Distributed Optimization of Multi-Class SVMs

Maximilian Alber 1 † Julian Zimmert 3 † Urun Dogan 2 Marius Kloft 3

¹ Berlin Big Data Center Berlin Institute of Technology 10587 Berlin, Germany Microsoft Research
 Cambridge CB1 2FB
 United Kingdom

Department of Computer Science Humboldt University of Berlin 10099 Berlin, Germany

Abstract

Training of one-vs.-rest SVMs can be parallelized over the number of classes in a straight forward way. Given enough computational resources, one-vs.-rest SVMs can thus be trained on data involving a large number of classes. The same cannot be stated, however, for the so-called all-in-one SVMs, which require solving a quadratic program of size quadratically in the number of classes. We develop distributed algorithms for two all-inone SVM formulations (Lee et al. and Weston and Watkins) that parallelize the computation evenly over the number of classes. This allows us to compare these models to one-vs.-rest SVMs on unprecedented scale. The results indicate superior accuracy on text classification data.

1 Introduction

Modern data analysis requires computation with a large number of classes. As examples, consider the following. (1) We are continuously monitoring the internet for new webpages, which we would like to categorize. (2) We are collecting data from an online feed of photographs that we would like to classify into image categories. (3) We add new articles to an online encyclopedia and intend to predict the categories of the articles.

The problems above involve a large number of classes, typically at least in the thousands. This motivates research on scaling up multi-class classification methods. In the present work, we address scaling up multi-class support vector machines (MC-SVMs) [1]. There are two major types of MC-SVMs:

- One-vs.-one (OVO) and one-vs.-rest (OVR) MC-SVMs decompose the problem into multiple binary subproblems that are subsequently aggregated [1, 2]. Training can be parallelized in a straight forward way.
- 2. All-in-one MC-SVMs extend the concept of the margin to multiple classes. Because there is no unique extension of the margin concept, multiple all-in-one MC-SVMs have been proposed, including the ones by Crammer and Singer (CS) [3], Lee, Lin, and Wahba (LLW) [4], and Weston and Watkins (WW) [5, 1]. See [2, 6, 7, 8, 9, 10, 11] for conceptual and empirical comparisons.

Recently, Dogan et al. [11] have compared the various all-in-one MC-SVM variants on rather moderately sized datasets and showed advantages of all-in-one MC-SVMs over OVR MC-SVM, but — so far — slow training time has prohibited comparisons on data involving a large number of classes.

The reason is that (linear) state of the art solvers require time complexity of $\mathcal{O}(\bar{d}\bar{n}\cdot\mathcal{C}^2)$ and memory complexity at least of $\mathcal{O}(\bar{n}\mathcal{C}^2)$, where d is the feature dimensionality, \bar{d} the average number of non-zeros ($\bar{d}=d$ for dense data), and \bar{n} the average number of samples per class. This quadratic dependence on the number of classes \mathcal{C} can be prohibitive for large \mathcal{C} , often leaving OVO and OVR as the only MC-SVM options in the big data setting.

In this paper, we develop distributed algorithms where up to $\mathcal{O}(\mathcal{C})$ nodes solve WW and LLW in parallel, dividing model and computation evenly.

The algorithm proposed for WW draws inspiration from a major result in graph theory: the solution to the 1-factorization problem of a graph [12]. The idea is that the optimization of a single coordinate $\alpha_{i,c}$ of the dual objective involves only the two hypotheses w_{y_i} and w_c . As in the 1-factorization problem, we can thus form pairs of classes where the corresponding blocks of coordinates can be optimized in parallel.

[†] Shared first co-authorship, sorted alphabetically. Preliminary work. Under review by AISTATS 2017. Do not distribute.

On the other hand, we parallelize LLW training by introducing an auxiliary variable \bar{w} into the dual problem that decouples the objective into a sum over \mathcal{C} many independent subproblems.

We provide both multi-core and distributed implementations of the proposed algorithms. We report on empirical runtime comparisons of the proposed solvers with the one-vs.-rest implementation by LIBLINEAR [13] on text classification data taken from the LSHTC corpus [14].

The main contributions of this paper are the following:

- We propose the first distributed, exact solver for WW and LLW.
- We provide both multi-core and truly distributed implementations of the solver.
- We give the first comparison of WW, LLW, and OVR on the DMOZ data from the LSHTC '10-'12 corpora using the full feature resolution.

We expect that the present work starts a line of research on parallelization of exact training of various all-in-one MC-SVMs, including Crammer and Singer, multi-class maximum margin regression [15], and the reinforced multicategory SVM[16], enabling comparison of all these methods on large datasets.

The paper is structured as follows. In the next section we discuss the problem setting and preliminaries. In Section 3, we present the proposed distributed algorithms for LLW and WW, respectively. We analyze their convergence empirically in Section 4. Section 5 discusses related work and Section 6 concludes.

2 Preliminaries

We consider the following problem. We are given data $(x_1, y_1), \ldots, (x_n, y_n)$ with $x_i \in \mathbb{R}^d$ and $y_i \in \{1, \ldots, \mathcal{C}\}$. Each class has in average \bar{n} samples. The largest number of samples for a single class is n_{max} . We are predicting using the model

$$\hat{y}(x) := \operatorname*{argmax}_{c} w_{c}^{T} x, \tag{1}$$

where $W = (w_1, ..., w_c) \in \mathbb{R}^{d \times c}$ are unknown parameters. The aim is to efficiently find good parameters in order to predict well on new data using (1).

To address this problem setting, a number of generalizations of the binary SVM [17] have been proposed. We are specifically studying the following two formulations, dropping the bias terms in both. Throughout this paper, $l(x) = \max\{0, 1-x\}$ will denote the hingeloss.

Lee, Lin, and Wahba (LLW) [4]

$$\min_{W} \sum_{c=1}^{C} \left[\frac{1}{2} ||w_c||^2 + C \sum_{i: y_i \neq c} l(-w_c^T x_i) \right]$$
s.t.
$$\sum_{c} w_c = 0$$
(2)

Weston and Watkins (WW) [5]

$$\min_{W} \sum_{c=1}^{\mathcal{C}} \left[\frac{1}{2} ||w_c||^2 + C \sum_{i: y_i \neq c} l(w_{y_i}^T x_i - w_c^T x_i) \right]$$
(3)

Both formulations lead to very similar dual problems, as shown below. For the dualization of WW, we refer to [18]. The LLW dual is given in Appendix A, where we introduce an auxiliary variable \bar{w} that is exploited by our distributed algorithm, as explained in the next section.

$$\max_{\alpha \in \mathbb{R}^{n \times C}, \bar{w} \in \mathbb{R}^{d}} \sum_{c=1}^{C} \left[-\frac{1}{2} || -X\alpha_{c} + \bar{w}||^{2} + \sum_{i: y_{i} \neq c} \alpha_{i,c} \right]$$
s.t.
$$\forall i: \ \alpha_{i,y_{i}} = 0, \forall c \neq y_{i}: 0 \leq \alpha_{i,c} \leq C$$
(LLW)

$$\max_{\alpha \in \mathbb{R}^{n \times c}} \quad \sum_{c=1}^{\mathcal{C}} \left[-\frac{1}{2} || -X\alpha_c ||^2 + \sum_{i: y_i \neq c} \alpha_{i,c} \right]$$
s.t.
$$\forall i: \ \alpha_{i,y_i} = -\sum_{c: c \neq y_i} \alpha_{i,c},$$

$$\forall c \neq y_i: \ 0 < \alpha_{i,c} < C$$
(WW)

3 Distributed Algorithms

In this section, we derive algorithms that solve (LLW) and (WW) in a distributed manner. With start by addressing LLW.

