

Towards Scalable Support Vector Machines using Squashing

Dmitry Pavlov

Dept. of Info. and Comp. Science
Univ. of California, Irvine
Irvine, CA 92697-3425, USA
pavlovd@ics.uci.edu

Darya Chudova

Dept. of Info. and Comp. Science
Univ. of California, Irvine
Irvine, CA 92697-3425, USA
dchudova@ics.uci.edu

Padhraic Smyth

Dept. of Info. and Comp. Science
Univ. of California, Irvine
Irvine, CA 92697-3425, USA
smyth@ics.uci.edu

March 3, 2000

Abstract

Support vector machines (SVMs) provide classification models with strong theoretical foundations as well as excellent empirical classification performance on a variety of applications. One of the major drawbacks of SVMs is the necessity to solve a large-scale quadratic programming problem. This paper combines likelihood-based squashing with a probabilistic formulation of SVMs, enabling their fast training on squashed data sets. We reduce the problem of training the SVMs on the weighted “squashed” data to a quadratic programming problem that can be solved using Platt’s sequential minimal optimization (SMO) algorithm. We compare performance of the SMO algorithm on the squashed and the full data, as well as simple random and boosted samples of the data. Experiments on a number of datasets show that squashing allows one to speed-up training, decrease memory requirements, and obtain parameter estimates close to that of the full data. More importantly, squashing produces classification accuracy that is close to optimal.

1 Introduction

Following the pioneering work of Vladimir Vapnik [Vap98], support vector machines (SVMs) are steadily gaining popularity in the machine learning community. SVMs have been proven to exhibit several attractive theoretical properties, including maximum margin separation

between the classes. In addition, SVMs have been empirically shown to outperform conventional classifiers on a variety of benchmarks [SBS99]. However, the applicability of SVMs to large datasets is limited because of the high computational cost involved in solving quadratic programming problem arising in their training.

Various training algorithms have been proposed to speed up the training, including chunking [SBS99], Osuna’s decomposition method [OFG97a], and Platt’s Sequential Minimal Optimization (SMO) [Pla99]. Although these algorithms accelerate the training, they do not scale well with the size of the training data.

Another approach to reducing the computational cost is to use approximation methods. The simplest method would be just to sample from the original dataset and use the sample to train the SVM. An extension of this idea was explored by Pavlov et. al. [PMD] who reported results on application of boosting to training SVMs. Boost-SMO algorithm trains a sequence of SVM classifiers on the samples of data so that each subsequent classifier concentrates mostly on the errors made by the previous ones [Sch99, PMD]. While this method is not optimal in general, it allows for a very fast training of SVMs, has substantially lower memory cost and yields performance close to that of the full SMO.

In this paper we suggest yet another method for scaling SVMs up based on *squashing*. DuMouchel et. al. [DVJ⁺99] and Madigan et. al. [MRD⁺] recently introduced squashing as a technique that allows one to scale the data down while preserving its statistical properties. In particular, likelihood-based squashing [MRD⁺] assumes a particular statistical model and tries to preserve the behaviour of the likelihood function of the original data (referred to as the “mother-data”) in the neighborhood of the maximum likelihood solution. A recent paper by Owen [Owe99] analyzes why and when squashing may reduce the complexity of learning parameters of the model and concludes that a method that uses global (as opposed to local) features of the data will be likely to benefit from squashing. Since support vectors and the margin are determined by the complete dataset we can expect linear SVMs to benefit from squashing.

To apply likelihood-based squashing to SVMs it is necessary to have their probabilistic interpretation. Several papers have explored the correspondence between SVMs and gaussian processes (see, e.g., [JH, Wah97]). In particular, [Sol99] shows how to interpret the training procedure as a problem of finding maximum a posteriori values for the parameters of the SVM.

