

Design Development and Performance Analysis of Distributed Least Square Twin Support Vector Machine for Binary Classification

Bakshi Rohit Prasad^{1*}, Sonali Agarwal²

^{1,2}Department of Information Technology, Indian Institute of Information Technology, Allahabad, India,

¹ORCID iD: <https://orcid.org/0000-0003-1921-9494>

²ORCID iD: <https://orcid.org/0000-0001-9083-5033>

Received: .201

Accepted/Published Online: .201

Final Version: ..201

Abstract: Machine Learning (ML) on Big Data has gone beyond the capacity of traditional machines and technologies. ML for large scale datasets is current focus of researchers. Most of the ML algorithms primarily suffer from memory constraints, complex computation and scalability issues. Least Square Twin Support Vector Machine (LSTSVM) technique is an extended version of Support Vector Machine (SVM). It is much faster as compared to SVM and is widely used for classification task. However, when applied to large scale datasets having millions or billions of samples and/or large number of classes, it causes computational and storage bottlenecks. This paper proposes a novel scalable design for LSTSVM named Distributed LSTSVM (DLSTSVM). This design exploits distributed computation on cluster of machines to provide scalable solution to LSTSVM. Very large datasets are partitioned and distributed in the form of Resilient Distributed Datasets on top of Spark cluster computing engine. LSTSVM is trained to generate two non-parallel hyper-planes. These hyper-planes are achieved by solving two systems of linear equations each of which involves data instances from either class. While designing DLSTSVM we employed distributed matrix operations using MapReduce paradigm of computing to distribute the tasks over multiple machines in the cluster. Thus, memory constraints with extremely large datasets are averted. Experimental results show the reduction in time complexity as compared to existing scalable solutions to SVM and its variants. Moreover, detailed experiments depict the scalability of the proposed design with respect to large datasets.

Key words: Distributed machine learning, Big Data, Cluster computing, LSTSVM, MapReduce, Parallel processing

1. Introduction

Existing ML techniques may be classified into two broad categories; supervised ML techniques and unsupervised ML techniques [1]. In supervised ML, classification rules are built on a training set which contains data instance with class labels. Then after, new data instances are classified using the learnt classification rules. Support Vector Machine (SVM) is an extensively used supervised ML algorithm developed by Vapnik et al. [1]. SVM generates two parallel hyper-planes during training. These hyper-planes separate data from two classes. The optimal hyper-plane is generated by solving a complex Quadratic Programming Problem (QPP). The overall training takes order of $O(n^3)$ time if there are n training data samples. A QPP size directly depends on the dataset size. Moreover, solving a QPP is computationally expensive job. Another faster variant of SVM has been proposed which solves two smaller QPPs unlike SVM. This technique is known as Twin Support Vector Machines (TWSVM) [2]. This causes significant reduction in time complexity involved in SVM training and speeds up the training four times faster than that of SVM. Though, the QPPs are smaller, solving the same is

*Correspondence: rs151@iiita.ac.in

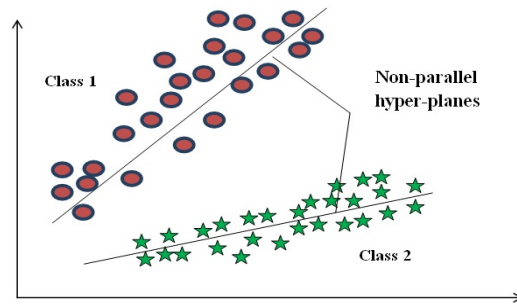


Figure 1: Hyper-planes generated by LSTSVM

1 still a complex task. A new faster version of TWSVM has been proposed developed by Kumar et al. [3] and
 2 is well known as LSTSVM which is widely used for the classification task. Unlike TWSVM, the LSTSVM is
 3 trained to generate two non-parallel hyper-planes. These hyper-planes are achieved by solving two systems of
 4 linear equations each of which involves data instances from either class. Thus, it is significantly faster than
 5 TWSVM. Figure 1 represents the categorization of two classes by using LSTSVM. As shown in the Figure 1
 6 there are two classes; class1 and class2, which are divided by using two non-parallel planes in such a way that
 7 each plane is nearer to the data points of one class while farther from the other class. Although LSTSVM is
 8 much faster as compared to TWSVM, yet very large datasets cause computation and storage bottlenecks in
 9 case of LSTSVM. Due to this fact, capability of LSTSVM technique is limited to smaller dataset only. For very
 10 large datasets, training becomes too expensive and system may go out of memory as well.

11 1.1. Research Contribution

12 To solve the underlying computation and storage challenges of LSTSVM on very large datasets, the proposed
 13 technique named DLSTSVM employs distributed computing over multiple machines in a parallel fashion on top
 14 of recent distributed parallel computing framework. Our research work makes following contribution:

- 15 • The way all computations are distributed over multiple machines makes DLSTSVM a scalable technique
 16 which is the prime objective of current research work.
- 17 • DLSTSVM not only handles scales well even with large number of input patterns, due to data parallel
 18 computing, it achieves significant speed-ups too.
- 19 • This work gives the designing of DLSTSVM based on MapReduce paradigm of computing along with
 20 detailed performance and scalability analysis.

21 Present research work is illustrated in six major sections. Section 1 gives a brief introduction to supervised
 22 ML techniques along with their underlying challenges in case of very large datasets. In section 2, state-of-
 23 the-art is presented around distributed frameworks for Big Data machine learning. Also, multiple variants of
 24 SVM, LSTSVM and distributed SVM are discussed along with their inherent challenges. Section 3 gives a
 25 brief background of mathematical preliminaries of LSTSVM as well as basic fundamentals of Spark framework.
 26 Section 4 describes the design of DLSTSVM algorithm and discusses its time and storage complexity. Detailed
 27 experimental results in section 5 compare the performance of the proposed algorithm with earlier techniques.
 28 Moreover, to assess the scalability of DLSTSVM, manifold experiments are performed on large-scale synthetically
 29 generated datasets. Finally, section 6 concludes the presented work with probable future prospects.

2. Literature Survey

Related research works are explored in three major dimensions; first is related to distributed frameworks and libraries for handling Big Data, second is related to multiple single machine variants of SVM and LSTSVM, and third is related to scalable solutions to SVM and its variants along with their limitations.

2.1. Distributed frameworks and libraries for Big Data

To utilize the advantages of distributed computing to enable machine learning systems to deal with large-scale datasets multiple distributed ML frameworks and libraries are developed. Apache Mahout includes distributed scalable ML algorithms implemented for Hadoop and provide free scalable ML libraries. GraphLab implements Machine Learning-Data Mining in the cloud-based environment for graph related operations in parallel. Microsoft's DryadLINQ [5] using LINQ on top of Dryad, is well suited for data parallel applications. Spark [6], a distributed computing engine, well suited for iterative nature of ML and facilitates vivid ML tasks. Multiple recent survey [7], discusses and compares pros and cons of several computing engines such as Flink, Spark, H2O, and Storm. Also, they compared ML libraries inside these engines like Mahout, MLlib, etc.[29]

2.2. SVM and LSTSVM variants

SVM and its various variants for binary class classifications are available. Several researchers extended the basic SVM and LSTSVM to multi-class scenarios too [9]. Multiple variants of SVMs have been proposed out of which some are exact solvers [11] [13], some fall under the categories of approximate solvers [14]-[17], divide and conquer SVM solvers [10] online SVM solvers [12]. Accurate solvers are time consuming whereas approximate solvers compromise with accuracy. Although, divide and conquer approach try to give faster solution but with large scale data instances they still exhibit high training time. In previous couple of years some more non-parallel variants of binary class LSTSVM have been proposed using Fuzzy Logic, KNN and Entropy based techniques [18]-[21] but with million scale data instances their computation time is extremely high.