3.1 Algorithm for Lee, Lin, and Wahba

First note the following optimality condition in (LLW):

$$\bar{w} = \frac{1}{\mathcal{C}} \sum_{c=1}^{\mathcal{C}} X \alpha_c. \tag{4}$$

Which was exploited by prevalent solvers to remove the variable \bar{w} from the optimization, yielding the objective $\text{obj}(\alpha) = \sum_{c=1}^{\mathcal{C}} \left[-\frac{1}{2} || -X\alpha_c + \bar{w}||^2 + \sum_{i:y_i \neq c} \alpha_{i,c} \right]$. In contrast, the core idea of our LLW solver is to actually

Algorithm 1 Lee, Lin, and Wahba

```
1: function SOLVE-LLW(C, X, Y)
          for c = 1..C do in parallel
 2:
 3:
               w_c \leftarrow 0
 4:
               \alpha_c \leftarrow 0
               for i \in I do
 5:
                    k_i \leftarrow x_i^T x_i
 6:
 7:
               while not optimal do
                    optimal \leftarrow True
 8:
                    shuffleData()
 9:
                    for i \in I \setminus I_c do
10:
11:
                         solve1DimLLW(i, c)
                    \bar{w} \leftarrow \text{Reduce}(\sum_{c} w_c / \mathcal{C})
12:
13:
                    w_c \leftarrow w_c - \bar{w}
```

keep this auxiliary variable, as it decouples the objective function into the following sum:

$$obj(\alpha) = \sum_{c=1}^{C} subobj(\alpha_c, \bar{w}).$$
 (5)

Then we perform dual block coordinate ascent (DBCA) [18, Algorithm 3.1] with a specially tailored block structure, considering as blocks \bar{w} as well as each single coordinate $\alpha_{i,c}$. As we observe from (5), the optimization of the columns $\alpha_{:,c}$ is mutually independent of each other, given fixed \bar{w} . Hence, it can be distributed evenly over \mathcal{C} many nodes. On the cth node, we run coordinate ascend on the subobjective subobj (α_c, \bar{w}) over $\alpha_{i,c}, i = 1, \ldots, n$, as described in the next paragraph. After one epoch of α computation, the variable \bar{w} is updated via (4). The final algorithm is shown in Algorithm 1.

As necessary step within Algorithm 1, we need to update every single $\alpha_{i,c}$. Optimizing $\alpha_{i,c}$ is solving the problem

$$\underset{\delta}{\operatorname{argmax}} D(\alpha + \delta e_{i,c}, \bar{w})$$
s.t. $0 \le \alpha_{i,c} + \delta \le C$, (6)

where $e_{i,c} \in \mathbb{R}^{n \times \mathcal{C}}$ is one at the (i,c)th coordinate and zero elsewise. Set $w_c := -X\alpha_c + \bar{w}$; then the gradient for δ is $\frac{\partial}{\partial \delta} \left[D(\alpha + \delta e_{i,c}) \right] = x_i^T w_c - x_i^T x_i \delta + 1$. Hence, the optimal solution of (6) is given by $\delta = \min\{C - \alpha_{i,c}, \max\{-\alpha_{i,c}, -\frac{x_i^T w_c - 1}{x_i^T x_i}\}\}$. The corresponding pseudo-code can be found in Algorithm 2.

3.1.1 Convergence

It is well known that the block coordinate ascent method converges under suitable regularity conditions [e.g., 19, 20]. Our objective is continuously differentiable and strictly convex. The constraints are solely

Algorithm 2 Solving 1-dim sub-problem of LLW

```
1: function SOLVE1DIMLLW(i,c)
 2:
           global C, X, k, \alpha_c, w_c, optimal
           g \leftarrow w_c^T x_i - 1
 3:
           if g < -\epsilon and \alpha_{i,c} < C then
 4:
 5:
                 \delta \leftarrow \min\{C - \alpha_{i,c}, -g/k_i\}
                 optimal \leftarrow False
 6:
 7:
           if g > \epsilon and \alpha_{i,c} > 0 then
 8:
                 \delta \leftarrow \max\{-\alpha_{i,c}, -g/k_i\}
                 optimal \leftarrow False
 9:
           w_c \leftarrow w_c + \delta x_i
10:
           \alpha_{i,c} \leftarrow \alpha_{i,c} + \delta
11:
```

box constraints, hence the feasible set decomposes as a Cartesian product over the blocks. Algorithm 1 traverses the two blocks in cyclic order. Under these conditions, the DBCA method provably converges [e.g, 20, Prop. 2.7.1].

Note that in practice, we observed speedups by updating \bar{w} in Algorithm 1 after each tenth of an epoch, breaking the cyclic order. The blocks of coordinates are then traversed in so-called essentially cyclic order [e.g., 19, Section 2], meaning that there exists $T \in \mathbb{N}$ such that each block is traversed at least once after T iterations. Closer inspection of the proof in [e.g., 20, Prop. 2.7.1] reveals that the result holds also under this slightly more general assumption.

Further, we drop variables $\alpha_{i,c}$ in the optimization (shrinking) if they are not updated in three subsequent epochs. Once the stopping condition holds, we run another epoch of optimization over all variables (including the ones that were dropped). If the stopping criterion is then met, we terminate the algorithm and continue optimization elsewise.

3.1.2 Implementation details

Our implementation uses OpenMPI for inter-machine [21, MPI] and OpenMP [22, MC] for intra-machine communication. Note that Algorithm 1 has very mild communication requirements: the only communication needed is the sum of all weight vectors $\bar{w} = \sum_c w_c$. Hence, MPI suffers very little from communication overhead between the various machines. In practice, we may not be able to fully parallelize to the maximum of \mathcal{C} cores; therefore our algorithm will divide the set of classes into number-of-cores many chunks and optimize each class sequentially.