We show that the probabilistic interpretation of SVM training in conjunction with the likelihood-based squashing allow one to scale the training up. We reduce the problem of training the SVMs on the weighted “squashed” data to a quadratic programming problem that can be solved using SMO algorithm. We compare performance of the SMO algorithm

on the squashed and the full data, as well as simple random and boosted samples of the data. Experiments on a number of datasets show that squashing allows one to speed-up training, decrease memory requirements, and obtain parameter estimates close to that of the full data. More importantly, squashing produces classification accuracy that is close to optimal. Note that although our results can be directly generalized to the large-scale datasets that do not fit in the main memory, we are using moderately sized data that resides in memory in our experiments (and we still use the term “scaling up” for this data).

The rest of the paper is organized as follows. We give a brief overview of SVMs in section 2. We then describe squashing and its application to SVMs in more detail in section 3. Section 4 describes results of empirical evaluation of the algorithms on 4 datasets. In section 5 we draw conclusions and discuss possible directions for the future work.

2 SVMs: A Brief Overview

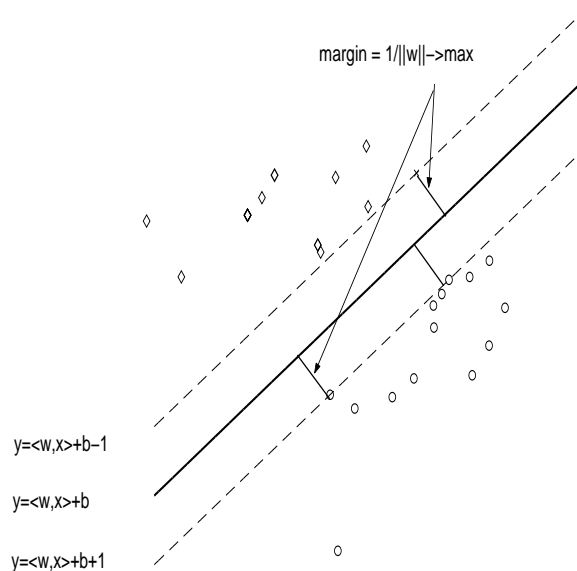


Figure 1: A linear support vector machine in 2 dimensional linearly separable case. From the infinite set of possible separating lines the one yielding maximum margin should be picked.

We give a brief overview of linear SVMs here. Suppose that we have a total of N labeled patterns in the n -dimensional space $D = \{(\mathbf{x}_i, y_i)\}$ belonging to two classes, so that label y_i is either -1 or 1 (see figure 1 for a 2-dimensional illustration). In linear SVMs the line separating the classes is sought in the form

$$y = \langle \mathbf{w}, \mathbf{x} \rangle + b, \tag{1}$$

where \mathbf{w} is the normal vector and b is the intercept of the hyperplane. The problem of finding the maximum margin separation between the classes is equivalent to the minimization of the $\|\mathbf{w}\|_{L_2}^2$ with the additional constraints $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1$ to ensure that all the patterns in the training set are classified correctly [CV97]. However, in most cases perfect separation is impossible and we need to trade errors on the individual training patterns for the maximum margin. The optimization problem becomes

$$\min_{\mathbf{w}, b} E = \min_{\mathbf{w}, b} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \right\}, \quad (2)$$

subject to

$$y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i, \quad (3)$$

where C is the trading constant and ξ_i are the non-negative errors. If ξ_i is equal to 0, the i -th example is correctly classified and doesn't make any contribution to the error in the Equation 2. Thus, learning parameters of the SVM is reduced to a linearly constrained quadratic programming problem that can be converted to a dual formulation using Lagrangian conditional optimization techniques. The dual problem is defined on N variables α_i and has the form

$$\max_{\alpha} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \right\} \quad (4)$$

subject to

$$0 \leq \alpha_i \leq C, \quad \text{and} \quad \sum_{i=1}^n \alpha_i y_i = 0, \quad (5)$$

where α_i is the Lagrange multiplier corresponding to the i^{th} pattern. The standard SMO algorithm solves the problem defined in 4 and 5 by decomposing it into the smaller problems that can be solved analytically. SMO is provably optimal [OFG97b]. Once the dual problem is solved it is possible to reconstruct the vector \mathbf{w} and the constant b from the set of Lagrange multipliers α_i using analytical formulas that we omit for brevity.

