interest and all other documents irrelevant. This changes abruptly (concept shift) in batch 10, where documents of category 3 are relevant and all others irrelevant. In the second scenario (*scenario B*), again

⁴Hence, in each trial, out of the 2608 documents, eight randomly selected texts are not considered.

first documents of category 1 are considered relevant for the user interest and all other documents irrelevant. This changes slowly (concept drift) from batch 8 to batch 12, where documents of category 3 are relevant and all others irrelevant. The third scenario (*scenario C*) simulates an abrupt concept shift in the user interest from category 1 to category 3 in batch 9 and back to category 1 in batch 11.

5.2 Results

Figure 1 compares the prediction error rates of the adaptive window size algorithm with the non-adaptive methods. The graphs show the prediction error on the following batch. In all three scenarios, the full memory strategy and the adaptive window size algorithm essentially coincide as long as there is no concept drift. During this stable phase, both show lower prediction error than the fixed size and the no memory approach. At the point of concept drift, the performance of all methods deteriorates. While the performance of no memory and adaptive size recovers quickly after the concept drift, the error rate full memory approach remains high especially in scenarios A and B. Like before the concept drift, the no memory and the fixed size strategies exhibit higher error rates than the adaptive window algorithm in the stable phase after the concept drift. This shows that the no memory, the fixed size, and the full memory approaches all perform suboptimally in some situation. Only the adaptive window size algorithm can achieve a relatively low error rate over all phases in all scenarios. This is also reflected in the average error rates over all batches given in Table 2. The adaptive window size algorithm achieves a low average error rate on all three scenarios. Similarly, precision and recall are consistently high.

The behavior of the adaptive window algorithm is best explained by looking at the window sizes it selects. Figure 2 shows the average training window ranges. The bottom of each graph depicts the time and extent of concept drift in the corresponding scenario. For