3.2 Algorithm for Weston and Watkins

In this section, we propose a distributed algorithm for WW, which draws inspiration from the 1-factorization

Algorithm 3 Solving 1-dim sub-problem of WW

```
1: function SOLVE1DIMWW(i,c)
            global C, X, w_{y_i}, w_c, \alpha_c, \text{optimal}
           g \leftarrow (w_{y_i}^T - w_c^T)x_i - 1
 3:
            if g < -\epsilon and \alpha_{i,c} < C then
  4:
                 \delta \leftarrow \min\{C - \alpha_{i,c}, -g/2k_i\}
  5:
                  optimal \leftarrow False
 6:
  7:
            if q > \epsilon and \alpha_{i,c} > 0 then
                 \delta \leftarrow \max\{-\alpha_{i,c}, -g/2k_i\}
 8:
 9:
                 optimal \leftarrow False
            w_{y_i} \leftarrow w_{y_i} + \delta x_i
10:
            w_c \leftarrow w_c - \delta x_i
11:
12:
            \alpha_{i,c} \leftarrow \alpha_{i,c} + \delta
```

problem of a graph. The approach is presented below.

3.2.1 Preliminaries

Our approach is based on running dual coordinate ascend [18, Algorithm 3.1] over $\alpha_{i,c}$ on the (WW) objective function as follows. Denote the objective in (WW) by $D(\alpha)$ and recall from [18] that optimizing $\alpha_{i,c}$ is solving the problem

$$\underset{\delta}{\operatorname{argmax}} D(\alpha + \delta \mathbf{e}_{i,c})$$

s.t. $0 \le \alpha_{i,c} + \delta \le C$.

Setting $w_c = \sum_{i:y_i \neq c} (-x_i \alpha_{i,c} + \sum_{c:c \neq y_i} x_i \alpha_{i,c}),$ the gradient for δ is given by $\frac{\partial}{\partial \delta} [D(\alpha + \delta \mathbf{e}_{i,c})] = -x_i^T (w_{y_i} - w_c) - x_i^T x_i \delta + 1$. Which is optimal at

$$\delta = \min \left(C - \alpha_{i,c}, \max \left(-\alpha_{i,c}, \frac{x_i^T(w_{y_i} - w_c) - 1}{2x_i^T x_i} \right) \right).$$
(7)

This computation is summarized in Algorithm 3.

3.2.2 Core Observation

We observe from above that optimizing with regard to $\alpha_{i,c}$ will require only the weight vectors w_{y_i} and w_c . In other words, given four different classes c_1, c_2, c_3, c_4 , the optimization of the block of variables $(\alpha_i, c_1)_{i:y_i=c_2}$ —according to (7)—is independent of the optimization of the block $(\alpha_i, c_3)_{i:y_i=c_4}$. Hence it can be parallelized. In the next section we describe how we exploit this structure to derive a distributed optimization algorithm.

3.2.3 Excursus: 1-Factorization of a Graph

Assume that C is even. The core idea now is to form $\frac{C}{2}$ many disjoint blocks $(\alpha_i, c_1)_{i:y_i=c_2}, \ldots, (\alpha_i, c_{C-1})_{i:y_i=C}$ of variables. Each of these blocks can be optimized in parallel. The challenge is to derive a maximally

distributed optimization schedule where each block $(\alpha_i, c_j)_{i:y_i=c_k}$ for any $j \neq k$ is optimized.

To better understand the problem, we consider the following analogy. We are given a football league with $\mathcal C$ teams. Before the season, we have to decide on a schedule such that each team plays any other team exactly once. Furthermore, all teams shall play on every matchday so that in total we need only $\mathcal C-1$ matchdays. This problem is the 1-factorization problem in graph theory [e.g., 12]. The solution to this problem, illustrated in Figure 1, is as follows.

We arrange one node centrally and all other nodes in a regular polygon around the center node. At round r, we connect the centered node with node r and connect all other nodes orthogonal to this line. The pseudocode to compute the partner of a given node c at a certain round r is given in Algorithm 5 in the appendix. Note that in case of an uneven number of classes, we introduce a dummy class $\mathcal{C}+1$, making the number of classes even. We run the algorithm, but skip all computations involving the dummy class.

3.2.4 Algorithm

The algorithm, shown in Algorithm 4, performs DBCA over the variables $\alpha_{i,c}$ using the schedule derived in Section 3.2.3 and the coordinate updates derived in Section 3.2.1.

3.2.5 Convergence and Implementation details

Furthermore, note that our algorithm performs the same coordinate updates as Algorithm 3.1 in [18]. Hence, they share the same favorable convergence behavior. Formally, convergence is guaranteed for exactly the same reasons discussed in Section 3.1.2. We also employ the same speedup tricks, i.e. shrinking and updating every tenth of an epoch.

In practice, because of limitations of computational resources, we might not be able to fully parallelize to the maximum of $\mathcal{C}/2$ cores. In that case, our algorithm divides the set of classes into number-of-cores many chunks and solves each bundle sequentially. For optimal speedup, it is advisable to arrange the classes

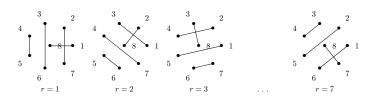


Figure 1: Illustration of the solution of the 1-factorization problem of a graph with C=8 many nodes.

Algorithm 4 Watkins-Weston

```
1: function Solve-WW(c,X,Y)
 2:
         for c = 1..C do in parallel
 3:
              w_c \leftarrow 0
              \alpha_c \leftarrow 0
 4:
         for i \in I do
 5:
              k_i \leftarrow x_i^T x_i
 6:
         while not optimal do
 7:
              optimal \leftarrow True
 8:
 9:
              shuffleData()
              for r = 1..C - 1 do
10:
                  for c = 1..C do in parallel
11:
                       \tilde{c} \leftarrow \text{matchClass}(\mathcal{C}, c, r)
12:
                       if \tilde{c} > c then
13:
                            for i \in I_c do
14:
                                solve1DimWW(i, \tilde{c})
15:
                            for i \in I_{\tilde{c}} do
16:
                                solve1DimWW(i, c)
17:
```

into chunks of equal number of classes and data points.

As with LLW, we implemented a mixed MPI-OpenMP solver for WW. However, note that, while LLW has mild communication needs, WW needs to pair the weight vectors of the matched classes c and \tilde{c} in each epoch, for which $\mathcal{C}/2$ weight vectors needs to communicated among computers. Therefore it is crucial to communicate efficiently.

We tackled the problem as follows. First of all, we use OpenMP for computations on a single machine (efficiently parallelizing among cores). Here, due to the shared memory, no weight vectors need to be moved. The more challenging task is to handle inter-machine communication efficiently. Our approach is based on two key observations.

If the data is high-dimensional data, yet sparse, we keep the full weight matrix in memory for fast access, yet communicating only the non-zero entries between computers. Regardless of the increased computational effort, this takes only a fraction of time compared to sending the dense data.