2.3. Distributed and Parallel Machine Learning Algorithms

Parallel machine learning is an aspect to scale up existing sequential ML algorithms [28]. Advancements in the field of multi-core and cloud computing architectures have caused wide accessibility to distributed and parallel computing systems [8]. Faster variants of SVM (including LSTSVM) too become incapable of handling very large-scale datasets. Several solutions have been tried to scale SVM using various approaches like dividing the dataset, training SVM on each subset in parallel, and then combining multiple classifiers, thereby iterating this process [32] [34]. Some used multiple re-partitioning of data which involves sharing of training data among nodes, thereby causing significant training overhead. Some required frequent communication of sub-solutions resulting in slow training time. Thus, drawback of these approaches is that reallocation process slows down the training process. Approach given in a research work [13] used cascading of support vectors in multiple stages of training till final set of support vectors is obtained which involves significant data transfer cost. Development of distributed computation on commodity nodes and MapReduce kind of frameworks, the problems of aforementioned techniques are tried to resolve [35]. Some researchers employed high-speed costly Graphics Processing Units (GPU) and applied MapReduce based technique. The major drawback of this approach is in tedious setting up specialized configurations which makes it difficult to use and develop applications using this approach [28]. Apart from state-of-the-art parallel SVM methods [22] [23] some recent parallel and distributed

1 SVM techniques are also proposed in previous couple of years [24]-[26]. One of these approaches provides a
 2 distributed and parallel mechanism for alternating direction method of multipliers for non-convex penalized
 3 SVMs [24] at cost of accuracy. Another technique performs subspace partitioning on the datasets by using
 4 decision tree on projection of data showing maximum variance [25] and shows approximately 150 times speedup
 5 over sequential SVM. One group of researchers proposed a parallelization strategy in training SVM which finds a
 6 low rank approximation of matrices and uses random projection mechanism [26]. This solves the approximated
 7 optimization problem and appears as scalable approach still, its run-time is very high for million scale datasets.

8 **3. Fundamental Concepts of LSTSVM and Spark**

9 This section specifies the basic formulation of LSTSVM and its incompetence to handle the large scale datasets
 10 with respect to storage and computation. Further, it briefs the basic description of Spark and its features.

11 **3.1. Mathematical Formulation of LSTSVM**

12 As discussed earlier, LSTSVM finds two hyper-planes. These hyper-planes are not necessary to be parallel. The
 13 non-parallel hyperplanes are obtained by following two systems of linear equations given by 1 and 2.

$$14 \quad \min(w_1, b_1, \xi) \frac{1}{2} \|X_1 w_1 + e_1 b_1\|^2 + \frac{c_1}{2} \xi^T \xi \left[(s.t.) - (X_2 w_1 + e_2 b_1) + \xi = e_2 \right] \quad (1)$$

$$15 \quad \min(w_2, b_2, \eta) \frac{1}{2} \|X_2 w_2 + e_2 b_2\|^2 + \frac{c_2}{2} \eta^T \eta \left[(s.t.) (X_1 w_2 + e_1 b_2) + \eta = e_1 \right] \quad (2)$$

16 Here, X_1 and X_2 are data matrices corresponding to class -1 and class +1 respectively and have total l_1 and
 17 l_2 instances. Here, w_1, b_1 and w_2, b_2 are respective hyper-plane parameters know as weights and bias. Here,
 18 $e \in R^{l_1}$ and $e \in R^{l_2}$ are the vectors with entries of 1 only, c_1 and c_2 are non-negative penalty parameters
 19 whereas $\xi \in R^{l_2}$ and $\eta \in R^{l_1}$ are known as slack variables for class -1 and class +1. Thus, optimization
 20 equations 1 and 2, for LSTSVM contain only equality constraints. Finally, upon solving these equations,
 hyper-plane parameters are attained as given in equation 3 and 4. Here, $G = [X_1 e_1]$ and $H = [X_2 e_2]$.

$$21 \quad \begin{bmatrix} w_1 \\ b_1 \end{bmatrix} = - \left(H^T H + \frac{1}{c_1} G^T G \right)^{-1} H^T e_2 \quad (3)$$

$$22 \quad \begin{bmatrix} w_2 \\ b_2 \end{bmatrix} = - \left(G^T G + \frac{1}{c_2} H^T H \right)^{-1} G^T e_1 \quad (4)$$

23 With the help above hyper-plane parameters, the obtained hyper-planes are given as per equation 5 and 6.

$$24 \quad X^T w_1 + b_1 = 0 \quad (5)$$

$$X^T w_2 + b_2 = 0 \quad (6)$$

When a new data instance arrives, it is assigned a class based on the decision function give in equation 7.

$$f(x) = \arg \min_{i=1,2} \frac{|w_i \cdot x + b_i|}{\|w_i\|} \quad (7)$$

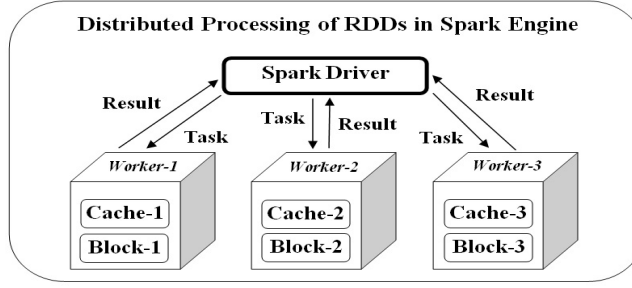


Figure 2: Spark distributed processing

1 With very large datasets, computing the equations 3 and 4 causes computation and storage bottlenecks, thereby
 2 training becomes too expensive and system may go out of memory as well. Thus, capability of LSTSVM
 3 technique is limited to smaller datasets only. Next sub-section briefs the characteristics of Spark distributed
 4 processing engine used in our research work for distributed parallel computation over cluster of machines.

5 3.2. Apache Spark - A Cluster Computing Framework

6 Apache Spark provides cluster computing which encompasses automatic locality-aware scheduling, load balanc-
 7 ing, and fault tolerance. Spark facilitates data-parallel computations on cluster of commodity machines using
 8 MapReduce computing paradigm [4]. Spark is also well suited for iterative computing and interactive analytics.
 9 As shown in Figure 2, Spark Engine Core takes care of prime functions. These functions cover Spark computing
 10 environment set up, generation, transformation and distribution of Resilient Distributed Datasets RDDs. Spark
 11 driver internally distributes RDDs onto multiple machines in the cluster. It can be explicitly cached in main
 12 memory for reuse if required in an application [6].

13 4. Proposed Methodology

14 Design of DLSTSVM includes three phases; First decomposes entire computation into different modules and
 15 identifies the parallel sections, Second specifies the Map and Reduce steps for each module and Third discusses
 16 time and space complexities of DLSTSVM to show how the challenges of LSTSVM is resolved by DLSTSVM.

17 4.1. Problem decomposition

18 Consider equation 3 and 4 which yields two hyper-planes each corresponding to data instances from a single
 19 class. Following sub-sections identifies parallel computable parts involved in these series of matrix calculations.