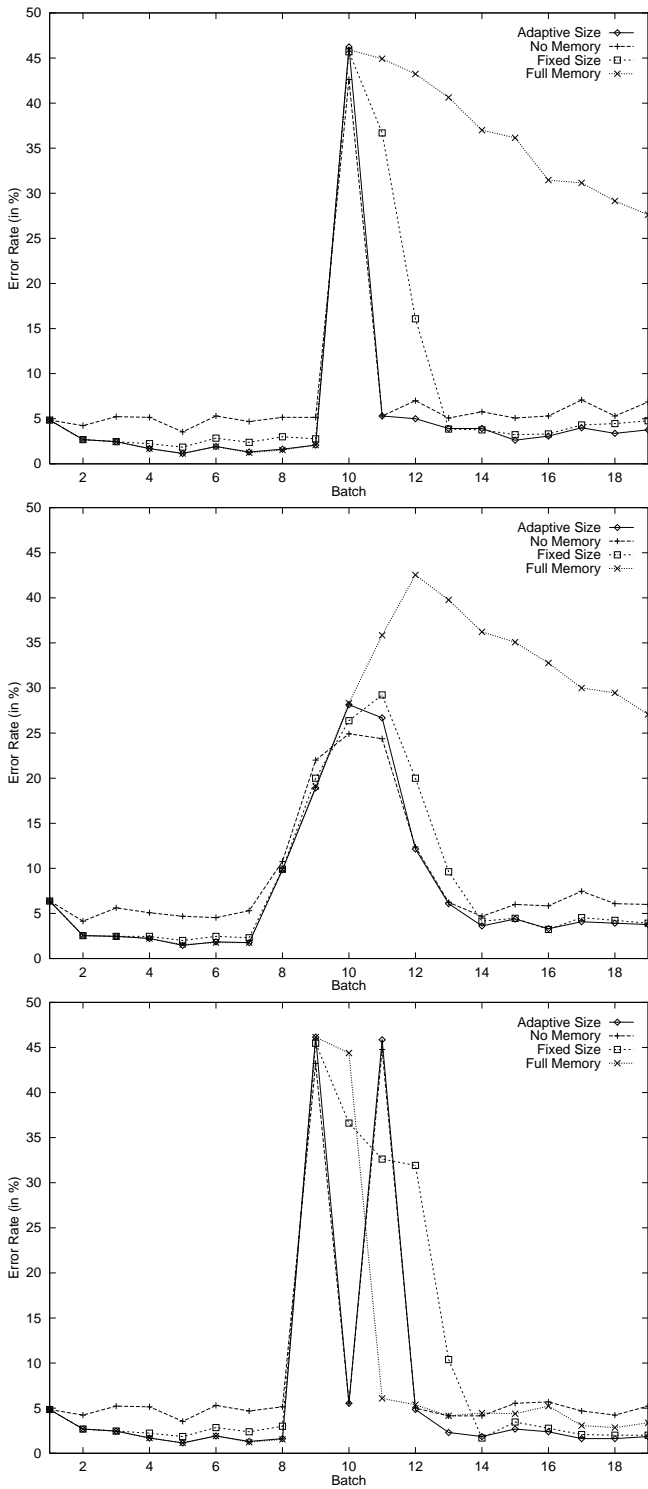


Figure 1. Comparison of the prediction error rates for scenario A (top), B (middle), and C (bottom). The x-axis denotes the batch number and the y-axis the average prediction error.

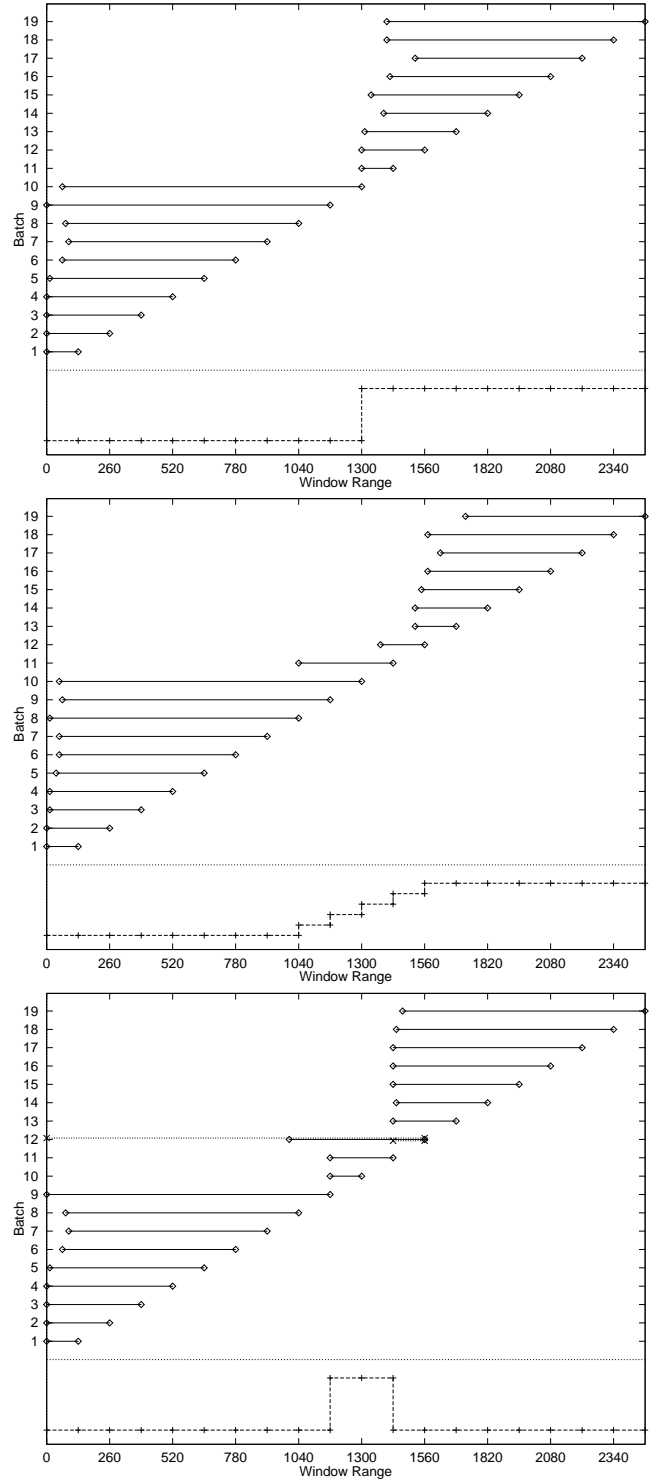


Figure 2. Window size and range for scenario A (top), B (middle), and C (bottom). The y-axis denotes the batch number. Each horizontal line marks the average training window range selected at that batch number. The bottom part of each graph depicts the location and type of the concept shift.

Table 2. Error, accuracy, recall, and precision of all window management approaches for all scenarios averaged over 10 trials with 20 batches each (standard sample error in parentheses).