Furthermore, we relax the WW matching scheme. Coming back to a football, consider each country hosts a league, and inside the league, we match the teams as known. Now we would like to match teams across leagues. In order to do so, we first match the countries with the scheme from Section 3.2.3. For each pair of countries, call them A and B, every team from country A plays any other team from country B. Coming back to classes and machines, this means we transfer bundles of classes (countries) between computers. This drastically reduces the network communication.

4 Experiments

This section is structured as follows. First we empirically verify the soundness of the proposed algorithms. Then we introduce the employed datasets, on which we investigate the convergence and runtime behavior of the proposed algorithms as well as the induced classification performance.

For the experimental setup we refer to B.1 in the supplement. Each training algorithm was run three times, using randomly shuffled data, and the results were averaged. Note that the training set is the same in each run, but the different order of data points can impact the runtime of the algorithms.

4.1 Validation of solver

In our first experiment, we validate the correctness of the proposed solvers. We downloaded data from the LIBLINEAR [13]¹ and UCI [23]² dataset repositories. Where training and test splits are unavailable, we split the data once into 90% train and 10% test sets. For each dataset, the optimal feature scaling was selected, in order to maximize the average accuracy on the test sets. Datapoints in iris and news were thus normalized to unit norm, letter and satimage were normalized to unit variance. All other data was considered unnormalized.

Then we compare our LLW and WW solvers with the state-of-the-art implementation contained in the ML library shark [24]. We implemented the same stopping criteria as [24]. The full results (averaged over 10 runs) are shown in Table B.1 in the supplement. A subset is given in Table 1. We observe good accordance of the results and model sparsity of the proposed solvers and the reference implementation from the Shark toolbox, thus confirming that our respective solvers are indeed exact solvers of LLW and WW.

At random we tested whether the duality-gap closes or not. We did this for both solvers with different C values and datasets. In any case the duality gap closed, i.e. decreased by an order of two magnitudes. Based on this we chose our stopping criteria ϵ equal to 0.1 for the LSHTC datasets.

4.2 Datasets

We experiment on large classification datasets, where the number of classes ranges between 451 and 27,875. The relevant statistics of the datasets are shown in 2. The LSHTC-* datasets are high-dimensional text datasets taken from the well-known LSHTC corpus

https://www.csie.ntu.edu.tw/~cjlin/liblinear/

²https://archive.ics.uci.edu/ml/datasets.html

Table 1: Error on the test set and density in % of the Shark solver (denoted S) and the proposed solver (denoted D). The complete table can be found in the supplement.

		Er	ror		$Model ext{-}Density$			
Dataset:	D-LLW	S-LLW	$\mathbf{D}\text{-}\mathbf{W}\mathbf{W}$	S-WW	D-LLW	S-LLW	\mathbf{D} - \mathbf{W} \mathbf{W}	S-WW
news20								
$\log(C)$: -1	29.23	29.23	15.32	15.30	97.24	97.24	51.16	49.72
0	22.97	22.97	14.80	14.80	97.24	97.24	44.74	42.70
1	16.15	16.15	15.98	15.98	97.17	97.04	45.97	43.47
rcv1								
$\log(C)$: -1	47.96	47.96	11.31	11.31	78.00	78.00	26.42	23.45
0	33.27	33.41	11.52	11.52	78.00	77.98	22.93	20.12
1	12.03	12.03	12.03	12.03	78.00	77.98	23.05	20.06

Table 2: The used datasets and their properties.

Dataset	n train	n test	$\mathcal C$	\mathbf{d}
LSHTC-small	4,463	1,858	1,139	51,033
LSHTC-large	128,710	34,880	12,294	381,581
LSHTC-2012	383,408	103,435	11,947	575,555
LSHTC-2011	394,754	104,263	27,875	594,158

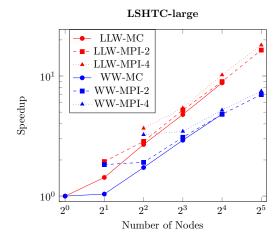


Figure 2: Speed-up averaged over 10 repetitions respectively in the number of cores.

[14]. The datasets belong to the released competition rounds 1 to 3, i.e. '10-'12. LSHTC-2011 and LSHTC-2012 originate from the DMOZ corpus. The features were extracted using TF/IDF representation and we use the full feature resolution for training.

4.3 Speedup

In order to measure the speedup provided by increasing the number of machines/cores, we run a fix amount of iterations over the whole LSHTC-large dataset. We use 10 runs over 10 iterations with a fixed parameter C equal 1 without shrinking. While the MC execution works on one machine, the MPI executes on two or four machines, i.e. spreading the used cores evenly on each node.

The results are shown in Figure 2. LLW exhibits linear speedup regardless if distributed or not, due to the

small communication cost.

For less than 4 cores WW has a similar, linear speedup. Then we assume that the cache bound is reached³. When distributing on two machines (MPI), i.e doubling the cache-throughput, despite the communication cost a higher speedup is reached. This again vanishes as more cpus per machine are used. This confirms our assumption and suggests that on a system with more cache capacity a better speedup can be reached.

4.4 Timing and Classification Results

Now we evaluate and compare the proposed algorithms on the LSHTC datasets for a range of C values, i.e. we perform no cross-validation. For comparison we use solvers from the well-known *LIBLINEAR* package, namely the multi-core implementation with L2L1-loss (OVR, [25]) and the Crammer-Singer implementation (CS, [13]). For the multi-core solvers, i.e. OVR and WW-MC, we use 16 cores. MPI spreads over 2 or 4 machines using 8 and 4 cores respectively at each node, thus trains the model distributed. Table 3 shows the error and the model sparsity for the compared solutions. We further provide the Micro-F1 and Macro-F1 score in Table B.2 in the supplement.

For all datasets the canonical multi-class formulations, i.e. CS and WW, perform significantly better than OVR. On one hand the error is smaller and the F1-scores better. On the other hand the learned models are much sparser, i.e. up to a magnitude. The results justify the increased solution complexity of canonical formulations.

Comparing CS and WW, CS performs as well or slightly better at classifying. Though WW leads to a sparser model. To the best of our knowledge this is the first comparison of these well-known multi-class SVMs on the studied LSHTC data.

From Figure 3, we observe that the runtime of our solver outperforms the one of OVR and CS by up to

 $^{^3{\}rm The}$ used machines have two cpu-socket, i.e. two major caches.

Table 3: Test set error and model density (in %) as achieved by the OVR, CS, WW, and LLW solvers on the LSHTC datasets. For each solver the result with the best error is in bold font. For LLW entries with a '*' did not converge within a day of runtime.