20
 21 *(i) Identification of parallel computations:* Equation 3 involves series of matrix operations. Entire
 22 computation is divided into two parts specified by equation 8 and 9. Each part is computed in parallel on the
 23 data residing on each machine. The serial and parallel portion of computation is clearly depicted in the flow
 24 diagram of Figure 3. Thus, the hyper-plane parameter u can be given as $u = R_{c \times c} \times S_{c \times 1}$ such that:

$$R = -((H^T H) + (1/c)(G^T G))^{-1} \quad (8)$$

$$S = H^T e_2 \quad (9)$$

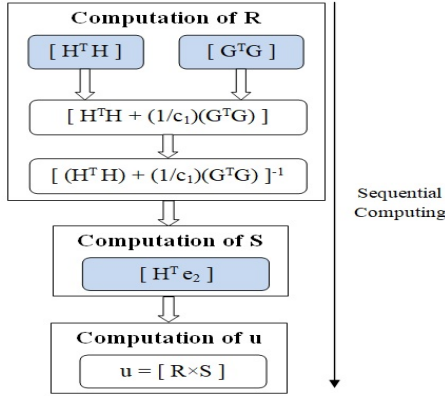


Figure 3: Computation flow of parallel (shaded block) and serial sections

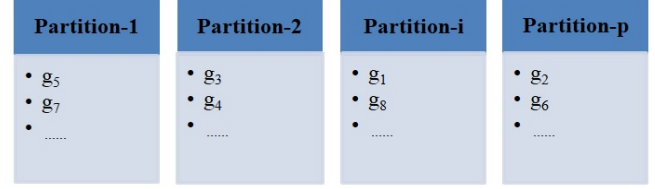


Figure 4: Data partitioning over cluster of machines

1 **(ii) Data Representation and Partitioning:** Matrix R requires computation of $G^T G$ and $H^T H$. For very
 2 large matrix G , computation of $G^T G$ needs to be distributed in such a way that may avert memory bottleneck as
 3 well as speedup the processing. Therefore, matrix G is distributed as RDD over cluster where each element of
 4 RDD represents a data instance from matrix G in mutually exclusive fashion. Figure 4 depicts this partitioning
 5 of matrix $G_{r \times c}$ having r rows (representing data instances) and c columns (representing dimensions of each
 6 data instance). Here, g_i specifies i^{th} data instance of matrix G . Then, $G^T = [g_1^T, g_2^T, g_3^T \dots g_r^T]_{c \times r}$, where
 7 each g_i^T is a column vector achieved by taking the transpose of corresponding row vector g_i .

8 4.2. Modeling parallel computations as MAP-REDUCE steps

9 **(i) For computing matrix R :** MAP-REDUCE steps for $G^T G$ are modeled in following sub-sections:

10 **MAP step:** Figure 5 shows the MAP step which multiplies j^{th} column of G^T to j^{th} row of G , i.e. $g_i^T \times g_i$.
 11 Here, g_i is a row vector of dimension $1 \times c$ and g_i^T is a column vector of dimension $c \times 1$. In each partition- i ,
 12 this MAP step is executed in parallel for each data instance g_j . Thus, $t_j = g_j^T \times g_j$ is a matrix of dimension
 13 $c \times c$, generated for very instance of matrix G .

14 **REDUCE step:** If each machine in cluster has m data instances then it generates m intermediate matrices.
 15 Since, matrix G has r such data instances, therefore total r intermediate matrices throughout the cluster.
 16 REDUCE step is designed to sum up all the temporary matrices t_j to produce the matrix $G^T G$ as $T = \sum_{j=1}^r t_j$
 17 i.e. all t_j are reduced on sum operation for corresponding elements at $index(i, j)$. A simple matrix addition
 18 is applied on two intermediate matrices at a time in each partition. Another intermediate matrix is added to
 19 the sum matrix produced in the previous step. Cascading the sum, matrix T is computed as shown in Figure
 20 6. Similarly, $H^T H$ is also computed in parallel fashion, added to $(1/c_1)G^T G$ and inverted to yield matrix R .

21
 22 **(ii) For computing $S = H^T e_2$:** For matrix H which is distributed over cluster, it has been evaluated
 23 that vector obtained by multiplying H^T to e_2 and the vector obtained by transposing the row-vector achieved
 24 after column-wise addition of elements of matrix H are same. Thus, MAP and REDUCE steps are designed as:
 25 **MAP step:** For the approach proposed in the work, this step does not make any changes to underlying matrix

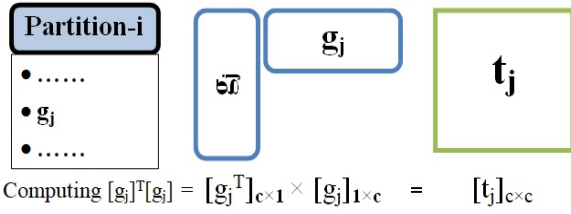


Figure 5: MAP step applied to each data instance

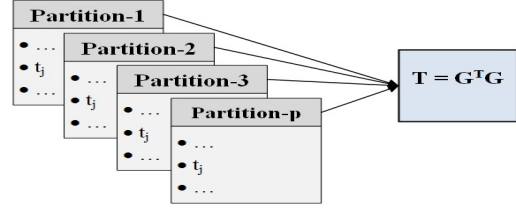


Figure 6: REDUCE step for summing all intermediate matrices

- 1 elements. It simply forwards the RDD of matrix H for reduction by the following REDUCER.
 - 2 **REDUCE step:** In reduce step, the column-vectors of the RDD corresponding to matrix H are reduced by
 - 3 sum operation defined on the $index(j)$. Here, j indicates a particular column and $1 \leq j \leq c$. This
 - 4 reduction is applied in parallel on each machine in the cluster in order to compute $H^T e_2$ in distributed parallel
 - 5 fashion. Internally, two row-vectors of H are picked up at a time in each partition and their corresponding
 - 6 column elements are added, resulting in new sum-vector. Next row-vector is added to the sum-vector produced
 - 7 in the previous step. Thus, all row-vectors on multiple machines are summed distributively to give vector S .
- Algorithm 1 specifies the required sequence of steps involved in DLSTSVM.

Algorithm 1 DLSTSVM

- 1: **procedure** SEGREGATE THE DATA OF CLASS-A AND CLASS-B AND TRANSFORM IT INTO RDD(as A, B)
 - 2: Define following two matrices: $G = [Ae_1]$ and $H = [Be_2]$ and two vectors as e_1, e_2
 - 3: Calculate matrix $R_1 = (H^T H + 1/c_1(G^T G))^{-1}$ and $R_2 = (G^T G + 1/c_2(H^T H))^{-1}$ as follows:
 - 4: a: Calculate $H^T H$ distributively as follows:
 - 5: MAP Phase: $t_1 = H.map(rowTrans_row)$
 - 6: REDUCE Phase: $t_1.reduce(sumMatrix)$
 - 7: b: Calculate $G^T G$ distributively:
 - 8: MAP Phase: $t_2 = G.map(rowTrans_row)$
 - 9: REDUCE Phase: $t_2.reduce(sumMatrix)$
 - 10: c: Compute $(H^T H + 1/c_1(G^T G))^{-1}$
 - 11: d: Compute $(G^T G + 1/c_2(H^T H))^{-1}$
 - 12: Calculate $S_1 = H^T e_2$ distributively as follows:
 - 13: MAP Phase: $d_1 = H.map(feedToReduce)$
 - 14: REDUCE Phase: $S_1 = d_1.reduce(sumRowVector)$
 - 15: Calculate $S_2 = G^T e_1$ distributively as follows:
 - 16: MAP Phase: $d_2 = G.map(feedToReduce)$
 - 17: REDUCE Phase: $S_2 = d_2.reduce(sumRowVector)$
 - 18: Find hyper-plane parameters as follows:
 - 19: $[W_1, b_1] = u = R_1 \times S_1$
 - 20: $[W_2, b_2] = u = R_2 \times S_2$
-