	Full Memory	No Memory	Fixed Size	Adaptive Size
Scenario A:				
Error	20.36% (4.21%)	7.30% (1.97%)	7.96% (2.80%)	5.32% (2.29%)
Recall	51.69% (8.37%)	74.42% (4.61%)	77.64% (6.07%)	85.35% (4.93%)
Precision	64.67% (8.38%)	91.29% (5.10%)	87.73% (5.93%)	91.61% (5.11%)
Scenario B:				
Error	20.25% (3.56%)	9.08% (1.57%)	8.44% (2.00%)	7.56% (1.89%)
Recall	49.35% (7.01%)	67.22% (5.04%)	73.85% (5.51%)	76.70% (5.42%)
Precision	65.09% (6.80%)	88.86% (3.67%)	87.19% (4.18%)	88.48% (3.89%)
Scenario C:				
Error	7.74% (3.05%)	8.97% (2.84%)	10.17% (3.30%)	7.07% (3.16%)
Recall	76.54% (6.26%)	63.68% (5.27%)	68.18% (7.05%)	78.17% (6.34%)
Precision	83.15% (6.69%)	87.67% (7.06%)	79.00% (8.09%)	87.38% (6.99%)

scenario A the training window increases up to the abrupt concept change after batch 10, covering almost all examples available for the current concept. Only in batches 5 to 10 the average training set size is slightly smaller than maximally possible. Our explanation is that for large training sets a relatively small number of additional examples does not always make a “noticeable” difference. After the concept change in batch 10 the adaptive window size algorithm now picks training windows covering only those examples from after the drift as desired. A similar behavior is found for scenario B (Figure 2, middle). Since the drift is less abrupt, the adaptive window size algorithm intermediately selects training examples from both concepts in batch 11. After sufficiently many training examples from the new distribution are available, those earlier examples are discarded. The behavior of the adaptive window size algorithm in scenario C is reasonable as well (Figure 2, bottom). A particular situation occurs in batch 12. Here the window size exhibits a large variance. For 8 of the 10 runs the algorithm selects a small training set size of one batch, while for the remaining 2 runs it selects all available training examples starting with batch 1. Here there appears to be a borderline decision between accepting 2 (out of 12) batches of “bad” examples or just training on a single batch.

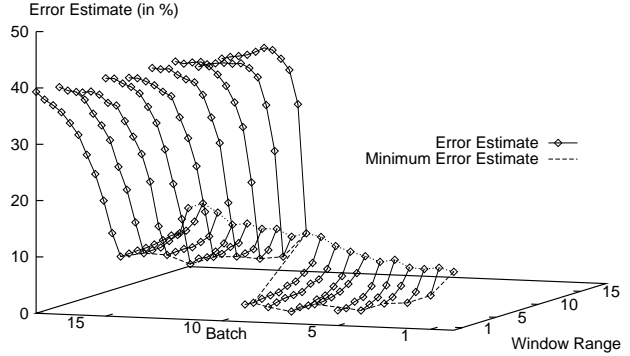


Figure 3. Average $\xi\alpha$ -estimates at different batches and for varying training window sizes for scenario A. The dashed curve marks the beginning of the window with the lowest error estimate.

Further insight on how the algorithm selects the window size is gained from figure 3. It plots the average $\xi\alpha$ -estimate in scenario A over all batches and for varying window size. The x_1 -axis denotes the number of the current batch (increasing from right to left) and the x_2 -axis the batch of the window start. The dashed line indicates the beginning of the window with the lowest estimate in the batch. The graph shows that the error estimate decreases with growing window size in batches 1 to 10. After batch 10, the estimate accurately reflects the concept change. The error estimate decreases with training windows growing towards the abrupt concept change. If the window is enlarged beyond this change, the estimated error increases steeply as expected.

6. Summary and Conclusions

In this paper, we proposed a new method for handling concept drift with support vector machines. The method directly implements the goal of discarding irrelevant data with the aim of minimizing generalization error. Exploiting the special properties of SVMs, we adapted $\xi\alpha$ -estimates to the window size selection problem. Unlike for the conventional heuristic approaches, this gives the new method a clear and simple theoretical motivation. Furthermore, the new method is easier to use in practical applications, since it involves less parameters than complicated heuristics. Experiments in an information filtering domain show that the new algorithm achieves a low error rate and selects appropriate window sizes over very different concept drift scenarios.

