

Distributed Sparse Support Vector Machine for Feature Selection on High Dimensional Datasets

Xiaowei KUANG

Department of Computer Science and Engineering

The Chinese University of Hong Kong

Abstract

Feature selection is very useful in reducing model complexity and gaining insights into data. As practical data often come in high volume and high dimension which exceed the computing capabilities and storage limit of a single machine, distributed implementations of feature selection algorithms are highly desirable. Incorporating an up-to-date feature selection algorithm called Feature Generating Machine (FGM) and an efficient distributed system named Husky, we propose a novel distributed implementation framework for feature selection in this report. The input data are partitioned among different workers and the workers can work jointly to solve model solutions with limited network communication. By cleverly caching important data, our method can handle large-scale high dimensional problems.

1 Introduction

It has been estimated that as of 2012, about 2.5 exabytes (10^{18}) of data are created each day, and that number is doubling every 40 months or so [1]. This exponential growth of data in terms of volume and dimension since the 21th century has spurred the development of feature selection methods and distributed computing systems. In the big data era, linear SVM has been one of the most popular tools for large-scale classification. Extensive efforts have also been put in applying SVM for feature selection. While many of the SVM algorithms proposed have shown state-of-the-art performance in certain applications, most of them are restricted to single machine training. This motivates us to develop a distributed solution that is more scalable and cost-efficient.

In a Machine Learning setting, feature selection refers to the process of selecting a subset of relevant features from the input data. Feature selection is useful in a number of ways [2]: First, it helps us visualize and understand the data better. Second, it alleviates the curse of dimensionality problem present in many real-world problems, which leads to a more accurate predictor. Third, as a dimension reduction technique, it significantly reduces the size of the input data by eliminating irrelevant features, which gives rise to faster and more cost-effective predictors in terms of time complexity and space requirement. To name one concrete example, consider the text classification problem: an input instance (document) lies in a feature space of dimension the size of the vocabulary containing word frequency count, which can count to millions! Recently, many features selection methods using SVM have been proposed. Many of them achieve the goal of feature selection through sparsity regularization. Instead of using the l_2 -norm regularizer($\|\mathbf{w}\|_2^2$), which results in non-sparse solutions, regularizers such as l_0 -norm regularizer($\|\mathbf{w}\|_0$) [3] and l_1 -norm regularizer($\|\mathbf{w}\|_1$) [4] are employed. While the above algorithms are good for enforcing sparsity in the solutions, they are computationally more expensive than l_2 regularization even after convex relaxation. Besides sparse regularization, an effective wrapper feature selection method Recursive Feature Elimination (SVM-RFE) was proposed [5], which recursively trains a classifier, ranks all features according to some criteria and eliminates features that have lowest ranking. While this method has achieved state of the art performance on gene selection, it is computationally expensive and may even be infeasible for high dimensional problem. To address such a problem, a Feature Generating Machine(FGM) method was proposed [6], which iteratively generates a subset of most violated features and solve the resulting Multiple Kernel Learning (MKL) problem. We will base our distributed support vector machine algorithms on the Feature Generating Machine Algorithm in this report.

To deal with the challenges presented by the huge volume of data, many general purpose distributed computing systems such as Spark [7], Dryad [8] and Hadoop [9] were developed. These systems allow programmers to express the application logic in simple coarse-grained primitives like *map* and *reduce*, which simplifies programming for distributed algorithms. But these coarse-grained programming paradigms often do not result in efficient program. For example, a MapReduce program often has to dump the data to the distributed file system, load them into the main memory, parse the data, do the computation, dump the data back to the distributed file system and repeat the process again. As this process involves a lot of unnecessary disk IO, network traffic and serialization/deserialization, the program becomes very inefficient. To address these inefficiencies, programmers often need to resort to more sophisticated domain-specific languages

(DSLs), which raises development cost. This problem motivates the development of Husky, an efficient and expressive distributed computing system designed to strike a better balance between the efficiency and development cost problems mentioned above. There are several features that make Husky attractive as a distributed computing framework. First, Husky is good for interactive data analysis. User can easily perform exploratory data analysis on Husky by using scripting language. Second, Husky is designed to cooperate with the Hadoop ecosystem. This makes Husky attractive because HDFS and Hive are now the dominating big data systems in industry. Third, Husky provides an expressive and user-friendly Application Programming Interface (API). User only needs to understand a few core concepts such as *Object List*, *Channel*, *Aggregator* and *list_execute* to write sophisticated and powerful distributed program. Fourth, Husky is fault-tolerant and easily scalable. Machine failure and growth of computing cluster are handled seamlessly by the runtime "Master" in Husky. Last but not least, Husky can easily outperform computing system that offers coarse-grained primitives and is able to achieve similar or even better performance compared with domain-specific systems[10].

In this report, we propose a distributed version of the FGM algorithm described in [11]. It iteratively generates a pool of violated sparse feature subsets and solves the resulting MKL problem. Adapting the algorithm in [12] to our problem, we are able to divide the original MKL problem into several MKL problems of smaller instances and solve them distributively using a simple and efficient coordinate descent method [13]. The rest of this report is organized as follows. Section 2 gives an introduction to SVM, Sparse SVM and Husky. Section 3 describe our algorithm Distributed Sparse SVM. The last section gives a summary and remark of our algorithm.

2 Background

2.1 Support Vector Machine

Support Vector Machine (SVM) [14] is one of the most often used supervised learning algorithms for pattern recognition. A SVM for binary classification aims to learn a separating hyperplane that maximizes the margin between two classes of data. In a typical machine learning setting, given a labeled training set $\{\mathbf{x}_i, y_i\}_{i=1}^n$, where $\mathbf{x}_i \in \mathbf{R}^m$ is the m-dimensional input and $y_i \in \{+1, -1\}_{i=1}^n$ is the output label, we learn a m-dimensional decision hyperplane $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ that maximizes the margin $\frac{1}{\|\mathbf{w}\|_2}$, or equivalently, minimizes the l_2 norm $\|\mathbf{w}\|_2^2$ of the normal of the hyperplane

such that all positive data lie on one side of the plane and all the negative data lie on the other side of the hyperplane (as shown in Figure 1). Then the parameter \mathbf{w} can be found by solving the following optimization problem.

$$\min_{\mathbf{w}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 \quad \text{s.t.} \quad y_i \mathbf{w}^T \mathbf{x}_i \geq 1, \quad i = 1, \dots, n \quad (1)$$

The inequality specified above constrains the hyperplane to pass through the origin. In order to solve this problem, we can add a bias term b in the hyperplane by changing the input instance \mathbf{x} and parameter \mathbf{w} as the following [13]:

$$\mathbf{x}_i^T \leftarrow [\mathbf{x}_i^T, 1] \quad \mathbf{w}_i^T \leftarrow [\mathbf{w}_i^T, b] \quad (2)$$

However, if the input data are not linearly separable, the constraints in equation (1) can not all be satisfied. To tackle this hard margin problem, slack variables are introduced to loosen the inequality constraints and we can reformulate the problem as a soft margin maximization problem:

$$\min_{\mathbf{w}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 \quad \text{s.t.} \quad y_i \mathbf{w}^T \mathbf{x}_i \geq 1 - \xi_i, \quad i = 1, \dots, n \quad (3)$$