8

4.3. Analysis of time and space complexity

- 10 Consider matrix G be of order $r \times c$. Hence, the resultant matrix $T = G^T \times G$ will be of order $c \times c$. Since,
- 11 LSTSVM applies usual matrix multiplication hence, takes $O(rc^2)$ time and requires $O(2rc + c^2) = O(rc + c^2)$
- 12 storage, $O(2rc)$ for G^T and G and $O(c^2)$ for matrix T . Very large value of r causes memory bottlenecks.
- 13 $O(rc^2)$ computations is time-expensive too. Further, we discuss the time and space complexity of DLSTSVM:
- 14 (i) **Time and space complexity for R:** Since, $G^T \times G$ is being computed in distributed parallel way, let
- 15 there are n machines in the cluster. Computation of each intermediate matrix t_j requires c^2 products and c^2

1 storage. Thus, for a single machine $O(mc^2)$ time and $O(mc)$ space required which can be handled by a single
 2 machine. Here, $m = (r/n)$ specifies number of data instances lying on each machine. It can be stated that
 3 in DLSTSVM, $O(rc^2)$ computation is now distributed among n machines and each machine performs only
 4 $O(mc^2)$ computations. Also, storage is distributed, reducing to $O(c^2)$ only on a each machine.

5 Since, matrix $G^T G$ and $H^T H$ are of dimension $c \times c$, adding $(1/c_1(G^T G))$ and $H^T H$ takes $O(c^2)$ addition
 6 operations and $O(c^2)$ storage. Resulting addition matrix is of order $c \times c$, hence, taking inverse will take $O(c^3)$
 7 computation and $O(c^2)$ storage only. Thus, matrix R specified is computed successfully in $O(mc^2 + c^3 + c^2) =$
 8 $O(mc^2 + c^3)$ operations and $O(c^2 + c^2 + c^2) = O(c^2)$ storage on each machine.

9 *(ii) Time and space complexity for S:* As discussed in the REDUCE step of sub-section 4.2 to compute
 10 vector S , sum of corresponding elements from two subsequent vectors performs c additions. Since, each machine
 11 keeps m data instances, total $O((m-1) \times c) = O(mc)$ sum operations are executed on each machine, all
 12 machines running in parallel in cluster. Furthermore, each machine stores only two row-vectors to be added
 13 thereby requiring $O(c)$ space.

14 In summary, computation of R requires $O(mc^2 + c^3)$ time and $O(c^2)$ space whereas computation of S
 15 requires $O(mc)$ time and $O(c)$ space. Thus, matrix $u = R \times S$, is computed without any bottleneck.

16 5. Results and Discussion

17 5.1. Experimental Setup and Configurations

18 System's configuration used for our experiment is listed in Table 1 specifying processor, memory, and operating
 19 system. Since, experimental setup requires distributed cluster hence, Table 1 also specifies the cluster size in
 20 terms of nodes, memory, and processor as well as the version of Spark and SBT (Simple Build Tool) used here.

Table 1: System and cluster configuration

System Configuration	Parameter Value	Cluster Configuration	Parameter Value
Processor	Intel(R) Core-i3 3.30 GHz	SBT version	0.11.3
Total number of cores per computer system	4	Spark version	1.5.0
RAM per computer system	6 GB	Number of nodes in the cluster	8
Operating System	Ubuntu 12.04 Linux	Total number of cores in cluster	24 (3 cores from each node in cluster /SPARK_WORKER_CORES = 3)
		Total executor memory in cluster	32 GB (4 GB from each node in cluster /SPARK_WORKER_MEMORY= 4g)

21

22 5.2. Comparison with existing solutions

23 Effective libraries have been developed for SVM such as LIBSVM [31] and SVMLight. However, SVM and
 24 their optimized implementations in these libraries still do not scale efficiently even for 1 million data samples.
 25 Computation cost is too high and memory requirements are also high. Though some approximate solvers [32]
 26 [33] [34] are proposed, but they compromise with accuracy. For the sake of comparison, datasets that have been
 27 used are listed in Table 2. In addition, Table 2 shows a comparison of results of other existing techniques with
 28 DLSTSVM that clearly depicts the time efficiency of DLSTSVM over other state-of-the-art techniques. Here,
 29 Table 2 specifies the best performances achieved out of multiple executions of various techniques over a grid

Table 2: Performance comparison with earlier techniques

Dataset	Ijenn1	Census	Webspam	Kddcup99	Cifar
#Instances	141691	199523	350000	5209460	60000
#Attributes	22	409	254	125	1024
Parameter Setting	$C=32, \gamma=2$	$C=512, \gamma=2^{-9}$	$C=8, \gamma=32$	$C=256, \gamma=2^{-1}$	$C=8, \gamma=2^{-22}$
Techniques	Time(s)/Acc(%)	Time(s)/Acc(%)	Time(s)/Acc(%)	Time(s)/Acc(%)	Time(s)/Acc(%)
DC-SVM-Early [10]	12/98.35	261/94.9	670/99.13	470/92.61	1977/87.02
DC-SVM [10]	41/98.69	1051/94.2	10485/99.28	2739/92.59	16314/89.5
LIBSVM [11]	115/98.69	2920/94.2	29472/99.28	6580/92.51	42688/89.5
LaSVM [12]	251/98.57	3514/93.2	20342/99.25	6700/92.13	57204/88.19
CascadeSVM [13]	17.1/98.08	849/93.0	3515/98.1	1155/91.2	6148/86.8
LLSVM [14]	38/98.23	1212/92.8	2853/97.74	3015/91.5	9745/86.5
FastFood [15]	87/95.95	851/91.6	5563/96.47	2191/91.6	3357/80.3
SpSVM [16]	20/94.92	3121/90.4	6235/95.3	5124/90.5	21335/85.6
LTPU [17]	248/96.64	1695/92.0	4005/96.12	5100/92.1	17418/85.3
	$c1=c2=10^{-1}$	$c1=c2=10^{-3}$	$c1=c2=10^{-3}$	$c1=c2=10^{-1}$	$c1=c2=10^{-2}$
LSTSVM [3]	743.65/98.62	31716.24/95.2	23041.36/99.15	*/NA	*/NA
DLSTSVM	10/98.62	238/95.2	121/99.15	190/94.4	5245/91.3

of points corresponding to a range of values for balancing parameters $C = [2^{-10}, 2^{-9}, \dots, 2^{10}]$ and kernel parameter $\gamma = [2^{-10}, 2^{-9}, \dots, 2^{10}]$ for RBF kernel. Moreover, the stopping criterion is set to the default value 10^{-3} in case of LIBSVM and DC-SVM. Rank, *#clusters* and *#features* (corresponding to LLSVM, LTPU and FastFood respectively) is taken as 3000. However, for DLSTSVM, value of $c1=c2$ has been set at the same value in the range $[10^{-5}, 10^{-4}, \dots, 10^5]$, on which LSTSVM gives the best performance. It is evident from the results that at same parameter setting, the predictive performance of LSTSVM and DLSTSVM is same. Another set of experiments are performed to compare the run-time and predictive performance of DLSTSVM with state-of-art recent variants of LSTSVM. Results are calculated on bench-marking large scale datasets like CovType and Mnist as well as David's NDC Data generator [27]. Table 3 lists the properties of these datasets, parameter settings used in experiments and performance comparison of various techniques. It is clearly observed that DLSTSVM outperforms all the non-parallel state-of-art LSTSVM and its variant techniques in terms of run-time performance and at the same time DLSTSVM exhibits comparable predictive performance too with state-of-the-art techniques.