Several authors (e.g., [Wah97, Sol99]) have shown that the right side of equation 2 can be treated as a log-posterior on the parameters \mathbf{w} and b , with the first term corresponding to a prior on \mathbf{w} and the second one to the likelihood. The prior on \mathbf{w} will have a normal distribution with zero mean:

$$p(\mathbf{w}) \sim \exp(-\|\mathbf{w}\|^2), \quad (6)$$

while the log-likelihood of the data will (up to a constant additive factor independent of \mathbf{w} and b) be proportional to:

$$l^f(\mathbf{w}, b) \propto \sum_{i=1}^N I(1 - y_i[\langle \mathbf{w}, \mathbf{x}_i \rangle + b]), \quad (7)$$

where $I(z)$ is z times the indicator function of z .

By solving the dual optimization problem via SMO one essentially finds the maximum a posteriori values for \mathbf{w} and b . We will use this probabilistic interpretation of the training procedure to develop a *squashing* procedure for SVMs.

3 Squashing for SVMs

To illustrate the idea of squashing let us assume that we have selected a probabilistic model $p(z|\theta)$ for the data D and our objective is to find the maximum likelihood estimate θ_{ML} for the vector of parameters θ . Assuming independence of observations, the parameters θ_{ML} can be found in the following manner:

$$\theta_{ML} = \arg \max_{\theta} \sum_{i=1}^N \log p((\mathbf{x}_i, y_i)|\theta). \quad (8)$$

Further, suppose that the terms in equation 8 can be grouped (or clustered) so that the likelihoods of the individual points within each group do not vary significantly in the neighborhood of θ_{ML} . Squashing procedure eliminates points that contribute similar amounts to the likelihood by replacing each cluster with a single point placed into its center of mass and a weight equal to the number of points in the cluster. The maximum likelihood estimate can now be found from expression 9 which approximates equation 8:

$$\theta_{ML} \simeq \arg \max_{\theta} \sum_{c=1}^{N_c} \beta_c \log p((\mathbf{x}_c, y_c)^{sq}|\theta). \quad (9)$$

Here N_c is the number of clusters, β_c is the weight (or equivalently the number of points in the original cluster c) and $(\mathbf{x}_c, y_c)^{sq}$ is the squashed data point placed at the center of the cluster c . Note that we typically look for at least a 50-100 times reduction in the amount of data so that conventional algorithms can be applied to the squashed data directly.

Notice that clustering at this step should not be computationally more intensive than solving the original maximum likelihood problem on the mother-data. In the formulation given in [MRD⁺] squashing makes only two passes through the mother-data. It evaluates the likelihoods of the training points for certain choices of θ and forms the so-called *likelihood profiles* in the first pass. After this N_c points are randomly selected to form initial centers of clusters. Finally, on the second pass each training point is assigned to the cluster whose center's likelihood profile is closest to the likelihood profile of the data point.

Important design issues for the squashing algorithm include the choice of the number and location of the points in the parameter space to evaluate the likelihoods at, and the number of squashed data points to ensure a sufficiently good approximation. Paper [MRD⁺] suggests to

use various *factorial designs* [BHH78] in the parameter space. While this method is universal, as we noted above, for SVMs we can sample the values of \mathbf{w} from the prior distribution in equation 6. The choice of the intercept term b can be made from the condition that the hyperplane goes through the cloud of training points. We evaluate the likelihood of each training point for the fixed pair of \mathbf{w} and b and repeat the selection procedure L times, where L is the length of likelihood profile. Our experiments show that for SVMs the performance of the squash-SMO method is almost insensitive to the choice of L as long as L is at least an order of 10^2 .