An open question is how sensitive the algorithm is to the size of individual batches. Since in the current version of the algorithm the batch size determines the estimation window, the variance of the window size is

likely to increase with smaller batches. It might be necessary to select the estimation window size independent of the batch size. A shortcoming of most existing algorithms handling concept drift (an exception is Lanquillon (1999)) is that they can detect concept drift only after labeled data is available. That is, after the learning algorithm starts making mistakes. While this appears unavoidable for concept drift with respect to $\Pr(y|\vec{x})$, it might be possible to detect concept drift in $\Pr(\vec{x})$ earlier by using transductive support vector machines.

Acknowledgments

This work was supported by the DFG Collaborative Research Center on Complexity Reduction in Multivariate Data (SFB 475) and by the DFG Collaborative Research Center on Computational Intelligence (SFB 531).

References

- Allan, J. (1996). Incremental relevance feedback for information filtering. *Proceedings of the Nineteenth ACM Conference on Research and Development in Information Retrieval* (pp. 270–278). New York: ACM Press.
- Balabanovic, M. (1997). An adaptive web page recommendation service. *Proceedings of the First International Conference on Autonomous Agents* (pp. 378–385). New York: ACM Press.
- Burges, C. (1998). A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery*, 2, 121–167.
- Fletcher, R. (1987). *Practical methods of optimization* (2nd edition). New York: Wiley.
- Helmbold, D. P., & Long, P. M. (1994). Tracking drifting concepts by minimizing disagreements. *Machine Learning*, 14, 27–45.
- Joachims, T. (1999). Making large-scale SVM learning practical. In B. Schölkopf, C. Burges, & A. Smola (Eds.), *Advances in kernel methods - Support vector learning*. Cambridge, MA, USA: MIT Press.
- Joachims, T. (2000). Estimating the generalization performance of a SVM efficiently. *Proceedings of the Seventeenth International Conference on Machine Learning*. San Francisco: Morgan Kaufman.
- Klinkenberg, R. (1998). *Maschinelle Lernverfahren zum adaptiven Informationsfiltern bei sich verändernden Konzepten*. Masters thesis, Fachbereich Informatik, Universität Dortmund, Germany.
- Klinkenberg, R., & Renz, I. (1998). Adaptive information filtering: Learning in the presence of concept drifts. *Workshop Notes of the ICML-98 Workshop on Learning for Text Categorization* (pp. 33–40). Menlo Park, CA, USA: AAAI Press.
- Kuh, A., Petsche, T., & Rivest, R. (1991). Learning time-varying concepts. *Advances in Neural Information Processing Systems* (pp. 183–189). San Mateo, CA, USA: Morgan Kaufmann.
- Kunisch, G. (1996). *Anpassung und Evaluierung statistischer Lernverfahren zur Behandlung dynamischer Aspekte in Data Mining*. Masters thesis, Fachbereich Informatik, Universität Ulm, Germany.
- Lanquillon, C. (1997). *Dynamic Neural Classification*. Masters thesis, Fachbereich Informatik, Universität Braunschweig, Germany.
- Lanquillon, C. (1999). Information filtering in changing domains. *Working Notes of the IJCAI-99 Workshop on Machine Learning for Information Filtering* (pp. 41–48). Stockholm, Sweden.
- Lunts, A., & Brailovskiy, V. (1967). Evaluation of attributes obtained in statistical decision rules. *Engineering Cybernetics*, 3, 98–109.
- Mitchell, T., Caruana, R., Freitag, D., McDermott, J., & Zabowski, D. (1994). Experience with a learning personal assistant. *Communications of the ACM*, 37, 81–91.
- Salton, G., & Buckley, C. (1988). Term weighting approaches in automatic text retrieval. *Information Processing and Management*, 24, 513–523.
- Syed, N. A., Liu, H., & Sung, K. K. (1999). Handling concept drifts in incremental learning with support vector machines. *Proceedings of the Fifth International Conference on Knowledge Discovery and Data Mining*. New York: ACM Press.
- Taylor, C., Nakhaeizadeh, G., & Lanquillon, C. (1997). Structural change and classification. *Workshop Notes of the ECML-97 Workshop on Dynamically Changing Domains: Theory Revision and Context Dependence Issues* (pp. 67–78).
- Vapnik, V. (1998). *Statistical learning theory*. Chichester, GB: Wiley.
- Widmer, G., & Kubat, M. (1996). Learning in the presence of concept drift and hidden contexts. *Machine Learning*, 23, 69–101.