Moreover, further experiments are conducted to compare DLSTSVM with some existing state-of-the-art distributed parallel techniques too. Benchmarking datasets like adults, gisette, covType, Mnist are used in our experiments. Table 4 specifies the properties of these datasets, parameter settings used in experiments and performance comparison of various techniques. Results clearly show better or comparable predictive performance of DLSTSVM with most of the other techniques. At the same time, run-time performance is much better than Parallel SVM, ADMM Parallel, xSVM and Parallel Distributed Penalized SVM techniques. Though, Projection-SVM shows better run-time here for large scale feature vector size however, since, DLSTSVM is scalable in nature, there is scope of adding more number of executors to reduce this computation time further.

5.3. Scalability Analysis of DLSTSVM

This section conducts experiments to test the scalability of DLSTSVM and observes its performance with respect to *#cores*, *#instances* and *#features* on synthetically generated datasets. Multiple datasets are generated with a different number of data instances; 1, 3, 5, 7, 9 and 11 million, a different number of attributes; 100, 300, 500, 700, 900. Run-times of DLSTSVM are estimated with varying number of cores, instances and features in following subsections.

Table 3: Performance comparison with state-of-art LSTSVM variants

Dataset	NDC-100k	NDC-500k	NDC-1m	CovType	Mnist
#Instances	100,000	500,000	1,000,000	464,810	2,000,000
#Attributes	32	32	32	54	784
Parameter Setting	$C=1, \gamma=1$	$C=1, \gamma=1$	$C=1, \gamma=1$	$C=32, \gamma=32$	$C=8, \gamma=2^{-3}$
Techniques	Time (s)/Acc.(%)	Time (s)/Acc.(%)	Time (s)/Acc.(%)	Time (s)/Acc.(%)	Time (s)/Acc.(%)
FLSTSVM [18]	21.67/87.48	89.5/86.24	209.28/86.10	242.34/94.74	9242.83/95.83
EFLSTSVM [19]	15.53/87.21	91.27/86.62	190.63/87.34	980.19/95.26	14855.62/96.18
ULSTPMSVM [20]	228.56/88.72	1144.95/87.13	2289.84/88.48	2082.46/95.42	*/NA
KNN-Based LSTSVM [21]	1919.5/87.69	119968.46/87.37	*/NA	*/NA	*/NA
	$c1=c2=10^{-1}$	$c1=c2=10^{-3}$	$c1=c2=10^{-3}$	$c1=c2=10^{-1}$	$c1=c2=10^{-2}$
LSTSVM [3]	421.76/88.14	2102.34/87.42	4263.78/87.58	15017.35/96.15	*/NA
DLSTSVM	6.31/88.14	30.22/87.42	63.47/87.58	91.26/96.15	4827.68/97.58

*represents Experiment stopped as computation time was very high, # represents Memory out of run

Table 4: Performance comparison with state-of-art distributed and parallel SVM techniques

Dataset	adult	gisette	CovType	Mnist
#Instances	32,561	6000	464,810	2,000,000
#Attributes	123	5000	54	784
Parameter Setting	$C=32, \gamma=2^{-7}$	$C=1, \gamma=2^{-4}$	$C=32, \gamma=32$	$C=1, \gamma=2^{-5}$
Techniques	Time (s)/Acc.(%)	Time (s)/Acc.(%)	Time (s)/Acc.(%)	Time (s)/Acc.(%)
Parallel SVM [22]	488.32/84.42	8241.47/97.56	5281.74/95.89	*/NA
ADMM Parallel [23]	59.87/85.67	1251.52/97.62	934.54/96.03	1855.12/97.21
Parallel and Distributed Penalized SVM [24]	33.51/84.31	745.20/96.43	567.44/96.24	1101.59/97.04
Projection-SVM [25]	58.86/84.83	1007.94/97.20	877.15/97.12	2986.21/97.59
xSVM [26]	39.11/85.14	712.23/97.41	543.56/96.03	4079.74/97.70
	$c1=c2=10^{-1}$	$c1=c2=10^{-3}$	$c1=c2=10^{-3}$	$c1=c2=10^{-2}$
DLSTSVM	2.87/85.74	4107.36/97.91	91.26/96.85	3827.68/97.58

1 **(i) Number of cores versus computation time:** Several graphs depicted in Figure 7 and 8 are drawn
2 between cores (at X-axis) and computation time (at Y-axis). Graphs are plotted for 1, 3, 5, 7, 9, and 11 million
3 data instances respectively. Multiple curves in each Figure 7a-8c shows the run-time for 100, 300, 500, 700 and
4 900 attributes. Run-time decreases with addition of executor-cores in the cluster. This decrement in time how-
5 ever decreases gradually because addition of more machines puts communication overheads, meta-information
6 management cost and task scheduling delays too. Also, more attributes results in more computation time.

7 **(ii) Number of data instances versus computation time:** Figure 9 and 10 analyzes the computation
8 time with increasing number of data instances where each graph corresponds to a fixed number of cores. It
9 establishes the scalability of DLSTSVM with respect to large-scale data instances. Machines in cluster do not
10 hang even for 11 million instances whereas LSTSVM is unable to handle even 1 million instances. Multiple
11 curves in each Figure 9a-9c and 10a-10c shows the run-time for 100, 300, 500, 700 and 900 attributes. It is
12 evident from Figure 9a to 10c that computation time of DLSTSVM behaves linearly with respect to increase in
13 data instances. However, for a higher number of attributes (≥ 300), the slope of the curve gets steeper.

14 **(iii) Number of attributes versus computation time:** Each graph in Figure 11 and 12 is plotted for
15 execution time against varying number of attributes and for a fixed value of data instances. Multiple curves
16 in each graph are drawn corresponding to different number of cores from among 24 cores available. All plots
17 show similar pattern of performance. It is also evident that greater the number of attributes, higher is the
18 computation time. It is also observed that computation time inflates as we go beyond 300 attributes which
19 supports the computational analysis for DLSTSVM given in subsection 4.2.