	Error			$Model ext{-}Density$				
Dataset:	OVR	\mathbf{CS}	$\mathbf{W}\mathbf{W}$	\mathbf{LLW}	OVR	\mathbf{CS}	$\mathbf{W}\mathbf{W}$	\mathbf{LLW}
LSHTC-small								
$\log(C)$: -3	93.00	59.74	72.82	93.00	92.74	11.11	69.73	92.74
-2	85.36	59.74	65.34	93.00	81.54	11.13	16.44	92.74
-1	74.54	59.74	57.59	93.00	46.76	11.12	6.06	92.74
0	64.37	55.49	54.57	93.00	38.20	11.76	5.74	92.74
1	57.75	54.57	54.41	93.00	38.63	11.69	5.73	92.74
LSHTC-large								
$\log(C)$: -3	88.12	58.57	66.47	95.86	75.26	2.53	18.50	100.0
-2	85.21	58.57	60.58	95.86	45.14	2.53	4.45	100.0
-1	77.96	57.82	55.28	95.86	25.28	2.55	1.71	100.0
0	63.11	53.61	53.98	95.86	18.33	2.69	1.61	100.0
1	57.18	54.18	54.41	*	18.55	2.67	1.66	*
LSHTC-2012								
$\log(C)$: -3	83.66	49.81	58.02	92.63	72.60	1.73	16.97	99.52
-2	75.15	49.65	50.20	92.63	46.20	1.71	4.06	99.52
-1	60.38	46.14	44.94	92.63	25.87	1.76	1.52	99.52
0	47.33	42.67	44.01	*	18.20	2.06	1.42	*
1	46.83	45.60	46.15	*	18.46	2.09	1.47	*
LSHTC-2011								
$\log(C)$: -3	87.95	59.09	68.19	96.18	72.38	1.57	13.49	100.0
-2	85.85	59.09	62.14	96.18	45.97	1.57	3.16	100.0
-1	76.78	58.18	57.31	96.18	25.97	1.55	1.19	100.0
0	63.11	55.58	56.94	*	18.24	1.69	1.11	*
1	60.01	57.78	58.32	*	18.46	1.70	1.14	*

Table 4: Error, Micro-F1, and Macro-F1 on the test set and model density in % of the LLW solver on the LSHTC-small dataset.

$\log(C)$:	2	3	4
Error:	87.73	66.74	59.31
Micro F1:	2.08	15.07	40.69
Macro F1:	12.27	33.26	24.58
W-Density:	92.74	92.74	92.74
α -Density:	99.88	99.87	99.90

two orders of magnitude. Even when distributed our solver outperforms multi-core OVR in all except one case.

All WW experiments use the same amount of cores, but with a varying degree of distribution. We observe that the communication imposes a modest overhead.

Figure 3 confirms our assumption from the previous section. First, please note that shrinking reduces the active training set, i.e. reduces the computational effort and puts stress on the communication overhead. Therefore WW-MPI is regardless the higher speedup (see Figure 2) slower than the multi-core version. Yet on a computationally intensive dataset as LSHTC-2011 the higher speedup cancels the overhead due to the communication for C equal to 0.1 and 1.

4.4.1 Lin, Lee, & Wahba

Knowing that LLW converges to the correct solution, as the duality-gap closes, the results indicate that the chosen C range is not suitable. For LSHTC-small we conducted experiments with much larger C values. And indeed, as shown in Table 4, LLW performs best in a nearly unconstrained setting. In any case, the model learned by LLW is never sparse. Not in the weight matrix W, nor in dual factors α . Resource limitations and slow convergence properties hindered us to conduct experiments with even larger C values. It is left to future work to explore this space or even develop a unconstrained version of LLW.

5 Discussion of Related Work

Most approaches to parallelization of MCSVM training are based on OVO or OVR [26], including a number of approaches that attempt to learn a hierarchy of labels [27, 28, 29, 30, 31, 32] or train ensembles of SVMs on individual subsets of the data [33, 34, 35, 36].

There is a line of research on parallelizing stochastic gradient (SGD) based training of MC-SVMs over multiple computers [37, 38]. SGD builds on iteratively approximating the loss term by one that is based on a subset of the data (mini-batch). In contrast, batch solvers (such as the ones proposed in the present pa-

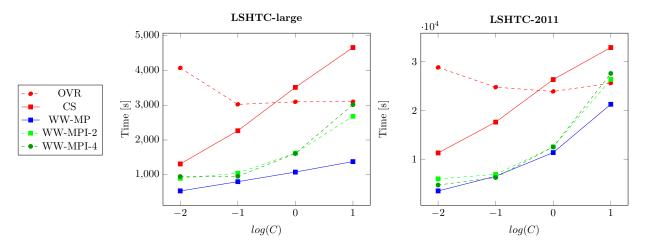


Figure 3: Training time averaged over 10 repetitions per C.

per) are based on the full sample. In this sense, our approach is completely different to SGD. While there is a long ongoing discussion whether the batch or the SGD approach is superior, the common opinion is that SGD has its advantages in the early phase of the optimization, while classical batch solvers shine in the later phase. In this sense, the two approaches are complementary and could also be combined.

The related work that is the most closest to the present work is by [39]. They build on the alternating direction method of multipliers (ADMM) [40] to break the Crammer and Singer optimization problem into smaller parts, which can be solved individually on different computers. In contrast to our approach, the optimization problem is parallelized over the samples, not the optimization variables. In our problem setting, high-dimensional sparse data, the model size is vary large. Because each node holds the whole model in memory, this algorithm hardly scales with large label spaces. E.g. consider Table 2; the model for LSHTC-2011 contains $\approx 16 * 10^9$ parameters. Note also that it is unclear at this point whether the approach of [39] could be adapted to LLW and WW, which are the object of study in the present paper.

Note that beyond SVMs there is a large body of work on distributed multi-class [e.g., 41, 42] and multi-label learning algorithms [43], which is outside of the scope of the present paper.

6 Conclusion

We proposed distributed algorithms for solving the multi-class SVM formulations by Lee et al. (LLW) and Weston and Watkins (WW). The algorithm addressing LLW takes advantage of an auxiliary variable, while our approach to optimizing WW in parallel is based

on the 1-factorization problem from graph theory.

The experiments confirmed the correctness of the solver (in the sense of an exact solver) and show linear speedup when the number of cores is increased. This speedup allows us to train LLW and WW on LSHTC datasets, for which results were lacking in the literature.

Our analysis contributed to comparing MC-SVM formulations on rather large data sets, where comparisons were still lacking. In comparison to OVR we showed that WW can achieve competitive classification results in less time, while still leading to a much sparser model. Unexpectedly, LLW shows clear disadvantages over the other MC-SVMs. Yet the favorable scaling properties make further research interesting, for instance regarding the development of an unconstrained algorithm. We ease further research by publishing the source code under https://github.com/albermax/xcsvm.