With ξ_i , we allow an input instance \mathbf{x}_i to violate the original hard margin constraint but penalize such violations by adding a positive term $\frac{C}{2} \xi_i^2$ to the objective function. The parameter C is called the slack trade off. A smaller C allows more points to be misclassified and a greater C forces the margin to become "harder". The loss function we use here is called square hinge loss, taking the form $loss(\mathbf{w}, \mathbf{x}_i, y_i) = \max(1 - y_i \mathbf{w}^T \mathbf{x}_i, 0)^2$. Another often used loss function in training a SVM is hinge loss, which is in the form $loss(\mathbf{w}, \mathbf{x}_i, y_i) = \max(1 - y_i \mathbf{w}^T \mathbf{x}_i, 0)$. These two functions are used because they are both convex functions, which is a desired property in an optimization problem. In addition to being a linear classifier, SVM can also be used to model nonlinear relationship by using the kernel trick. The idea stems from the fact that when formulated in its dual form (see the next subsection), the objective function can be written as a linear combination of the dot product $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$, which can be viewed as the similarity between instance \mathbf{x}_i and instance \mathbf{x}_j [15]. Therefore, we can replace $\mathbf{x}_i^T \mathbf{x}_j$ with a Kernel function $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$, where $\phi: \mathcal{X} \rightarrow \mathcal{Z}$ is a mapping from feature space X onto another feature space Z . Then a linear SVM corresponds to the Kernel $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$. One commonly used Kernel is the Gaussian Kernel $K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)$. Since the Taylor expansion of a function allows us to indefinitely

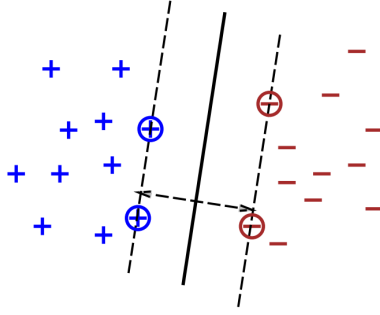


Figure 1: Figure 1. The solid line in the middle is the hyperplane that separates the positive(+) data points from the negative(-) data points. The circled points that lie on the dashed line are called support vectors.

approximate a continuous function by adding high order polynomial term, with each polynomial representing a new dimension, the Gaussian Kernel allows us to map our our data from the original n-dimensional feature space onto a feature space of infinite dimension without explicitly calculating the dot product in the infinite-dimensional space. One important consequence of this mapping is that if the input data are not linearly separable in its original feature space, we can always map it onto a feature space of higher dimension to make it separable by a hyperplane in the higher dimensional feature space. But consideration must be given to input data with noisy features, in which case such a method may not generalize well.

2.2 Sparse SVM

2.2.1 Primal and Dual formulation of Sparse SVM

Following the notation in [11], we denote $A \odot B$ to be the elementwise product between two vectors/matrices A and B. To obtain a sparse parameter \mathbf{w} of SVM, a 0-1 control vector $\mathbf{d} = [d_1, \dots, d_m]^T \in \mathcal{D}$ is introduced to control which features get selected. The decision hyperplane thus becomes: $f(x) = \mathbf{w}^T x = (\tilde{\mathbf{w}} \odot \mathbf{d})^T x = \tilde{\mathbf{w}}^T (\mathbf{x} \odot \mathbf{d})$, where $\mathcal{D} = \{\mathbf{d} | \sum_{j=1}^m d_j \leq B, d_j \in \{0, 1\}, j = 1, \dots, m\}$ is the domain of \mathbf{d} , and B controls the sparsity of \mathbf{d} . Using square hinge loss, the objective function of Sparse SVM can be formulated as:

$$\min_{\mathbf{d} \in \mathcal{D}} \min_{\tilde{\mathbf{w}}, \xi} \frac{1}{2} \|\tilde{\mathbf{w}}\|_2^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 \quad s.t. \quad y_i \tilde{\mathbf{w}}^T (\mathbf{x}_i \odot \mathbf{d}) \geq 1 - \xi_i, i = 1, \dots, n \quad (4)$$

To derive its dual form, we introduce one Lagrange multiplier for each constraint. Let $\mathcal{L}(\tilde{\mathbf{w}}, \xi, \alpha) = \frac{1}{2} \|\tilde{\mathbf{w}}\|_2^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 - \sum_{i=1}^n \alpha_i (y_i \tilde{\mathbf{w}}^T (\mathbf{x}_i \odot \mathbf{d}) - 1 + \xi_i)$, where $\alpha = [\alpha_1, \dots, \alpha_n]^T$ is a vector of dual

variables for the inequality constraints and $\boldsymbol{\alpha} \in \mathcal{A}$, $\mathcal{A} = \{\boldsymbol{\alpha} | \alpha_i \geq 0, i = 1, \dots, n\}$, it can be seen that,

$$\max_{\boldsymbol{\alpha} \in \mathcal{A}} \mathcal{L}(\tilde{\boldsymbol{w}}, \boldsymbol{\xi}, \boldsymbol{\alpha}) = \begin{cases} \infty & \text{if some constraints are violated} \\ \frac{1}{2} \|\tilde{\boldsymbol{w}}\|_2^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 & \text{otherwise} \end{cases}$$

If a parameter candidate $\tilde{\boldsymbol{w}}$ violates some constraints, then $\mathcal{L}(\tilde{\boldsymbol{w}}, \boldsymbol{\xi}, \boldsymbol{\alpha}) = \infty$ and such a solution will not be chosen when we want to minimize over $\tilde{\boldsymbol{w}}$. Hence (4) is equivalent to

$$\min_{\boldsymbol{d} \in \mathcal{D}} \min_{\tilde{\boldsymbol{w}}, \boldsymbol{\xi}} \max_{\boldsymbol{\alpha} \in \mathcal{A}} \frac{1}{2} \|\tilde{\boldsymbol{w}}\|_2^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 - \sum_{i=1}^n \alpha_i (y_i \tilde{\boldsymbol{w}}^T (\boldsymbol{x}_i \odot \boldsymbol{d}) - 1 + \xi_i)$$

Since the maximum of the minimum of a function is always less than or equal to the minimum of the maximum of the same function, we have

$$\max_{\boldsymbol{\alpha} \in \mathcal{A}} \min_{\tilde{\boldsymbol{w}}, \boldsymbol{\xi}} \frac{1}{2} \|\tilde{\boldsymbol{w}}\|_2^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 - \sum_{i=1}^n \alpha_i (y_i \tilde{\boldsymbol{w}}^T (\boldsymbol{x}_i \odot \boldsymbol{d}) - 1 + \xi_i) \leq \min_{\tilde{\boldsymbol{w}}, \boldsymbol{\xi}} \max_{\boldsymbol{\alpha} \in \mathcal{A}} \frac{1}{2} \|\tilde{\boldsymbol{w}}\|_2^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 - \sum_{i=1}^n \alpha_i (y_i \tilde{\boldsymbol{w}}^T (\boldsymbol{x}_i \odot \boldsymbol{d}) - 1 + \xi_i)$$