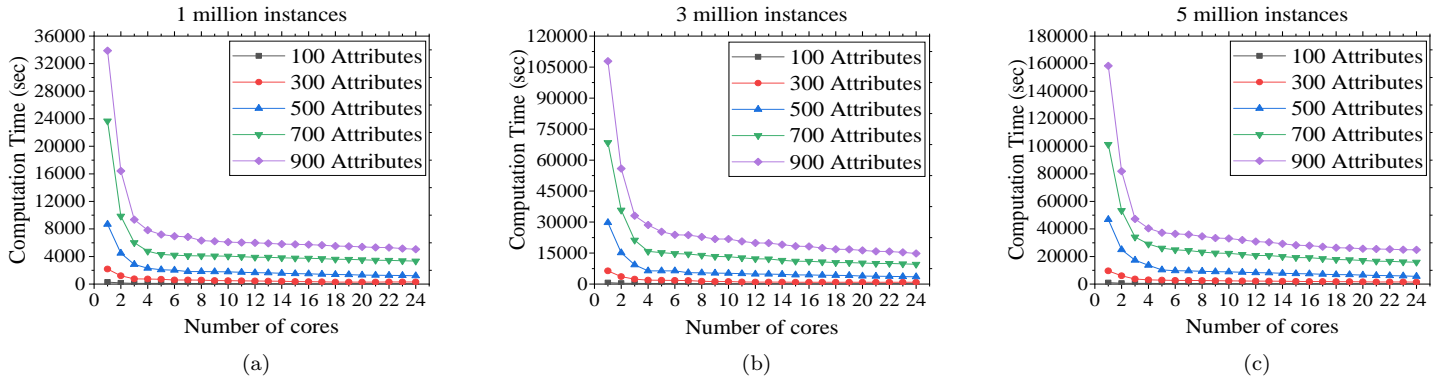


Figure 7: Run-time performance for different cores

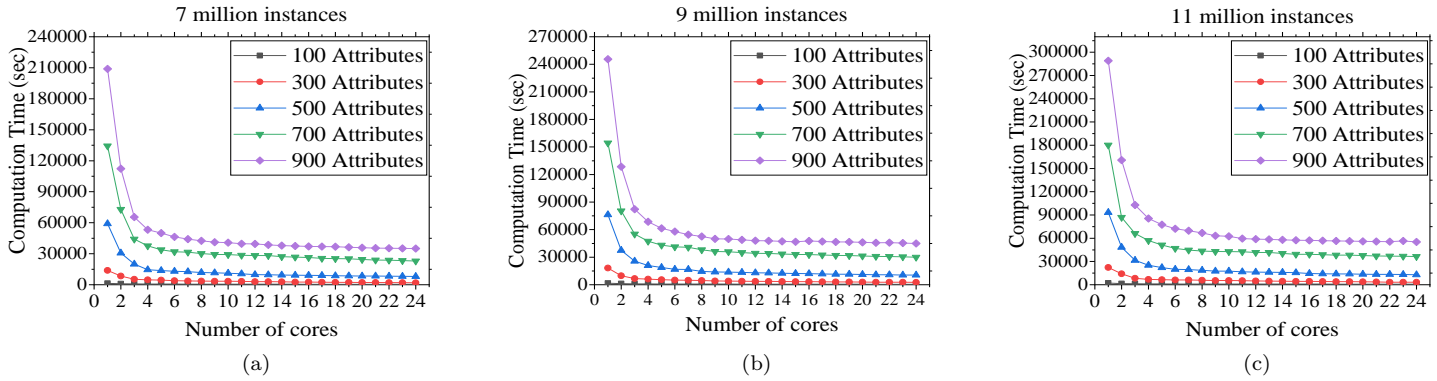


Figure 8: Run-time performance for different cores

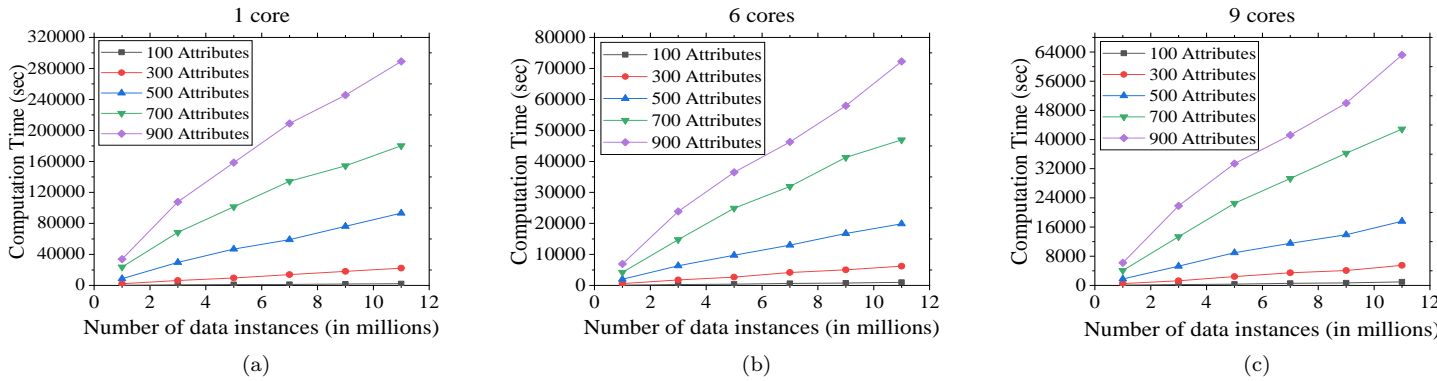


Figure 9: Run-time performance for different data instances

- 1 *(iv) 3S based performance evaluation of DLSTSV*
- 2 *M*: Performance of distributed parallel algorithms
- 3 can be assessed along 3S evaluators; Sizeup, Speedup and Scaleup. These parameters specify three different
- 4 aspects of a distributed parallel algorithm. Description and significance of each parameter along with the mea-
- 5 surements for DLSTSV