Overcoming the limitations of a single machine, i.e. distribution, is a key problem and a key enabler in large scale learning. To best of our knowledge, we are the first to train an exact, all-in-one MC-SVMs in a distributed manner. We hope this first step inspires further research in this context.

In the future, we would like to study extensions of the concepts presented in this paper to various more MC-SVMs, including the Crammer and Singer MC-SVM [44], multi-class maximum margin regression [15], and the reinforced multicategory SVM[16].

Acknowledgments

We thank Rohit Babbar, Shinichi Nakajima, and Klaus-Robert Müller for helpful discussions. We thank Giancarlo Kerg for pointing us to the graph 1-factorization problem. We thank Ioannis Partalas for help regarding the LSHTC datasets. MK acknowledges support by the German Research Foundation (DFG) under KL 2698/2-1 and by the Federal Ministry of Education and Research (BMBF) under 031L0023A and 031B0187B. MA acknowledges support by the Federal Ministry of Education and Research (BMBF) under 01IS14013A.

References

- V. Vapnik, Statistical Learning Theory. John Wiley and Sons, 1998.
- [2] R. Rifkin and A. Klautau, "In defense of one-vs-all classification," *Journal of Machine Learning Research*, vol. 5, pp. 101–141, 2004.
- [3] K. Crammer and Y. Singer, "On the algorithmic implementation of multiclass kernel-based vector machines," *Journal of Machine Learning Research*, vol. 2, pp. 265–292, 2002.
- [4] Y. Lee, Y. Lin, and G. Wahba, "Multicategory support vector machines: Theory and application to the classification of microarray data and satellite radiance data," *Journal of the American Statistical Association*, vol. 99, no. 465, pp. 67–82, 2004.
- [5] J. Weston and C. Watkins, "Support vector machines for multi-class pattern recognition," in *Proceedings of* the Seventh European Symposium On Artificial Neural Networks (ESANN) (M. Verleysen, ed.), pp. 219–224, Evere, Belgium: d-side publications, 1999.
- [6] E. L. Allwein, R. E. Schapire, and Y. Singer, "Reducing multiclass to binary: A unifying approach for margin classifiers," *Journal of Machine Learning Research*, vol. 1, pp. 113–141, 2001.
- [7] C. W. Hsu and C. J. Lin, "A comparison of methods for multiclass support vector machines," *IEEE Trans*actions on Neural Networks, vol. 13, no. 2, pp. 415– 425, 2002.
- [8] S. I. Hill and A. Doucet, "A framework for kernel-based multi-category classification," *Journal of Artificial Intelligence Research*, vol. 30, no. 1, pp. 525–564, 2007.
- [9] Y. Liu, "Fisher consistency of multicategory support vector machines," in *Eleventh International Confer*ence on Artificial Intelligence and Statistics (AIS-TATS) (M. Meila and X. Shen, eds.), vol. 2 of JMLR W&P, pp. 289–296, 2007.
- [10] Y. Guermeur, "VC theory for large margin multicategory classifiers," *Journal of Machine Learning Re*search, vol. 8, pp. 2551–2594, 2007.
- [11] Ü. Doğan, T. Glasmachers, and C. Igel, "A Unified View on Multi-class Support Vector Classification," *Journal of Machine Learning Research*, vol. 17, no. 45, pp. 1–32, 2016.
- [12] J. A. Bondy and U. S. R. Murty, Graph theory with applications, vol. 290. Macmillan London, 1976.
- [13] R. Fan, K. Chang, C. Hsieh, X. Wang, and C. Lin, "LIBLINEAR: A library for large linear classification," *Journal of Machine Learning Research*, vol. 9, pp. 1871–1874, 2008.

- [14] I. Partalas, A. Kosmopoulos, N. Baskiotis, T. Artières, G. Paliouras, É. Gaussier, I. Androutsopoulos, M. Amini, and P. Gallinari, "LSHTC: A benchmark for large-scale text classification," CoRR, vol. abs/1503.08581, 2015.
- [15] S. Szedmak, J. Shawe-Taylor, and E. Parado-Hernandez, "Learning via linear operators: Maximum margin regression," tech. rep., PASCAL, Southampton, UK, 2006.
- [16] Y. Liu and M. Yuan, "Reinforced multicategory support vector machines," *Journal of Computational and Graphical Statistics*, vol. 20, no. 4, pp. 901–919, 2011.
- [17] C. Cortes and V. Vapnik, "Support-vector networks," Machine Learning, vol. 20, no. 3, pp. 273–297, 1995.
- [18] S. S. Keerthi, S. Sundararajan, K.-W. Chang, C.-J. Hsieh, and C.-J. Lin, "A sequential dual method for large scale multi-class linear syms," in *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '08, (New York, NY, USA), pp. 408–416, ACM, 2008.
- [19] P. Tseng, "Convergence of a block coordinate descent method for nondifferentiable minimization," *Journal* of optimization theory and applications, vol. 109, no. 3, pp. 475–494, 2001.
- [20] D. P. Bertsekas, M. L. Homer, D. A. Logan, and S. D. Patek, "Nonlinear programming," Athena scientific, 1995.
- [21] W. Gropp, E. Lusk, N. Doss, and A. Skjellum, "A high-performance, portable implementation of the mpi message passing interface standard," *Parallel* computing, vol. 22, no. 6, pp. 789–828, 1996.
- [22] L. Dagum and R. Enon, "Openmp: an industry standard api for shared-memory programming," Computational Science & Engineering, IEEE, vol. 5, no. 1, pp. 46–55, 1998.
- [23] A. Asuncion and D. Newman, "Uci machine learning repository," 2007.
- [24] C. Igel, T. Glasmachers, and V. Heidrich-Meisner, "Shark," *Journal of Machine Learning Research*, vol. 9, pp. 993–996, 2008.
- [25] W.-L. Chiang, M.-C. Lee, and C.-J. Lin, "Parallel dual coordinate descent method for large-scale linear classification in multi-core environments," in *Proceedings* of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16, (New York, NY, USA), pp. 1485–1494, ACM, 2016.
- [26] R. Babbar, K. Maundet, and B. Schölkopf, "Tersesvm: A scalable approach for learning compact models in large-scale classification," in *Proceedings of the* 2016 SIAM International Conference on Data Mining, pp. 234–242, SIAM, 2016.
- [27] S. Bengio, J. Weston, and D. Grangier, "Label embedding trees for large multi-class tasks," in Advances in Neural Information Processing Systems, pp. 163–171, 2010.
- [28] J. Deng, S. Satheesh, A. C. Berg, and F. Li, "Fast and balanced: Efficient label tree learning for large scale object recognition," in *Advances in NIPS*, pp. 567–575, 2011.