Suppose that we have performed squashing of the original data set using the likelihood defined in equation 7. We now need to classify the squashed dataset $D^{sq} = \{(\mathbf{x}_c, y_c)^{sq}\}$ containing N_c weighted points. An expression for the likelihood of the squashed dataset (up to a constant additive factor independent of \mathbf{w} and b) is proportional to

$$l^{sq}(\mathbf{w}, b) \propto \sum_{c=1}^{N_c} \beta_c I(1 - y_c^{sq}[\langle \mathbf{w}, \mathbf{x}_c^{sq} \rangle + b]), \quad (10)$$

since the squashed dataset contains β_c points coinciding with $(\mathbf{x}_c, y_c)^{sq}$. Going back to the non-probabilistic formulation we will have the following optimization problem for the squashed data:

$$\min_{\mathbf{w}, b} E^{sq} = \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C^{sq} \sum_{i=1}^{N_c} \beta_i \xi_i, \quad (11)$$

with the constraints

$$y_i^{sq}(\langle \mathbf{w}, \mathbf{x}_i^{sq} \rangle + b) \geq 1 - \xi_i, \quad (12)$$

where the slack variables ξ_i are non-negative as before. It is easy to see that the substitution of slack variables $\eta_i = \beta_i \xi_i$ transforms the objective function for the squashed problem into exactly the same form as we had for the mother-data, while the inequalities for constraints will have slack variables η_i divided by the corresponding weights. The latter observation results in only one difference between the dual problems for the mother- and the squashed data: the allowed regions for the Lagrange multipliers in the squashed data will be bounded by $C^{sq} \beta_i$ as opposed to C in the mother-data. Thus, once the SMO code is updated to account for new bounds on Lagrange multipliers, it can be used directly on the weighed data from squashing.

4 Experiments

We have run experiments on four datasets, one of which is synthetic and the rest are publicly available at either UCI machine learning [BM98] or UCI KDD repository [Bay99]. We

evaluated the performance of the full-SMO (SMO on the full training data), srs-SMO (SMO on a simple random sample), squash-SMO (SMO on the squashed data) and boost-SMO (SMO on the boosted samples) with respect to the misclassification rate, the time to learn the parameters of the model and the memory used during the training.

We performed all experiments on Sun Solaris 2.7 UNIX workstation with two 168MHz Sparc CPUs and 256Mb of main memory. We averaged results for squash-SMO, boost-SMO and srs-SMO over 100 runs to account for variability in the data resulting from sampling. Our implementation of the SMO algorithm used caching [Pla99] to store frequently used dot products.

4.1 Results on Synthetic Data

Synthetic data was obtained by sampling from the specially designed piece-wise gaussian distributions in 2D. Each distribution is composed of two parts as follows. Let $g_1(x, y)$ and $g_2(x, y)$ be two gaussian distributions, with the equal means $\mu = (\mu_x, \mu_y) = (-0.75, 0)$ and covariance matrices

$$\Sigma_1 = \begin{bmatrix} 1 & 0 \\ 0 & 9 \end{bmatrix},$$

and

$$\Sigma_2 = \begin{bmatrix} 3 & 0 \\ 0 & 9 \end{bmatrix},$$

respectively. Let the probability density for one of the classes be $g_1(x, y)$ for $x < \mu_x$ and $g_2(x, y)$ for $x > \mu_x$ so that there is a discontinuity at $x = \mu_x$. Another distribution is a mirror symmetrical image of the first one relative to the line $x = 0$. For clarity we have placed the plot of the cross-section $y = 0$ of these distributions in the figure 2.

We generated 5000 examples of each class for both training and test sets.