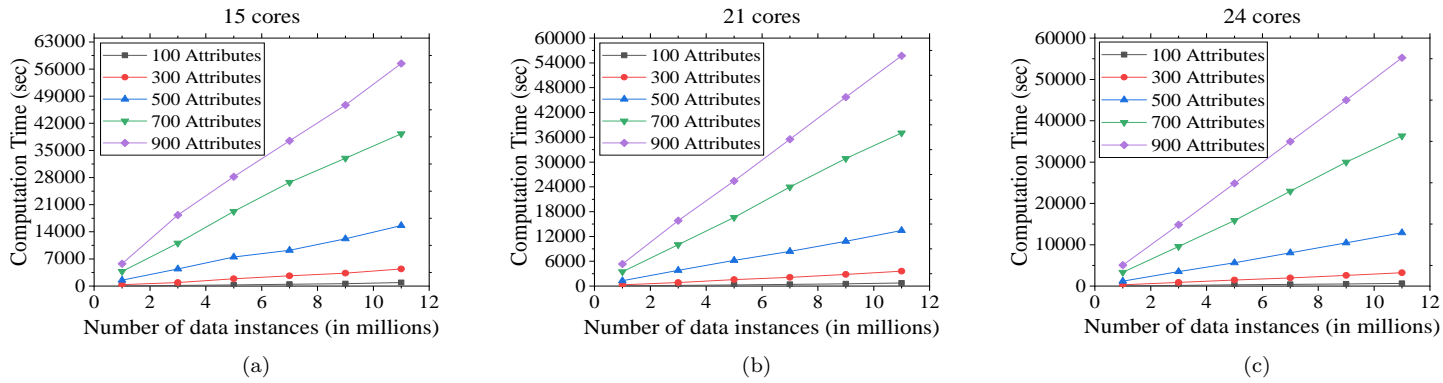


Figure 10: Run-time performance for different data instances

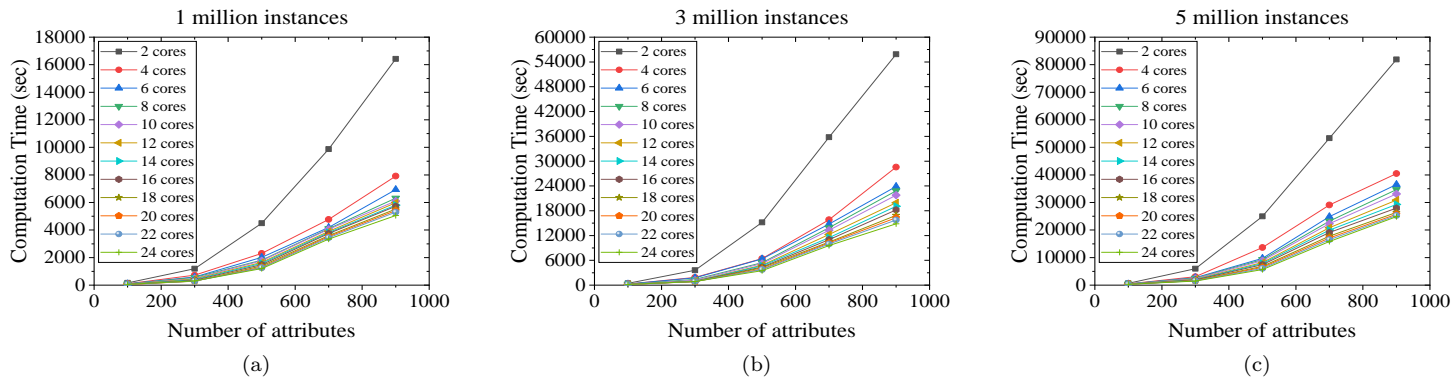


Figure 11: Run-time performance for different attributes

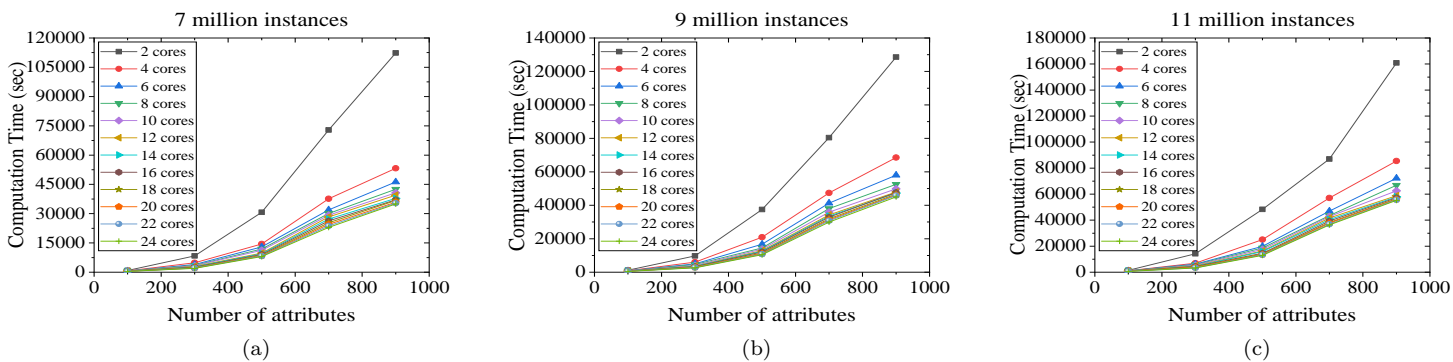


Figure 12: Run-time performance for different attributes

- 1 **Sizeup:-** Performance evaluation parameter Sizeup is a measure of computation time of the algorithm with
- 2 respect to increment in input size keeping number of computation unit fixed. For ideal situation graph, depict-
- 3 ing Sizeup of an algorithm should be linear. However, practically Sizeup may be sub-linear. Figure 13a shows
- 4 Sizeup of DLSTSVM for different computation units (executor-cores) where multiple curves show the Sizeup

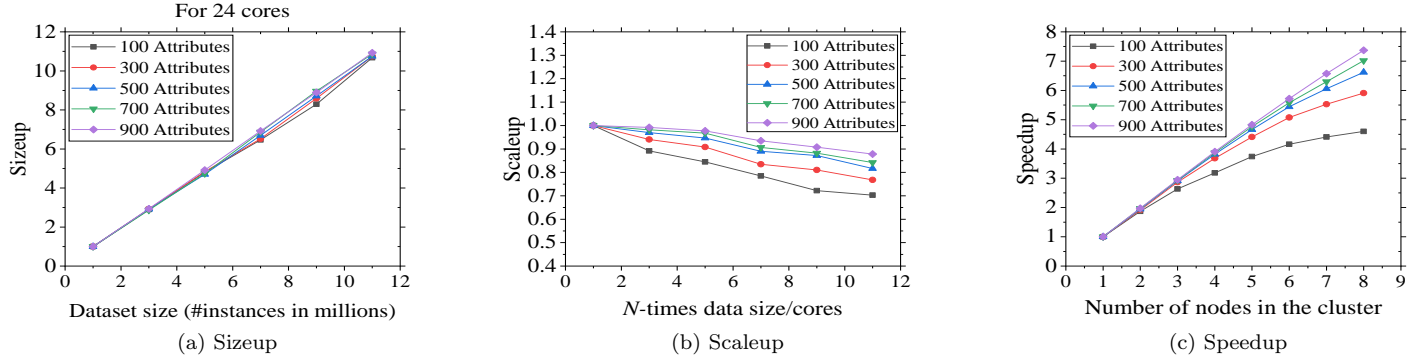


Figure 13: 3S based performance evaluation of DLSTSVM

- 1 for 100, 300, 500, 700 and 900 attributes respectively. Evidently, Sizeup of DLSTSVM is almost linear.
- 2 **Scaleup:-** This parameter captures the behavior of algorithm in case when data size and computation units
 3 are scaled in same ratio. Let a system with 1 computation unit takes time t on data size d . If k times larger
 4 computing system takes t time on kd size data, then system scales perfectly. However, perfect scalability is
 5 practically typical to achieve. Figure 13b specifies Scaleup of DLSTSVM for different settings of attributes. To
 6 measure Scaleup, number of executor cores in the cluster is set to 1, 3, 5, 7, 9 and 11 corresponding to 1 million,
 7 3 million, 5 million, 7 million, 9 million and 11 million data instances respectively. It is observed that Scaleup
 8 degrades with rise of executor cores (or #machines) due to inter-machine communication cost, task scheduling
 9 delays, uneven machine performance and skewed data distribution. In case of lesser attributes such as 100, all
 10 afore-mentioned overheads are dominant over computation cost, hence, degrade in Scaleup is maximum for it.
 11 However, in case of more attributes (300 and above) and with 9 to 11 million data instances, Scaleup varies
 12 from 88% to 70% whereas with 1 to 9 million data instances 100% to 88% which is quite good.
- 13 **Speedup:-** Speedup describes the behavior of a distributed parallel algorithm with increase in number of ma-
 14 chines in the cluster keeping the data size fixed. Presented work evaluates speedup over 1 to 8 nodes/machines
 15 (i.e. number of cores from 1 to 24) in the cluster. Figure 13c depicts speedup where different curves corresponds
 16 to number of attributes 100, 300, 500, 700 and 900 respectively. Here, number of instances is fixed at 11 million.
 17 DLSTSVM exhibits prominent speedups. For number of attributes ≥ 300 speed up is extremely good.
 18 Speedup decreases with adding more number of nodes in the cluster which is expected behavior of any scal-
 19 able distributed algorithm. For 100 attributes, the decrement in speedup is more because of dominance of
 20 communication and scheduling overheads over computation costs. All discussions regarding the performance of
 DLSTSVM for the 3S evaluators are summarized in Table 5.

Table 5: Summary of performance of DLSTSVM for 3S evaluators

Performance Evaluators	# data instances	Attribute Range		
		100	100-300	300
Sizeup	1-11 million	Almost Linear		
		Range of Scaleup achieved		
Scaleup	1-7 million	100%-79%	100%-89%	
	7-11 million	84%-70%	93%-82%	
		Speedup achieved		
Speedup	1-11 million	Average	Good	Very Good

6. Conclusion and Future Work

In this research work a novel technique DLSTSVM is proposed which achieves scalability by using MapReduce to process distributed data partitions on cluster of machines in parallel fashion. It not only averts memory and computational bottlenecks of existing non-parallel version LSTSVM, but also achieves significant speedups too without losing predictive accuracy. Detailed experimental comparison with existing non-parallel versions of SVM and LSTSVM shows that DLSTSVM exhibits speedups in running-time and achieve scalability with respect to large scale datasets. Number of attributes has more impact on run-time. DLSTSVM shows speedup 3x to 20x performance gains. Speedups Performance of DLSTSVM is evaluated along well known 3S criterion, i.e. Sizeup, Scaleup and Speedup, and found to be effective. Sizeup of DLSTSVM is almost linear. For attributes above 300 and data instances 7 to 11 million, Scaleup of 77% to 53% whereas for 1 to 7 million data instances, Scaleup of 100% to 70% is achieved effectively. DLSTSVM exhibits prominent speedups and is extremely good for number of attributes ≥ 300 .