According to Lagrange Duality theory, the above inequality satisfies with equality and we can solve (4) by solving the following:

$$\min_{\boldsymbol{d} \in \mathcal{D}} \max_{\boldsymbol{\alpha} \in \mathcal{A}} \min_{\tilde{\boldsymbol{w}}, \boldsymbol{\xi}} \frac{1}{2} \|\tilde{\boldsymbol{w}}\|_2^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 - \sum_{i=1}^n \alpha_i (y_i \tilde{\boldsymbol{w}}^T (\boldsymbol{x}_i \odot \boldsymbol{d}) - 1 + \xi_i) \quad (5)$$

Taking the above objective function with respect to $\tilde{\boldsymbol{w}}$ and $\boldsymbol{\xi}$, setting them to 0, we have

$$\tilde{\boldsymbol{w}} = \sum_{i=1}^n \alpha_i y_i (\boldsymbol{x}_i \odot \boldsymbol{d}) \quad \alpha_i = C \xi_i, i = 1, \dots, n$$

Substitute the above back to (5), we get the dual form of the SSVM problem as follows. Notice how the objective function is expressed as a linear combination of the dot product $\langle \boldsymbol{x}_i, \boldsymbol{x}_j \rangle$ as mentioned in section 2.1.

$$\min_{\boldsymbol{d} \in \mathcal{D}} \max_{\boldsymbol{\alpha} \in \mathcal{A}} -\frac{1}{2} \left\| \sum_{i=1}^n \alpha_i y_i (\boldsymbol{x}_i \odot \boldsymbol{d}) \right\|^2 - \frac{1}{2C} \boldsymbol{\alpha}^T \boldsymbol{\alpha} + \boldsymbol{e}^T \boldsymbol{\alpha} \quad (6)$$

According to the convex relaxation in [11], (6) can be written as a MKL problem

$$\max_{\boldsymbol{\mu} \in \mathcal{M}} \min_{\boldsymbol{\alpha} \in \mathcal{A}} f^{\mathcal{D}}(\boldsymbol{\alpha}, \boldsymbol{\mu}) = \frac{1}{2} \boldsymbol{\alpha}^T Y \left(\sum_{\boldsymbol{d}^t \in \mathcal{D}} \mu_t X_t X_t^T + \frac{1}{C} I \right) Y \boldsymbol{\alpha} - \boldsymbol{e}^T \boldsymbol{\alpha} \quad (7)$$

Or equivalently,

$$\max_{\boldsymbol{\mu} \in \mathcal{M}} \min_{\boldsymbol{\alpha} \in \mathcal{A}} f^{\mathcal{D}}(\boldsymbol{\alpha}, \boldsymbol{\mu}) = \frac{1}{2} \boldsymbol{\alpha}^T \bar{Q} \boldsymbol{\alpha} - \mathbf{e}^T \boldsymbol{\alpha}$$

where $\mathcal{M} = \{\boldsymbol{\mu} | \mathbf{e}^T \boldsymbol{\mu} = 1, \boldsymbol{\mu} \geq 0\}$, $\bar{Q} = Y(\sum_{\mathbf{d}^t \in \mathcal{D}} \mu_t X_t X_t^T + \frac{1}{C} I) Y$, $X_t = [\mathbf{x}_1 \odot \mathbf{d}^t, \dots, \mathbf{x}_n \odot \mathbf{d}^t]^T$ and Y is a diagonal matrix such that $Y_{ii} = y_i$.

2.2.2 Cutting Plane Algorithm

The overall algorithm for solving (7) is described in Algorithm 1 [11]. We denote the subset of constraints by $\mathcal{C} \subset \mathcal{D}$. First the vector of dual variables $\boldsymbol{\alpha}$ is initialized to $\frac{1}{n} \mathbf{1}$. We then find the most violated $\hat{\mathbf{d}} \in \mathcal{D}$ and initialize the working set $\mathcal{C} = \{\hat{\mathbf{d}}\}$. To solve the MKL problem, we use an efficient method called SimpleMKL[16]. We fix $\boldsymbol{\mu}$ and run SVM solver to update $\boldsymbol{\alpha}$ one step. Then we fix $\boldsymbol{\alpha}$ and update $\boldsymbol{\mu}$ one step. This iterative updating procedure of $\boldsymbol{\alpha}$ and $\boldsymbol{\mu}$ is repeated until the MKL problem converges. Then the next most violated constraint \mathbf{d} is calculated and added to \mathcal{C} . This process is repeated until the termination criterion is met. To find the most violated constraint, the constraint that has the highest feature score given the current values of $\boldsymbol{\alpha}$, we denote $c_j = \sum_{i=1}^n \alpha_i y_i x_{ij}$ and sort the c_j s in descending order and set the first B numbers corresponding to \mathbf{d}_j to 1 and the rests to 0.

2.2.3 Prediction

When the algorithm converges, the decision function can be obtained by $f(x) = \sum_{\mathbf{d}^t \in \mathcal{D}} \sum_{i=1}^n \alpha_i y_i (\mathbf{x}_i \odot \mathbf{d}^t)^T \mathbf{x} = \sum_{i=1}^n \alpha_i y_i (\mathbf{x}_i \odot \tilde{\mathbf{d}})^T \mathbf{x}$, where $\tilde{\mathbf{d}} = \sum_{\mathbf{d}^t \in \mathcal{D}} \mu_t \mathbf{d}^t$

Algorithm 1 The cutting plane algorithm for FGM

- 1: Initialize $\boldsymbol{\alpha} = \frac{1}{n} \mathbf{1}$. Find the most violated $\hat{\mathbf{d}}$, and set $\mathcal{C} = \{\hat{\mathbf{d}}\}$
 - 2: Run MKL solver to solve for $\boldsymbol{\alpha}$ and $\boldsymbol{\mu}$
 - 3: Find the most violated $\hat{\mathbf{d}}$ and set $\mathcal{C} = \mathcal{C} \cup \hat{\mathbf{d}}$
 - 4: Repeat step 2-4 until convergence
-