- [29] T. Gao and D. Koller, "Discriminative learning of relaxed hierarchy for large-scale visual recognition," in *ICCV*, pp. 2072–2079, IEEE, 2011.
- [30] A. E. Choromanska and J. Langford, "Logarithmic time online multiclass prediction," in *Advances in Neural Information Processing Systems* 28, pp. 55–63, Curran Associates, Inc., 2015.
- [31] D. Zhou, L. Xiao, and M. Wu, "Hierarchical classification via orthogonal transfer," in *ICML*, pp. 801–808, 2011.
- [32] S. Gopal and Y. Yang, "Recursive regularization for large-scale classification with hierarchical and graphical dependencies," in ACM SIGKDD, pp. 257–265, ACM, 2013.
- [33] A. Govada, B. Gauri, and S. K. Sahay, "Distributed multi class sym for large data sets," in *Proceedings* of the Third International Symposium on Women in Computing and Informatics, WCI '15, (New York, NY, USA), pp. 54–58, ACM, 2015.
- [34] A. Govada, S. Ranjani, A. Viswanathan, and S. Sahay, "A novel approach to distributed multi-class svm," arXiv preprint arXiv:1512.01993, 2015.
- [35] S. Lodi, R. Nanculef, and C. Sartori, "Single-pass distributed learning of multi-class syms using core-sets," methods, vol. 14, no. 27, p. 2, 2010.
- [36] R. Jenssen, M. Kloft, A. Zien, S. Sonnenburg, and K.-R. Müller, "A scatter-based prototype framework and multi-class extension of support vector machines," *PloS one*, vol. 7, no. 10, p. e42947, 2012.
- [37] M. R. Gupta, S. Bengio, and J. Weston, "Training highly multiclass classifiers.," *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1461–1492, 2014.
- [38] T.-N. Do, "Parallel multiclass stochastic gradient descent algorithms for classifying million images with very-high-dimensional signatures into thousands classes," Vietnam Journal of Computer Science, vol. 1, no. 2, pp. 107–115, 2014.
- [39] X. Han and A. C. Berg, "Dcmsvm: Distributed parallel training for single-machine multiclass classifiers," in Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on, pp. 3554–3561, IEEE, 2012.
- [40] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," Foundations and Trends® in Machine Learning, vol. 3, no. 1, pp. 1–122, 2011.
- [41] A. Agarwal, O. Chapelle, M. Dudík, and J. Langford, "A reliable effective terascale linear learning system," CoRR, vol. abs/1110.4198, 2011.
- [42] S. Gopal and Y. Yang, "Distributed training of large-scale logistic models.," in ICML (2), pp. 289–297, 2013.
- [43] Y. Prabhu and M. Varma, "Fastxml: A fast, accurate and stable tree-classifier for extreme multi-label learning," in *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 263–272, ACM, 2014.
- [44] K. Crammer and Y. Singer, "On the learnability and design of output codes for multiclass problems," Machine Learning, vol. 47, no. 2, pp. 201–233, 2002.

- [45] S. Van Der Walt, S. C. Colbert, and G. Varoquaux, "The numpy array: a structure for efficient numerical computation," *Computing in Science & Engineering*, vol. 13, no. 2, pp. 22–30, 2011.
- [46] S. Behnel, R. Bradshaw, C. Citro, L. Dalcin, D. S. Seljebotn, and K. Smith, "Cython: The best of both worlds," Computing in Science & Engineering, vol. 13, no. 2, pp. 31–39, 2011.
- [47] L. D. Dalcin, R. R. Paz, P. A. Kler, and A. Cosimo, "Parallel distributed computing using python," Advances in Water Resources, vol. 34, no. 9, pp. 1124– 1139, 2011.

Supplementary Material

A Derivation of Lagrandgian Dual Problems

A.1 Lin, Lee, and Wahba

Using slack variables, the primal LLW problem reads

$$\min_{W} \sum_{c=1}^{C} \left[\frac{1}{2} ||w_c||^2 + C \sum_{i:y_i \neq c} \xi_{i,c} \right]$$
s.t.
$$\sum_{c} w_c = 0$$

$$\forall i: \begin{cases} \xi_{i,c} \ge 1 + w_c^T x_i \\ \forall c \ne y_i: & \xi_{i,c} \ge 0. \end{cases}$$

We introduce Lagrangian multipliers $\alpha \in \mathbb{R}^{n \times C}$, $\beta \in \mathbb{R}^n$, and $\bar{w} \in \mathbb{R}^d$ with $\alpha_{i,c}, \beta_i \geq 0$. The Lagrangian is given in Figure A.1.

Slater's condition holds; therefore, we have strong duality and can use the dual

$$\begin{aligned} & \max_{\alpha,\beta,\bar{\alpha}} \min_{W,\xi} & L(W,\xi,\alpha,\beta,\bar{w}) \\ & \text{s.t.} & \forall i \ \forall c: \ \alpha_{i.c},\beta_{i.c} \geq 0. \end{aligned}$$

The partial derivatives are given by

$$\frac{\partial}{\partial \xi_{i,c}} L(W, \xi, \alpha, \beta, \bar{w}) = C - \alpha_{i,c} - \beta_{i,c}$$
$$\frac{\partial}{\partial w_c} L(W, \xi, \alpha, \beta, \bar{w}) = w_c + \sum_{i: y_i \neq c} \alpha_{i,c} x_i + \bar{w}.$$

Setting those to zero leads to

$$\begin{aligned} \forall\,i\,\forall\,c: &\quad 0 \leq \alpha_{i,c} \leq C \\ &\quad w_c = -\sum_{i:y_i \neq c} \alpha_{i,c} x_i + \bar{w} \\ &\quad = -X\alpha_c + \bar{w}. \end{aligned}$$

And plugging in into the lagrangian, finally gives the dual

$$\max_{\alpha \in \mathbb{R}^{n \times \mathcal{C}}, \bar{w} \in \mathbb{R}^d} \quad \sum_{c=1}^{\mathcal{C}} \left[-\frac{1}{2} || -X\alpha_c + \bar{w}||^2 + \sum_{i: y_i \neq c} \alpha_{i,c} \right]$$

$$\forall i: \begin{array}{l} \alpha_{i,y_i} = 0 \\ \forall c \neq y_i: \ 0 \leq \alpha_{i,c} \leq C. \end{array}$$
(A.1)

B Experiments

B.1 Setup

For our experiments we use two different types of machines. Type A has 20 physical cpu cores, 128 GB of memory and a 10 GigaBit Ethernet network. Type B has 24 physical cpu cores and 386 GB of memory. On type B we ran the experiments involving CS due to the memory requirements.

Training repetitions were run on training sets with a random order of the data (note that the training set is the same in each run; only the order of points is shuffled, which can impact the DBCA algorithm). For *LIBLINEAR* solvers we use the newest available version as of April 2016 with the default settings.

We implemented our solveres using OpenMP, Open-MPI, and the Python-ecosystem. In more detailed we used [45], [46], and [47].