Recall that we seek the boundary in the form $w_1x + w_2y + b = 0$, here $\mathbf{w} = (w_1, w_2)$ is a slope vector and b is an intercept. Table 1 illustrates how dramatically different was the performance of squash-SMO from srs-SMO on this dataset. In the first column we report the average relative error in parameter estimates over 100 runs, where the error is defined as follows:

$$E^2 = \frac{\|\mathbf{w}_f - \mathbf{w}_s\|^2 + (b_f - b_s)^2}{\|w_f\|^2 + b_f^2}. \quad (13)$$

Here (\mathbf{w}_f, b_f) are the parameters estimated by the SMO algorithm on the full data and (\mathbf{w}_s, b_s) are obtained from either the squashed or sampled data. Squash-SMO obtains parameter estimates that are almost 4 times more accurate than srs-SMO. The next 3 columns contain the standard deviation of the individual estimates for w_1 , w_2 and b in 100 runs of

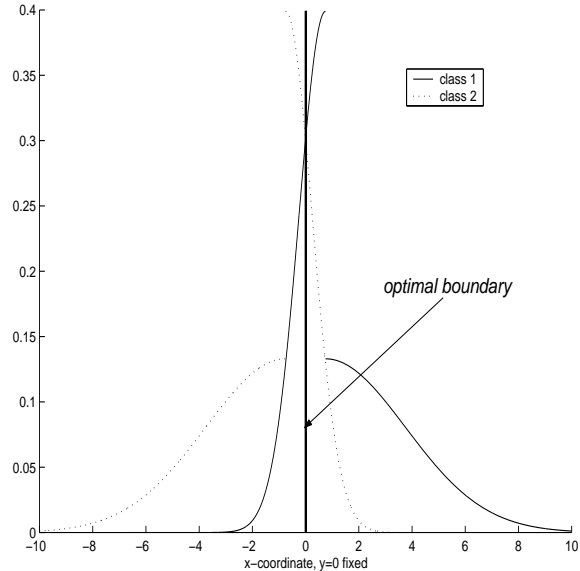


Figure 2: Cross-section of the Two Gaussian Distributions Generating the Data for Synthetic Dataset.

Table 1: Average Relative Error and Standard Deviation in the SVM Parameter Estimates for Squash-SMO and Srs-SMO on the Synthetic Example. See Equation 13 for the definition of E .

	E	σ_{w_1}	σ_{w_2}	σ_b
Srs-SMO	0.467	0.52	0.47	0.38
Squash-SMO	0.124	0.1	0.09	0.08

srs-SMO and squash-SMO. Notice that the parameters estimated by the squash-SMO are 5 times more stable than estimates obtained by srs-SMO.

Note that the two distributions have a much higher density near the optimal boundary $x = 0$ separating them. We analyzed the structure of the errors made by different models and concluded that srs-SMO made most errors in the neighborhood of the optimal decision boundary. Since Squash-SMO was able to obtain more stable and accurate parameter estimates, it exhibited a much better performance in that region.

Table 2 shows the error on the test set for various models on the synthetic data. Clearly, srs-SMO provides the worst improvement over the baseline model, full-SMO is the best with an error close to the Bayes error rate, and squash-SMO is a close runner-up.

Table 2: Error on the Test Set for Various Models On the Synthetic Data.

Baseline	1% srs-SMO	1% squash-SMO	1% boost-SMO	full-SMO
50%	41%	34.04%	34.84%	33.54%

4.2 Results on Benchmark Data

The public datasets were “The Microsoft Anonymous Web” and the “Forest Cover Type” datasets available at UCI KDD archive [Bay99] and Adult dataset available at UCI machine learning repository [BM98]. *Web data* reflects the Web pages of www.microsoft.com that each user visited during one week of February 1998. The classification task for this dataset, as we pose it, is to predict whether a user will visit the most popular site S based on his/her visiting pattern of all other sites. This dataset contains binary data and average record has only 3% of its attributes set to 1. For *Adult data* the task is to predict if the income of a person is greater than 50K based on several census parameters, such as age, education, marital status and so forth. Finally, for the Forest data the task is to distinguish between several forest cover types based on the information about the location of the place, type of soil, distance to the water etc. For this dataset we experimented with 2 most populated classes for which the misclassification rate is the greatest: Spruce-Fir and Lodgepole Pine. The parameters of all datasets are summarized in the Table 3.

Table 3: Characteristics of the Data Sets.