2.3 Simple MKL

To solve Multiple Kernel Learning problem like problem 7, an efficient algorithm called SimpleMKL [16] has been proposed. Formally, the Multiple Kernel Learning problem corresponds to the

following primal optimization problem:

$$\begin{aligned}
& \underset{f_m, b, \xi, \mathbf{d}}{\text{minimize}} && \frac{1}{2} \sum_m \frac{1}{d_m} \|f_m\|_{\mathcal{H}_m}^2 + C \sum_{i=1}^n \xi_i \\
& \text{subject to} && y_i \left(\sum_m f_m(x_i) + b \right) \geq 1 - \xi_i, \quad \forall i \\
& && \xi_i \geq 0 \quad \forall i \\
& && \sum_m d_m = 1, \quad d_m \geq 0, \quad \forall m
\end{aligned}$$

where \mathcal{H}_m defines a space that is associated with a particular kernel K_m and d_m can be understood to be the weight of a particular kernel in the decision function. To solve the above, a constrained optimization problem is proposed:

$$\underset{\mathbf{d}}{\text{minimize}} \quad J(\mathbf{d}) \quad \text{subject to} \quad \sum_{m=1}^M d_m = 1, d_m \geq 0 \quad m = 1, \dots, M \quad (8)$$

where $J(\mathbf{d})$ is the optimal value of the following SVM problem

$$\begin{aligned}
& \underset{f_m, b, \xi}{\text{minimize}} && \frac{1}{2} \sum_m \frac{1}{d_m} \|f_m\|_{\mathcal{H}_m}^2 + C \sum_{i=1}^n \xi_i \\
& \text{subject to} && y_i \left(\sum_m f_m(\mathbf{x}_i) + b \right) \geq 1 - \xi_i, \quad \forall i \\
& && \xi_i \geq 0 \quad \forall i
\end{aligned} \quad (9)$$

Thus the above can be solved using the following iterative procedure. First a SVM solver is invoked to solve the SVM problem (9), then reduced gradient method is used to find a feasible descent direction for \mathbf{d} , followed by a line search for the optimal step size. The above iterative process is repeated until the duality gap of the Multiple Kernel Learning problem is within a tolerable threshold. Note in order to ensure the equality constraint $\sum_{m=1}^M d_m = 1$ and the positivity constraint $d_m \geq 0$, reduced gradient method is used instead of a normal gradient method. In particular, the reduced gradient of $J(\mathbf{d})$ denoted by $\nabla_{red} J$ has the following components:

$$\begin{aligned}
[\nabla_{red} J]_m &= \left(\frac{\partial J}{\partial d_m} - \frac{\partial J}{\partial d_u} \right) \quad \forall m \neq u \\
[\nabla_{red} J]_u &= \sum_{m \neq u} \left(-\frac{\partial J}{\partial d_m} + \frac{\partial J}{\partial d_u} \right)
\end{aligned}$$

where u can be any one of $1, \dots, M$. Since $\gamma \sum_{m \neq u} \frac{\partial J}{\partial d_m} - \frac{\partial J}{\partial d_u} + \gamma \sum_{m \neq u} -\frac{\partial J}{\partial d_m} + \frac{\partial J}{\partial d_u} = 0$ for all possible step size γ , the equality constraint is satisfied. To ensure that the positivity constraint is

preserved, a feasible descent direction D can be found by projecting the reduced gradient to the feasible domain:

$$D_m = \begin{cases} 0, & \text{if } d_m = 0 \text{ and } \frac{\partial J}{\partial d_m} - \frac{\partial J}{\partial d_u} \geq 0 \\ -\frac{\partial J}{\partial d_m} + \frac{\partial J}{\partial d_u}, & \text{if } d_m > 0 \text{ and } m \neq u \\ \sum_{v \neq u, d_v \geq 0} (\frac{\partial J}{\partial d_v} - \frac{\partial J}{\partial d_u}) & \text{if } m = u \end{cases} \quad (10)$$

The above simply says that if an entry d_m is equal to 0 and its reduced gradient is non-negative (negative of reduced gradient is positive), then we need to project its descent direction to be 0 indicating that it is already at its optimum. Otherwise we set its descent direction to be its reduced gradient. With the descent direction available, one can perform a simple step size search to find the optimal step size and update \mathbf{d} . A detailed algorithm adapted to the setting of this report will be shown in section 3.

2.4 Core Concepts in Husky

Husky follows a master-slave architecture. One Husky cluster is made up of a master, which is responsible for load balancing and coordination between workers, and multiple workers which can reside in one or multiple machines. Below is a description of the four most important concepts in Husky.

Object List is the primary data abstraction in Husky, where user-defined objects can be stored. In a distributed setting, *Object List* is a distributed abstraction of objects. When writing the program, users can treat the *Object List* as a single array of object when in fact it is distributed among different workers, with each worker storing a portion of the list. This is a very useful abstraction as users can write a distributed program as if he/she is writing a single machine program.

Channel is another important concept in Husky. It defines how *Object Lists* communicate with each other. Messages, which can be any user-defined data types, are sent through channels to a distant location using the *push* function and are received through channels by a distant target using the *get* function.

Aggregator is a powerful *allreduce* like object designed for easily aggregating values from executors and broadcasting the aggregated values back to the executors.

`list_execute` is the most important function provided by Husky. Synchronization of the above three concepts across different workers are done through `list_execute`. Consider the following C++ example

```
for (int i = 0; i < objList.size(); i++)
    std::cout << objList[i].id() << std::endl;
```

The piece of code above can be easily expressed in a distributed setting in Husky as follows.

```
// objects in objList reside in different workers
list_execute(objList, [] (Obj& obj) {
    log_msg(obj.id());
});
```

As can be seen, even though the objects are distributed among different workers, we can use `list_execute` to easily access the objects as if we are writing a simple *for-loop*. And most importantly, we have fine grained control at the level of objects instead of a coarse grained control at the level of tables like that provided by Spark[17] and Hive[18].

3 Distributed Feature Generating Machine

In this section, we present a distributed implementation framework for the Sparse SVM algorithm. Building on the communication pattern of Husky, we carefully investigate the mathematical relationship between different parameters and design a scheme that can coordinate different workers to jointly solve the Sparse SVM problem with low communication cost. We start by distributing the input data $\{\mathbf{x}_i, y_i\}_{i=1}^n$ among K workers. To make it clear, we introduce the following notation: Instances in machine k are denoted as $\{(\mathbf{x}_i, y_i)\}_{i \in J_k}$ where J_k are disjoint sets such that $\cup_{k=1}^K J_k = \{1, \dots, n\}$. $\mathcal{P}(\boldsymbol{\alpha}) \equiv [-\alpha_1, \infty] \times \dots \times [-\alpha_n, \infty]$ is the feasible domain of $\boldsymbol{\alpha}$. $\|\mathbf{u}\|_A^2 \equiv \mathbf{u}^T A \mathbf{u}$ for $\mathbf{u} \in R^t$ and $A \in R^{t \times t}$. \mathbf{v}_{J_k} denotes the sub-vector of \mathbf{v} that contains the coordinates in J_k and similarly for A_{J_k, J_m} .