B.2 Validation

The following Table B.1 contains the complete results of the validation experiment in Section 4.1.

B.3 LSHTC F1-Scores

The following Table B.2 contains the F1-Scores achieved by the solvers on the LSHTC1 datasets.

C Algorithms

The following routine complement the mainalgorithms in the paper.

Algorithm 5 Solving the graph 1-factorization problem. Indices start with one.

```
1: function MatchClass(C,c,r)
      if C is even and c = C then
2:
3:
         return r
4:
      if c = r then
         if C is even then
5:
6:
            return C
7:
         else
8:
            return c
      return mod(2r-c, C-1)
```

$$L(W, \xi, \alpha, \beta, \bar{w}) = \sum_{c=1}^{C} \left[\frac{1}{2} ||w_c||^2 + \sum_{i: y_i \neq c} \left(C\xi_{i,c} + \alpha_{i,c} (1 + w_c^T x_i - \xi_{i,c}) - \beta_{i,c} \xi_{i,c} \right) \right] - \bar{w}^T (\sum_c w_c).$$

Figure A.1: Lagrangian for LLW.

Table B.1: Error on the test set and density in % of the Shark solver (denoted S) and the proposed solver (denoted D).

	D-LLW		S-L	LW	D-WW		S-WW	
Dataset:	Err.	Den.	Err.	Den.	Err.	Den.	Err.	Den.
SensIT v.(com.)								
$\log(C)$: -1	21.34	100.0	21.34	100.0	19.88	100.0	19.88	100.0
0	20.95	100.0	20.95	100.0	19.51	100.0	19.51	100.0
1	20.78	100.0	20.78	100.0	19.38	100.0	19.38	100.0
glass								
$\log(C)$: -1	66.67	100.0	66.67	100.0	38.10	100.0	38.10	100.0
0	61.90	100.0	61.90	100.0	19.05	100.0	19.05	100.0
1	33.33	100.0	33.33	100.0	19.05	100.0	19.05	100.0
iris								
$\log(C)$: -1	13.33	100.0	13.33	100.0	6.67	100.0	6.67	100.0
0	26.67	100.0	26.67	100.0	13.33	100.0	13.33	100.0
1	26.67	100.0	26.67	100.0	13.33	100.0	13.33	100.0
letter								
$\log(C)$: -1	87.04	100.0	87.04	100.0	28.25	100.0	28.26	100.0
0	87.24	100.0	87.24	100.0	29.04	100.0	29.03	100.0
1	61.91	100.0	87.24	100.0	28.92	100.0	28.93	100.0
news20								
$\log(C)$: -1	29.23	97.24	29.23	97.24	15.32	51.16	15.30	49.72
0	22.97	97.24	22.97	97.24	14.80	44.74	14.80	42.70
1	16.15	97.17	16.15	97.04	15.98	45.97	15.98	43.47
rcv1	47.00	7 0.00	47.00	7 0.00	11.01	00.40	11.01	00.45
$\log(C)$: -1	47.96	78.00	47.96	78.00	11.31	26.42	11.31	23.45
0	33.27	78.00	33.41	77.98	11.52	22.93	11.52	20.12
1	12.03	78.00	12.03	77.98	12.03	23.05	12.03	20.06
satimage	00.75	100.0	00.70	100.0	15 00	100.0	15 00	100.0
$\log(C)$: -1	26.75	100.0	26.73	100.0	15.80	100.0	15.80	100.0
0	26.80	100.0	26.80	100.0	15.47	100.0	15.53	100.0
splice	26.90	100.0	26.90	100.0	15.96	100.0	16.00	100.0
$\log(C)$: -1	16.29	100.0	16.37	100.0	16.16	100.0	16.16	100.0
$\log(C)$: -1	16.29 16.09	100.0	16.37 16.15	100.0	16.16 16.37	100.0	16.16 16.28	100.0
1	16.09 16.34	100.0	16.13 16.28	100.0	16.37 16.32	100.0 100.0	16.28 16.24	100.0
usps	10.54	100.0	10.28	100.0	10.52	100.0	10.24	100.0
$\log(C)$: -1	31.84	100.0	31.84	100.0	8.17	100.0	8.17	100.0
$0 \log(C)$.	30.09	100.0	30.04	100.0	9.37	100.0	9.37	100.0
1	28.00	100.0	28.00	100.0	10.51	100.0	10.51	100.0
	20.00	100.0	20.00	100.0	10.01	100.0	10.01	100.0

Table B.2: Micro-F1 and Macro-F1 scores (in %) as achieved by the OVR, CS, WW, and LLW solvers on the LSHTC datasets. For each solver and each metric the best result across C values is in bold font. For LLW entries with a '*' did not converge within a day of runtime.

	Micro-F1				$Macro ext{-}F1$			
Dataset:	OVR	\mathbf{CS}	$\mathbf{W}\mathbf{W}$	\mathbf{LLW}	OVR	\mathbf{CS}	$\mathbf{W}\mathbf{W}$	\mathbf{LLW}
LSHTC-small								
$\log(C)$: -3	7.00	40.26	27.18	7.00	0.61	22.08	10.73	0.61
-2	14.42	40.26	34.66	7.00	2.70	22.08	16.15	0.61
-1	25.46	40.26	42.41	7.00	8.72	22.08	24.71	0.61
0	35.47	44.46	45.43	7.00	16.42	26.70	28.75	0.61
1	42.41	45.48	45.59	7.00	25.09	28.73	29.15	0.61
LSHTC-large								
$\log(C)$: -3	11.77	41.35	33.53	4.14	0.88	25.43	15.05	0.09
-2	14.80	41.52	39.42	4.14	1.51	25.41	20.83	0.09
-1	22.02	42.19	44.72	4.14	3.35	25.83	27.90	0.09
0	36.86	46.41	46.02	*	14.76	30.99	31.29	*
1	42.80	45.83	45.59	*	25.87	31.13	31.12	*
LSHTC-2012								
$\log(C)$: -3	16.34	50.19	41.98	7.37	0.28	20.55	8.08	0.01
-2	24.85	50.35	49.80	7.37	0.69	20.72	16.17	0.01
-1	39.62	53.86	55.06	7.37	2.64	23.76	25.94	0.01
0	52.67	57.33	55.99	*	12.46	32.57	32.06	*
1	53.17	54.40	53.85	*	24.41	31.84	30.95	*
LSHTC-2011								
$\log(C)$: -3	12.05	40.91	31.81	3.82	0.46	22.44	10.47	0.05
-2	14.15	40.91	37.86	3.82	0.62	22.46	16.48	0.05
-1	23.22	41.82	42.69	3.82	1.89	23.37	23.17	0.05
0	36.89	44.42	43.06	*	10.60	26.97	27.25	*
1	39.99	42.22	41.86	*	21.30	26.31	26.97	*