Parameter	Web	Adult	Forest
$N_{Attributes}$	294	14	54
Attributes	Binary	Categ. and Cont.	Categ. and Cont.
Training set	10057(+)	7508(+)	56404(+)
	21875(-)	22654(-)	42480(-)
Test set	1561(+)	3700(+)	226897(+)
	3325(-)	11360(-)	169360(-)

Note that one of the distinct advantages of the methods that use sampled or squashed data is that they allow to do cross-validation to optimize parameters of the SVM algorithm. Since m -fold cross-validation will increase the computational cost of training by a factor of m , using it for the full-SMO is out of question. In our experiments we used 5-fold cross-validation to adjust the value of the regularization constant C for the squash-SMO, srs-SMO and boost-SMO. For the full-SMO C was set to 1 in all runs.

We analyzed the dependence of the SVM misclassification rate on the length of the likelihood profile and the number N_c of squashed points. In Figure 3 the parameter P refers to the squashing ratio, i.e. the number of points in the squashed data relative to the mother-data. As P increases from 0.5% to 7.5% we see a gradual decrease in the misclassification rate so that it approaches the misclassification rate of SVMs on the full data. It is also apparent that for this data (Adult dataset) likelihood profiles of size 50-100 will yield the best rate. Similar results were obtained on the Web data. In all experiments described below

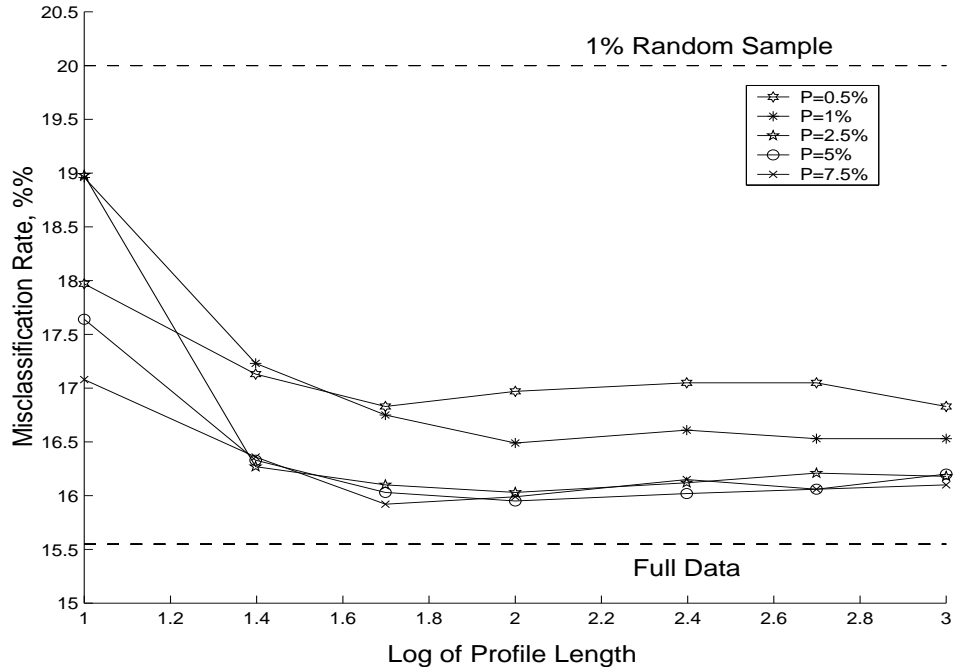


Figure 3: Dependence of the SVM Misclassification Rate on the Length of Likelihood Profile for Various Squashing Ratios P for the Adult data.

we were using a fixed length of the likelihood profile set to 100.

Table 4: Baseline Test Set Error and % Decrease from Baseline Test Set Error for Various Models and Datasets.