3.1 Finding the Most Violated Constraints

Finding the most violated constraints involves finding the constraint $\hat{\mathbf{d}}$ that yields the largest objective value in (6), which can be expressed as a function of $\hat{\mathbf{d}}$: $\frac{1}{2} \sum_{j=1}^m (\sum_{i=1}^n \alpha_i y_i x_{ij})^2 d_j =$

$\frac{1}{2} \sum_{j=1}^m c_j^2 d_j$, where $c_j = \sum_{i=1}^n \alpha_i y_i x_{ij}$. Hence the most violated constraint corresponds to the set of d_j s that has the B -largest c_j s. We can calculate c_j as

$$c_j = \bigoplus_{k=1}^K X_{J_k, j}^T Y_{J_k} \alpha_{J_k} \quad (11)$$

where $X_{J_k, j}$ stands for the vector obtained by taking the (i, j) entry of the input matrix X with $i \in J_k$. The symbol \bigoplus represent the operation of receiving the information from all workers and broadcasting the results back to all workers. In Husky, this can be done using *Aggregator*, which implements the *allreduce* functionality of broadcasting the values to all the workers and collecting the aggregated values from all the workers. Hence, each worker will update c_j using the input data available in its local storage and receive a global version of c_j . With the availability of c_j , each worker will sort the c_j in descending order and set the first B numbers corresponding to d_j to 1 and the rests to 0. Then the new generated $\hat{\mathbf{d}}$ is added to the constraint set \mathcal{C} . Notice because each machine has the same copy of c_j , the resulting $\hat{\mathbf{d}}$ is also the same for all workers, which obviates the need for a centralized master to send the new $\hat{\mathbf{d}}$ to every other workers.

3.2 Solving the MKL Problem

As mentioned in section 2, we use SimpleMKL to solve the MKL problem (7). A detailed description of the algorithm is given in Algorithm 2. where $K_t(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i \odot \mathbf{d}_t)^T (\mathbf{x}_j \odot \mathbf{d}_t)$ is the kernel

Algorithm 2 The SimpleMKL algorithm

- 1: Set $\mu_t = \frac{1}{T}$ for $T = 1, \dots, T$
 - 2: while stopping criterion not met do
 - 3: Compute $J(\boldsymbol{\mu})$ using a SVM solver with kernel $K = \sum_t \mu_t K_t$
 - 4: Compute $\frac{\partial J}{\partial \mu_t}$ for $\mu = 1, \dots, T$ and descent direction D
 - 5: Set $u = \operatorname{argmax}_t \mu_t$, $J' = 0$, $\boldsymbol{\mu}' = \boldsymbol{\mu}$, $\mathbf{D}' = \mathbf{D}$ while $J' < J(\boldsymbol{\mu})$ do
 - 6: $\boldsymbol{\mu} = \boldsymbol{\mu}'$, $\mathbf{D} = \mathbf{D}'$
 - 7: $v = \operatorname{argmin}_{\{t | D_t < 0\}} -\frac{\mu_t}{D_t}$, $\gamma_{max} = -\frac{\mu_v}{D_v}$
 - 8: $\boldsymbol{\mu}' = \boldsymbol{\mu} + \gamma_{max} \mathbf{D}$, $D'_u = D_u + D_v$, $D'_v = 0$ // normalization
 - 9: Compute J' using a SVM with kernel $K = \sum_t \mu_t K_t$
 - 10: end while
 - 11: line search along \mathbf{D} for $\gamma \in [0, \gamma_{max}]$
 - 12: $\boldsymbol{\mu} \leftarrow \boldsymbol{\mu} + \gamma \mathbf{D}$
 - 13: end while
-

corresponding to the t -th control variable. As described in the SimpleMKL algorithm, we need a SVM solver to solve for $\boldsymbol{\alpha}$. In the following sections, we describe a distributed SVM solver called Box constrained Quadratic Optimization algorithm[12], a single machine SVM solver called Dual

Coordinate Descent algorithm[13] and adapt these two algorithms to our setting.

3.2.1 Updating α

The main computations in solving (7) are the computation of the Kernel Matrix $\sum_{t=1}^T \mu_t X_t X_t^T$, which is easily computed in a single machine setting since all input instances are easily available. However, in a distributed setting, computing the exact Kernel matrix is not easy. One naive approach is that each worker can have in their local storage all the input instances $\{\mathbf{x}_i, y_i\}_{i=1}^n$. Another approach is that all machines receive the necessary data from other machines by means of communication. It is easily seen that the first method will hit a bottle neck when the volume of the input data gets larger and larger and the second method will incur too much communication cost. Therefore, we consider an approximation of the Kernel matrix $\sum_{t=1}^T \mu_t X_t X_t^T$ in each local machine using only the input data available locally.

Computing the Update direction: The algorithm starts with an initial point α^0 for (7) and iteratively generates a sequence of solutions $\{\alpha^l\}_{l=0}^\infty$ until convergence. At the l -th iteration, we update the current α^l by

$$\alpha^{l+1} = \alpha^l + \eta^l \Delta \alpha^l \quad (12)$$

where $\eta^l \in R$ is the step size and $\Delta \alpha^l \in R^n$ is the update direction. To compute $\Delta \alpha^l$, let $\mathbf{d} \in \mathcal{P}(\alpha^l)$ be any direction in the n -dimensional space (\mathbf{d} not to be confused with the control variable), according to Taylor series expansion, we have with μ fixed,

$$f^{\mathcal{D}}(\alpha^l + \mathbf{d}, \mu) = f^{\mathcal{D}}(\alpha^l, \mu) + \nabla f^{\mathcal{D}}(\alpha^l, \mu)^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T H \mathbf{d} \quad (13)$$

where $H = \tilde{Q} + \frac{1}{C} I$ and

$$\tilde{Q}_{i,j} = \begin{cases} Q_{i,j} & \text{if } i \in J_k \text{ and } j \in J_k \text{ for some } k \\ 0 & \text{otherwise} \end{cases}$$

where $Q = Y(\sum_{d^t \in \mathcal{D}} \mu_t X_t X_t^T) Y$ and H is a symmetric, block-diagonal matrix with K blocks, where the k -th block is

$$H_{J_k} = Y_{J_k} \left(\sum_{t=1}^T \mu_t X_{t,J_k} X_{t,J_k}^T + \frac{1}{C} I \right) Y_{J_k}$$

Then in order to minimize $f^{\mathcal{D}}(\alpha^l, \mu)$, we find an update direction \mathbf{d} such that (13) is minimized,

which results in the following problem

$$\Delta \boldsymbol{\alpha}^l = \underset{\mathbf{d} \in \mathcal{P}(\boldsymbol{\alpha}^l)}{\operatorname{argmin}} \nabla f^{\mathcal{D}}(\boldsymbol{\alpha}^l, \boldsymbol{\mu})^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T H \mathbf{d} \quad (14)$$

Since each machine have some input instances available in their local storage, they each can solve a part of $\Delta \boldsymbol{\alpha}^l$ locally

$$\Delta \boldsymbol{\alpha}_{J_k}^l = \underset{\mathbf{d}_{J_k} \in \mathcal{P}(\boldsymbol{\alpha}^l)_{J_k}}{\operatorname{argmin}} \nabla f^{\mathcal{D}}(\boldsymbol{\alpha}^l, \boldsymbol{\mu})_{J_k}^T \mathbf{d}_{J_k} + \frac{1}{2} \|\mathbf{d}_{J_k}\|_{H_{J_k}}^2 \quad (15)$$

We can calculate $\nabla f^{\mathcal{D}}(\boldsymbol{\alpha}^l, \boldsymbol{\mu})_{J_k}$ as follows

$$\nabla f^{\mathcal{D}}(\boldsymbol{\alpha}^l, \boldsymbol{\mu})_{J_k} = Y_{J_k} X_{J_k} \mathbf{w}^l + \frac{1}{C} \boldsymbol{\alpha}_{J_k}^l - \mathbf{e}$$

where \mathbf{w}^l can be obtained by gathering information from all workers.