References

- [1] Cortes C, Vapnik V. Support-vector networks. *Machine learning* 1995; 20(3):273-97.
- [2] Khemchandani R, Chandra S. Twin support vector machines for pattern classification. *IEEE Transactions on pattern analysis and machine intelligence* 2007; 29(5):905-10.
- [3] Kumar MA, Gopal M. Least squares twin support vector machines for pattern classification. *Expert systems with applications* 2009; 36(4):7535-43.
- [4] Dean J, Ghemawat S. MapReduce: simplified data processing on large clusters. *Communications of the ACM* 2008; 51(1):107-13.
- [5] Isard M, Budi M, Yu Y, Birrell A, Fetterly D. Dryad: distributed data-parallel programs from sequential building blocks. In *Proceedings of the 2nd ACM SIGOPS/EuroSys European Conference on Computer Systems* 2007, pp. 59-72.
- [6] Zaharia M, Chowdhury M, Franklin MJ, Shenker S, Stoica I. Spark: cluster computing with working sets In: *Proceedings of the 2nd USENIX Conference on Hot Topics in Cloud Computing*, 10–10. USENIX Association, Berkeley, CA, USA 2010.
- [7] Inoubli W, Aridhi S, Mezni H, Maddouri M, Nguifo EM. An experimental survey on big data frameworks. *Future Generation Computer Systems* 2018; 86 :546-64.
- [8] Bal H, Pal A. Parallel and Distributed Machine Learning Algorithms for Scalable Big Data Analytics. *Future Generation Computer Systems* 2020; 108 :1159-1161.
- [9] Zhang M L, Zhou Z H. A review on multi-label learning algorithms. *IEEE transactions on knowledge and data engineering* 2013; 26(8):1819-37.
- [10] Hsieh C J, Si S, Dhillon I. A divide-and-conquer solver for kernel support vector machines. In *International conference on machine learning* 2014, pp. 566-574.
- [11] Chang C C, Lin C J. LIBSVM: A library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)* 2011 May 6; 2(3):1-27.
- [12] Bordes A, Ertekin S, Weston J, Bottou L. Fast kernel classifiers with online and active learning. *Journal of Machine Learning Research* 2005; 6(Sep):1579-619.
- [13] Graf H, Cosatto E, Bottou L, Dourdanovic I, Vapnik V. Parallel support vector machines: The cascade svm. *Advances in neural information processing systems* 2004; 17:521-8.

- 1 [14] Zhang K, Lan L, Wang Z, Moerchen F. Scaling up kernel svm on limited resources: A low-rank linearization
2 approach. InArtificial intelligence and statistics 2012 Mar 21, pp. 1425-1434.
- 3 [15] Le Q, Sarlós T, Smola A. Fastfood-computing hilbert space expansions in loglinear time. InInternational Conference
4 on Machine Learning 2013 Feb 13, pp. 244-252.
- 5 [16] Selvaraj S K, Decoste D M, inventors; Altaba Inc, assignee. Building support vector machines with reduced classifier
6 complexity. United States patent US 7,630,945. 2009 Dec 8.
- 7 [17] Moody J, Darken CJ. Fast learning in networks of locally-tuned processing units. *Neural computation*. 1989 Jun; 1
8 (2):281-94.
- 9 [18] Sartakhti J S, Afrabandpey H, Ghadiri N. Fuzzy least squares twin support vector machines. *Engineering Applica-*
10 *tions of Artificial Intelligence*. 2019 Oct 1; 85:402-9.
- 11 [19] Chen S, Cao J, Chen F, Liu B. Entropy-Based Fuzzy Least Squares Twin Support Vector Machine for Pattern
12 Classification. *Neural Processing Letters*. 2020 Feb; 51(1):41-66.
- 13 [20] Richhariya B, Tanveer M. Universum least squares twin parametric-margin support vector machine. In2020 Inter-
14 national Joint Conference on Neural Networks (IJCNN) 2020 Jul 19, pp. 1-8.
- 15 [21] Mir A, Nasiri J A. KNN-based least squares twin support vector machine for pattern classification. *Applied*
16 *Intelligence*. 2018 Dec 1; 48(12):4551-64.
- 17 [22] Zhu K, Wang H, Bai H, Li J, Qiu Z, Cui H, Chang EY. Parallelizing support vector machines on distributed
18 computers. InAdvances in neural information processing systems 2008, pp. 257-264.
- 19 [23] Boyd S, Parikh N, Chu E. Distributed optimization and statistical learning via the alternating direction method of
20 multipliers. Now Publishers Inc; 2011.
- 21 [24] Guan L, Sun T, Qiao L B, Yang Z H, Li D S, Ge K S, Lu X C. Anefficient parallel and distributed solution to
22 nonconvex penalized linear SVMs. *Frontiers of Information Technology & Electronic Engineering*. 2019 Sep, 5:1-7.
- 23 [25] Singh D, Mohan C K. Projection-SVM: Distributed Kernel Support Vector Machine for Big Data using Subspace
24 Partitioning. In2018 IEEE International Conference on Big Data (Big Data) 2018 Dec 10, pp. 74-83.
- 25 [26] Shah R, Zhang S, Lin Y, Wu P. xSVM: Scalable Distributed Kernel Support Vector Machine Training. In2019 IEEE
26 International Conference on Big Data (Big Data) 2019 Dec 9, pp. 155-164.
- 27 [27] Musicant D R. NDC: normally distributed clustered datasets. Computer Sciences Department, University of Wis-
28 consin, Madison. 1998.
- 29 [28] Tavara S. Parallel computing of support vector machines: a survey. *ACM Computing Surveys (CSUR)* 2019;
30 51(6):1-38.
- 31 [29] Landset S, Khoshgoftaar TM, Richter AN, Hasanin T. A survey of open source tools for machine learning with big
32 data in the Hadoop ecosystem. *Journal of Big Data* 2015; 2(1):24.
- 33 [30] [University of California, Irvine, Datasets](#)
- 34 [31] Chang C C, Lin C J. LIBSVM: A library for support vector machines. *ACM transactions on intelligent systems*
35 *and technology (TIST)* 2011; 2(3):1-27.
- 36 [32] Zhang K, Lan L, Wang Z, Moerchen F. Scaling up kernel svm on limited resources: A low-rank linearization
37 approach. InArtificial intelligence and statistics 2012, pp. 1425-1434.
- 38 [33] Le V L, Lauer F, Bako L, Bloch G. Learning nonlinear hybrid systems: from sparse optimization to support vector
39 regression. InProceedings of the 16th international conference on Hybrid systems: computation and control 2013,
40 pp. 33-42.
- 41 [34] Hsieh C J, Si S, Dhillon I. A divide-and-conquer solver for kernel support vector machines. InInternational conference
42 on machine learning 2014, pp. 566-574.
- 43 [35] Çatak F Ö, Balaban M E. A MapReduce-based distributed SVM algorithm for binary classification. *Turkish Journal*
44 *of Electrical Engineering & Computer Sciences* 2016; 24(3):863-73.