	Web	Adult	Forest
Baseline Error	0.32	0.25	0.43
by using 1% srs-SMO	10%	18.7%	40%
1% squash-SMO	20%	31.8%	46.7%
1% boost-SMO	20.5%	33.2%	45.6%
full-SMO	22.5%	36.8%	48.1%

Table 4 shows the misclassification rate on the test set for various datasets and algorithms. The first row, baseline error, shows the misclassification rate of the trivial prediction that classifies all examples into the most populated class. The rest of the rows show what percentage of decrease in error relative to baseline a particular model is able to provide. Notice that srs-SMO algorithm is the worst and full-SMO algorithm is the best. The difference in misclassification rate between the two extremes ranges from 8% to 18%. It is remarkable that performance of both squash-SMO and boost-SMO is much closer to full-SMO than performance of srs-SMO.

Table 5 reports timing results for various algorithms. Srs-SMO turns out to be the fastest

Table 5: CPU Time to Learn Parameters of Various Models for Various Datasets. For squash-SMO (srs-SMO) the first number is the time to do squashing (sampling) and the second is the time to do SMO on the squashed (sampled) data.

	Web	Adult	Forest
1% srs-SMO	1.5 + 0.2	1.5 + 0.2	13 + 2
1% squash-SMO	361.7 + 0.2	132.4 + 0.2	682.9 + 2
1% boost-SMO	10	7.5	85.3
full-SMO	4018	1151	1225

with times to do SMO below 2 seconds. The runner up is boost-SMO with times ranging from 7.5 for the synthetic problem to 85 seconds for the forest data. For both srs-SMO and squash-SMO the time to perform sampling or squashing dominates over the time to do SMO and it takes much longer to do squashing than sampling. The full-SMO algorithm is the slowest, being on average 10-100 times slower than the rest of the algorithms. While the optimality is by all means a big advantage of the full-SMO, its time and memory requirements leave much to be desired even for the modestly sized problems that we considered here. As an illustration, in our experiments we only allowed to caching of up to 9 million dot products. When this upper limit was reached cache was flushed. For all non-synthetic datasets this upper limit was reached and the program overall took from 80 to 120Mb of memory during training (the training and test data resided in main memory as well). Note that in SMO we could trade memory for speed if we did not allow to use cache to store dot products, but in this case the training will be even slower.

For the really large datasets squash-SMO might outperform boost-SMO since the latter needs a) to store the current boosted sample and the distribution over all training examples and b) multiple paths through the dataset to update the distribution on examples and sample from it. The Squash-SMO on the other hand will only need 2 full paths over the data as we discussed before during the squashing phase and the mother-data will not be used at all during SVM training.

Yet another advantage of squash-SMO over boost-SMO is a better interpretability of the model that might be essential for exploratory data analysis and visualization.

5 Conclusions and Future work

We describe how the use of squashing makes the training of SVMs applicable to large datasets. Comparison with the optimal full-SMO algorithm shows that both squash-SMO and boost-SMO are able to reach a near-optimal performance with much lower time and

memory requirements. On the other hand, the simplest and fastest standard random sampling SMO on average has a considerably higher misclassification rate. We conclude that statistical techniques based on the ideas of sampling (boost-SMO) and likelihood preservation (squash-SMO) promise high computational advantages for large datasets over the standard SMO algorithm.

Both squash- and boost-SMO allow one to use cross-validation to tune parameters of the SVM, such as regularization constant C or the width of the kernel in non-linear SVMs. Such tuning is usually impossible on the full data due to the high computational cost.

Although the performance of boost-SMO and squash-SMO is similar on the benchmark problems we outlined when and why one algorithm should be preferred over another. In particular, squash-SMO offers a better interpretability of the model and can be expected to run faster than SMO on the datasets that do not reside in the main memory. For the datasets that fit in the main memory we saw that boost-SMO is faster than squash-SMO.

We see several possibilities to extend this work. We think that gains offered by squash-SMO and boost-SMO might be even more significant for non-linear SVMs that take much longer to train.

Several authors have studied the question of feature selection for SVMs (see, e.g. [BO98]). The main idea is to interpret the norm of w as a norm in L_1 space and reduce the training task to linear programming optimization problem. We think that squashing might offer a principled way to do feature selection for the mother-data. More precisely, one might first squash the data, solve the linear programming problem on the squashed data, figure out what features are redundant and then do feature selection on the mother data.