$$\mathbf{w}^l = \bigoplus_{k=1}^K \sum_{\mathbf{d}^l \in \mathcal{D}} \mu_t X_{t, J_k}^T Y_{J_k} \boldsymbol{\alpha}_{J_k}^l \quad (16)$$

or alternatively,

$$\begin{aligned} \nabla f^{\mathcal{D}}(\boldsymbol{\alpha}^l, \boldsymbol{\mu})_{J_k} &= \sum_{t=1}^T \mu_t Y_{J_k} X_{t, J_k} \mathbf{w}_t + \frac{1}{C} \boldsymbol{\alpha}_{J_k}^l - \mathbf{e} \\ \mathbf{w}^l &= \sum_{t=1}^T \mu_t \mathbf{w}_t^l \end{aligned} \quad (17)$$

where $\mathbf{w}_t = \sum_{k=1}^K X_{t, J_k}^T Y_{J_k} \boldsymbol{\alpha}_{J_k}^l$ is the weight vector associated with the t-th kernel. In this report, we will take the second approach of computing \mathbf{w}^l . In fact this is a key step of reducing the communication cost and we will discuss the rationale for choosing this method later in the end of section 3.

As with (11), \mathbf{w}^l can also be maintained in local storage of each worker using *Aggregator*. Notice with $\nabla f^{\mathcal{D}}(\boldsymbol{\alpha}^l, \boldsymbol{\mu})_{J_k}$ and H_{J_k} available, (11) is in the same form as a standard SVM dual problem, with $\boldsymbol{\alpha}$, \bar{Q} , $-\mathbf{e}$ and $\mathcal{P}(0)$ replaced by \mathbf{d}_{J_k} , H_{J_k} , $\nabla f^{\mathcal{D}}(\boldsymbol{\alpha}^l, \boldsymbol{\mu})_{J_k}$ and $\mathcal{P}(\boldsymbol{\alpha}^l)_{J_k}$ respectively. We describe the dual coordinate descent method for solving (15) in the below.

Algorithm for solving (15) locally: For SVM problem in which the bias term b is not appended at the end of the decision hyperplane parameter \mathbf{w} as in equation (2), a linear constraint $\sum_{i=1}^n \alpha_i y_i = 0$ will arise when deriving the dual form. This linear constraint forces any optimization method for the dual problem to update at least two variables at a time. For such kind of

problem, the Sequential Minimal Optimization (SMO) method [19] is an often-used solution. But in our case since we do not have this linear constraint, we are able to optimize α one variable at a time, which leads us to a simple dual coordinate descent method [13]. For ease of notation, we rewrite (15) as the following

$$\min_{\gamma \in \mathcal{F}} h(\gamma) = \beta^T \gamma + \frac{1}{2} \gamma^T A \gamma \quad (18)$$

where \mathbf{d}_{J_k} , $\mathcal{P}(\boldsymbol{\alpha}^l)_{J_k}$, $\nabla f^{\mathcal{D}}(\boldsymbol{\alpha}^l, \boldsymbol{\mu})_{J_k}$ and H_{J_k} are replaced with γ , \mathcal{F} , β and A respectively. Since we are dealing with SVM with square hinge loss, $\mathcal{P}(\boldsymbol{\alpha}^l)_{J_k}$ only has lower bound. Let $\mathcal{F}_i = -\alpha_i, i = 1, \dots, |J_k|$ be the lower bound of the i -th α in J_k . Because we do not have the equality constraint, we can update $\boldsymbol{\alpha}$ one variable at a time, suppose we are updating the i -th α , (18) becomes

$$\min_d h(\gamma + d e_i) \equiv \nabla_i h(\gamma) d + \frac{1}{2} A_{ii} d^2 + \text{const} \quad (19)$$

which is a quadratic function of d , where d is a scalar, $\nabla_i h$ is the i -th component of the gradient ∇h . Since the variable γ is constrained by a feasible set \mathcal{F} with lower bound, we need to project the gradient to a feasible direction when optimizing γ . We only need to consider two cases: When $\gamma_i = \mathcal{F}_i$ and $\nabla_i h(\gamma) > 0$, we can not further decrease γ_i and have to clip $\nabla_i h(\gamma)$ to 0. When $\gamma_i = \mathcal{F}_i$ and $\nabla_i h(\gamma) \leq 0$, we can move in the negative direction of the gradient to decrease (19). Denote the i -th component of the projected gradient as $\nabla_i^P h(\gamma)$, we have

$$\nabla_i^P h(\gamma) = \begin{cases} \nabla_i h(\gamma) & \text{if } \mathcal{F}_i < \gamma_i \\ \min(0, \nabla_i h(\gamma)) & \text{if } \mathcal{F}_i = \gamma_i \end{cases}$$

where $\nabla_i h(\gamma)$ can be computed as

$$\nabla_i h(\gamma) = (A\gamma)_i + \beta_i$$

or equivalently,

$$\nabla_i h(\gamma) = \mathbf{y}_i \boldsymbol{\kappa}^T \mathbf{x}_i + \beta_i + \frac{1}{C} \gamma_i = \mathbf{y}_i \sum_{t=1}^T \boldsymbol{\kappa}_t^T \mathbf{x}_{t,i} + \beta_i + \frac{1}{C} \gamma_i \quad (20)$$

where

$$\boldsymbol{\kappa} = \sum_{t=1}^T \mu_t X_{t,J_k}^T Y_{J_k} \gamma, \quad \boldsymbol{\kappa}_t = \mu_t X_{t,J_k}^T Y_{J_k} \gamma, \quad \mathbf{x}_{t,i} = \mathbf{x}_i \odot \mathbf{d}_t$$

The above calculation of $\nabla_i h(\gamma)$ using (20) is only valid in linear SVM where the kernel $K(\mathbf{x}_i, \mathbf{x}_j)$ is simply the dot product $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$. By computing $\nabla_i h(\gamma)$ in this way and caching $\boldsymbol{\kappa}_t$, we have

avoided storing the dense kernel matrix A in the memory, which is more practical. If the projected gradient is 0, then (19) is already at its optimum. If not, we need to update γ_i and clip it to the feasible set. The algorithm is described in Algorithm 2. The general idea is that while γ is not optimal, we loop over all its coordinate entries and try to optimize them with respect to the objective function.

Algorithm 3 Dual Coordinate Descent Method

- 1: Given the initial value and feasible set of γ , compute $\kappa_t = \mu_t X_{t,J_k}^T Y_{J_k} \gamma$, $t = 1, \dots, T$
 - 2: For $i = 1, \dots, |\gamma|$
 - 3: $\nabla_i h(\gamma) = y_i \sum_{t=1}^T \kappa_t^T \mathbf{x}_{t,i} + \beta_i + \frac{1}{C} \gamma_i$
 - 4: $PG \leftarrow \nabla_i^P h(\gamma)$
 - 5: If $PG \neq 0$
 - 6: $\tilde{\gamma}_i \leftarrow \gamma_i$
 - 7: $\gamma_i \leftarrow \max(\gamma_i - \frac{\nabla_i h(\gamma)}{A_{ii}}, \mathcal{F}_i)$
 - 8: $\kappa_t \leftarrow \kappa_t + \sum_{t=1}^T \mu_t (\gamma_i - \tilde{\gamma}_i) y_i \mathbf{x}_{t,i}$, $t = 1, \dots, T$
 - 9: While γ is not optimal, go to step 2
-

Computing the Step Size: According to [12], the step size η^l can be computed as the following

$$\eta^l = \min((\eta^l)^*, \min_{1 \leq i \leq n} \lambda_i) \quad (21)$$

where

$$(\eta^l)^* = \frac{-\nabla f^D(\boldsymbol{\alpha}^l, \boldsymbol{\mu}) \Delta \boldsymbol{\alpha}^l}{(\Delta \boldsymbol{\alpha}^l)^T \bar{Q} \Delta \boldsymbol{\alpha}^l}$$

$$\lambda_i = \begin{cases} \frac{-\alpha_i^l}{\Delta \alpha_i^l} & \text{if } \Delta \alpha_i^l < 0, \\ \infty & \text{if } \Delta \alpha_i^l \geq 0. \end{cases}$$

The parameters needed for computing η^l can be computed from

$$\Delta \mathbf{w}_t^l = \bigoplus_{k=1}^K X_{t,J_k}^T Y_{J_k} \Delta \boldsymbol{\alpha}_{J_k}^l, t = 1, \dots, T$$

$$-\nabla f^D(\boldsymbol{\alpha}^l, \boldsymbol{\mu}) \Delta \boldsymbol{\alpha}^l = \sum_{t=1}^T (\mathbf{w}_t^l)^T \Delta \mathbf{w}_t^l + \frac{1}{C} \bigoplus_{k=1}^K (\boldsymbol{\alpha}_{J_k}^l)^T \Delta \boldsymbol{\alpha}_{J_k}^l - \bigoplus_{k=1}^K \mathbf{e}^T \Delta \boldsymbol{\alpha}_{J_k}^l$$

$$(\Delta \boldsymbol{\alpha}^l)^T \bar{Q} \Delta \boldsymbol{\alpha}^l = \sum_{t=1}^T \|\Delta \mathbf{w}_t^l\|^2 + \frac{1}{C} \bigoplus_{k=1}^K \|\Delta \boldsymbol{\alpha}_{J_k}^l\|^2$$

To compute η^l on Husky, we use *Aggregator* and we use *Parameter Server* to compute $\Delta \mathbf{w}_t$, $t = 1, \dots, T$. All the workers can then individually compute the step size and do the updates using (12) because every worker sees the same value of the parameters in the *Aggregator* and the *Parameter*

Server.

Convergence of α : With the availability of the new α , all the workers then check if α has already been optimized by comparing the primal objective value with the dual objective value. Computing the primal objective value involves iterating all the input instances once, which is expensive in a single machine setting but this problem is greatly alleviated in a distributed setting. If the duality gap is smaller than a pre-determined threshold, then the problem of updating α is solved. If not, all the workers will continue to optimize α . The algorithm is described in Algorithm 4.

Algorithm 4 Box-constrained Quadratic Optimization Algorithm

- 1: $l \leftarrow 0$, given α^0 , computes $w_t^0, t = 1, \dots, T$
 - 2: while α^l is not optimal:
 - 3: Obtain $\Delta\alpha^l$ by distributively solving (15) on K workers in parallel
 - 4: compute $\Delta w_t^l = \bigoplus_{k=1}^K X_{t,J_k}^T Y_{J_k} \Delta\alpha_{J_k}^l, t = 1, \dots, T$
 - 5: compute $\nabla f^{\mathcal{D}}(\alpha^l, \mu)\Delta\alpha^l$ and $(\Delta\alpha^l)\bar{Q}\Delta\alpha^l$
 - 6: obtain η^l by (13)
 - 7: $\alpha^{l+1} \leftarrow \alpha^l + \eta^l \Delta\alpha^l, w_t^{l+1} \leftarrow w_t^l + \eta^l \Delta w_t^l$
 - 8: $l \leftarrow l + 1$.
-

3.2.2 Updating μ

Once α is solved, as described in 2, we first obtain the gradient $\frac{\partial J}{\partial \mu_t} = -\frac{1}{2}\alpha^T(YX_tX_t^TY)\alpha = -\frac{1}{2}w_t^T w_t$ and project this gradient to obtain a feasible descent direction D . Then we use a simple binary search procedure to sample the objective values using half the maximum step size $\frac{\gamma_{max}}{2}$ and repeat this procedure until the step size to be applied is smaller than a pre-determined threshold.

3.2.3 Convergence of SimpleMKL

Once μ has been updated, we test whether the duality gap of the MKL problem is small enough. In particular, the duality gap is measured by:

$$\max_t w_t^T w_t - \sum_{t=1}^T \mu_t w_t^T w_t \leq \epsilon$$

3.3 Convergence of Distributed Feature Generating Machine

Once SimpleMKL converges, we find a new violated constraint $\hat{\mathbf{d}}$ using (11). With this new $\hat{\mathbf{d}}$, we first check if it is in the constraint set \mathcal{C} . If yes, then the algorithm has converged and we can stop. If not, we add $\hat{\mathbf{d}}$ to the constraint set \mathcal{C} and repeat the above procedure for solving the MKL problem. During implementation, other stopping criteria can also be used, for example, one can measure the relative difference between the current MKL objective value and the last MKL objective value.

3.4 Complexity

Let n be the number of input instances, m be the number of features in an instance, l be the average number of non-zero elements in each instance, K be the number of workers in a cluster and B be the number of non-zero entries in each kernel, which typically is set to $\{2, 5, 10, 20, 50, 100\}$, T be the cardinality of the violated constraints set \mathcal{D} (According to [11], a maximum of 10 iterations are enough for Algorithm 1 to converge, which means T is a relatively small number). Since the primary objective of this algorithm is to deal with high dimensional problems, which typically has hundreds of thousands or even millions of features, we assume in this section $l \gg TB$. We discuss in this section about the complexity of this algorithm in lieu of the implementation choices that we mentioned in section 3.2.1 about how to compute \mathbf{w}^l . In particular, we will focus on how different implementations affect the complexity of the algorithms in the Dual Coordinate Descent part and the communication part.