As we outlined above squashing itself takes considerable time. For the high-dimensional data one will need to increase the profile length to get acceptable estimates for parameters of the model. This will result in the time increase and performance degradation. Thus, it will be useful to study the conditions which would allow to apply squashing to the high dimensional data.

An interesting theoretical question is whether it is possible to give a bound on the distance between the parameter estimates obtained by the full-SMO and squash-SMO.

6 Acknowledgements

We would like to thank Jianchang Mao for making his SMO code available to us. We would also like to thank Jock Blackard for his help with the Forest Cover Data.

References

- [Bay99] S.D. Bay. UCI kdd archive, 1999.
- [BHH78] G.E.P. Box, W.G. Hunter, and J.S. Hunter. *Statistics for Experimenters: An Introduction to Design, Data Analysis, and Model Building*. John Wiley & Sons, New York, 1978.
- [BM98] C.L. Blake and C.J. Merz. UCI repository of machine learning databases, 1998.
- [BO98] P.S. Bradley and O.L. Mangasarian. Feature selection via concave minimization and support vector machines. In *Machine Learning Proceedings of the Fifteenth International Conference (ICML 1998)*, pages 82–90, 1998.
- [CV97] C. Cortes and V. Vapnik. Support-vector networks. *Machine Learning*, pages 273–297, 1997.
- [DVJ+99] W. DuMouchel, C. Volinsky, T. Johnson, C. Cortes, and D. Pregibon. Squashing flat files flatter. *Proceedings of the KDD-99*, 1999.
- [JH] T. S. Jaakola and D. Haussler. Probabilistic kernel regression methods; <http://www.ai.mit.edu/~tommi/publications/probker.ps.gz>.
- [MRD⁺] D. Madigan, N. Raghavan, W. DuMouchel, M. Nason, C. Posse, and G. Ridgeway. Likelihood-based data squashing: A modeling approach to instance construction; <http://www.stat.washington.edu/greg/docs/lds.ps>.
- [OFG97a] E. Osuna, R. Freund, and F. Girosi. An improved training algorithm for support vector machines. In J. Principe, L. Gile, N. Morgan, and E. Wilson, editors, *Proc. of IEEE Workshop on Neural Networks for Signal Processing*, pages 276–285. IEEE Press, New York, 1997.
- [OFG97b] E. Osuna, R. Freund, and F. Girosi. Training support vector machines: An application to face detection. In *Proc. of CVPR 1997*. 1997.
- [Owe99] A. Owen. Data squashing by empirical likelihood; <http://www-stat.stanford.edu/~owen/reports/>. 1999.
- [Pla99] J. Platt. Fast training of support vector machines using sequential minimal optimization. In B. Scholkopf, C.J.C. Burges, and A.J. Smola, editors, *Advances in Kernel Methods – Support Vector Learning*, pages 185–208. MIT Press, Cambridge, MA, 1999.

- [PMD] D. Pavlov, J. Mao, and B. Dom. Scaling-up support vector machines using boosting algorithm. *Submitted to ICPR 2000*.
- [SBS99] B. Scholkopf, C.J.C. Burges, and A.J. Smola (eds). *Advances in Kernel Methods – Support Vector Learning*. MIT Press, Cambridge, MA, 1999.
- [Sch99] R. E. Schapire. An introduction to boosting algorithm. In *Proc. of the sixteenth Intl. Joint Conf. on Artificial Intelligence*, 1999.
- [Sol99] P. Sollich. Probabilistic interpretation and bayesian methods for support vector machines. In *Proceedings of ICANN 1999*, pages 91–96, 1999.
- [Vap98] V. N. Vapnik. *Statistical Learning Theory*. John Wiley & Sons, New York, 1998.
- [Wah97] G. Wahba. Support vector machines, reproducing kernel hilbert spaces and the randomized gacv; <ftp://ftp.stat.wisc.edu/pub/wahba/nips97.ps.gz>. University of Wisconsin Statistics Dept TR 984, 1997.