3.4.1 Implementation scheme 1: Do not cache $\mathbf{x}_i \odot \mathbf{d}$

If $\mathbf{x}_i \odot \mathbf{d}$ is not cached, we face a hurdle when computing $\sum_{t=1}^T \mu_t \boldsymbol{\alpha}^T Y X_t X_t^T Y \boldsymbol{\alpha}$. But the feature extraction operator " \odot " has an interesting property: $(\mathbf{x}_i \odot \mathbf{d})^T (\mathbf{x}_j \odot \mathbf{d}) = \mathbf{x}_i^T (\mathbf{x}_j \odot \mathbf{d})$, which is a result of the fact that elements in \mathbf{d} can either be 0 or 1. This property allows us to factor

$\sum_{t=1}^T \mu_t \boldsymbol{\alpha}^T Y X_t X_t^T Y \boldsymbol{\alpha}$ as the following:

$$\begin{aligned}
\sum_{t=1}^T \mu_t \boldsymbol{\alpha}^T Y X_t X_t^T Y \boldsymbol{\alpha} &= \sum_{t=1}^T \mu_t \sum_i \sum_j \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \odot \mathbf{d}_t)^T (\mathbf{x}_j \odot \mathbf{d}_t) \\
&= \sum_{t=1}^T \mu_t \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T (\mathbf{x}_j \odot \mathbf{d}_t) \\
&= \sum_{t=1}^T \mu_t \left(\sum_i \alpha_i y_i \mathbf{x}_i \right)^T \left(\sum_j \alpha_j y_j \mathbf{x}_j \odot \mathbf{d}_t \right) \\
&= \left(\sum_i \alpha_i y_i \mathbf{x}_i \right)^T \left(\sum_j \alpha_j y_j \mathbf{x}_j \odot \left(\sum_{t=1}^T \mu_t \mathbf{d}_t \right) \right) \\
&= \mathbf{w}^T (\mathbf{w} \odot \tilde{\mathbf{d}})
\end{aligned}$$

where $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$ and $\tilde{\mathbf{d}} = \sum_{t=1}^T \mu_t \mathbf{d}_t$ can be cached to speed up the algorithm.

3.4.2 Implementation scheme 2: Cache $\mathbf{x}_i \odot \mathbf{d}$

If $\mathbf{x}_i \odot \mathbf{d}$ is cached, then inside the SVM solver, we can treat it as if we were solving T SVM problems and each one of the SVM problem is associated with a particular kernel $K_t(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i \odot \mathbf{d}_t)^T (\mathbf{x}_j \odot \mathbf{d}_t)$. And for each one of these problems we will train a weight vector \mathbf{w}_t . Then the complete weight vector \mathbf{w} is just the sum of all these T weight vectors, weighted by μ_t .

3.4.3 Time Complexity

Implementation Scheme 1: Using the implementation scheme 1, $\tilde{\mathbf{d}}$ will be a dense vector of size m where each \mathbf{d}_t is only a vector of size B and $B \ll m$. This also means that \mathbf{w} will also be a dense vector of size m . In each Dual Coordinate Descent iteration, we update \mathbf{w} with $\mathbf{w} \leftarrow \mathbf{w} + (\alpha_i - \bar{\alpha}_i) y_i (\mathbf{x}_i \odot \tilde{\mathbf{d}})$. Therefore, each inner iteration of Dual Coordinate Descent algorithm will need approximately $O(\frac{nl}{K})$ operations. After local α_{J_k} is solved, we need an *allreduce* step to synchronize all the workers. In this scheme, we need to synchronize \mathbf{w} and several scalar quantities, resulting in a communication cost of $O(m)$. Notice the use of sparse representation will not help here because \mathbf{w} will be dense vector of size m .

Implementation Scheme 2: Using the implementation scheme 2, we need an extra step of caching the kernel each time we add a new control variable \mathbf{d}_t , which results in a cost of $O(\frac{n(B+l)}{K})$. For the representation of \mathbf{w}_t , we can store \mathbf{w}_t using a sparse representation because

we have the information about \mathbf{d}_t , which tells us exactly which features will likely be non-zero. Since each instance \mathbf{x}_i will have on-average l non-zero elements. $\mathbf{x}_i \odot \mathbf{d}_t$ will have on-average $\frac{l}{m}B$ non-zero elements. Therefore, each iteration of Dual Coordinate Descent method will need approximately $O(\frac{l}{m} \frac{n}{K} TB)$ operations. And synchronizing $\mathbf{w}_t, t = 1 \cdots T$ has a communication cost of $O(TB)$. We thus obtain a speed up of $O(\frac{m}{TB})$. Thus we managed to cut down the running time cost and the communication cost dramatically. Notice in a lot of sense this algorithm is independent of the number of features of the input data because all we are interested in are TB number of useful features. Since we adopt a scheme that caches $\mathbf{x} \odot \mathbf{d}$, we reduce the original high dimensional problem to a low dimensional problem with dimension equal to TB during runtime.

Analysis for both Schemes: The above analysis mainly focuses on Dual Coordinate Descent method and the communication in the Box constrained Quadratic Optimization method. Now we analyze other aspects of the algorithm. In a distributed setting, when finding the most violated constraint, we need to synchronize the weight vectors \mathbf{w} for all the features, which lead to a communication cost of $O(m)$. When we are dealing with high dimensional datasets, this is a very expensive step. And this step becomes more expensive as we increase the number of workers in the cluster. But luckily we do not need to do this step very often as T ranges from 1 to 20. It is proven in [12] that our distributed SVM solver Box constrained Quadratic Optimization method has linear convergence rate. As for SimpleMKL, according to our experiment, typically each iteration composes of 2 search calls for descent direction and 4 search calls for step size. Although this may seem very expensive because for each call of descent direction search or step size search we need to invoke the SVM solver, we can do more caching and initialize the solver with previous values of α to speed up the algorithm. In summary, our algorithm takes advantage of the 0-1 control variable and caching to reduce the a high dimensional problem into a low dimensional one, thus enabling our algorithm to handle data of ultra-high dimension. In addition, we take advantage of a distributed computing framework to scale up our algorithm in the presence of input data with huge volume of input instances.

3.4.4 Space Complexity

Implementation Scheme 1: The space requirement of Scheme 1 is only $O(\frac{nl}{K})$

Implementation Scheme 2: The space requirement of Scheme 2 is $O(\frac{nl+n\frac{l}{m}TB}{K})$, which is not too worse than Scheme 1 considering that TB is a relatively small number. Besides, one of

the main advantages of a distributed algorithms is that we have more memory to spare and we can always add more machines to scale up the total memory available in the cluster. Therefore, we are able to use the nature of a distributed algorithm to cancel out the potential drawbacks caused by caching. In summary, this algorithm is computationally efficient as the Box constrained Quadratic Optimization method has linear convergence rate[12] and each iteration of the algorithms involves only tens of SVM solver call.

4 Experiments

In this section, we evaluate the performance of the single machine FGM algorithm and our distributed FGM algorithm on two synthetic dataset and a collection of real world datasets.

4.1 Datasets

4.1.1 Synthetic Datasets

The synthetic datasets consist of one small dataset (two variants) and one large dataset. Both datasets are binary classification problems with the ground truth known.

SYN-SMALL For the small dataset, we generate the first variant with 10,000 samples and 200 features with each value initialized with a uniform distribution $U(-1, 1)$ and a sparsity ratio of 15%. A ground truth weight vector \mathbf{w} is then generated with 200 features with each value initialized with a uniform distribution $U(-1, 1)$ and sparsity ration of 15%. The dot product $\langle \mathbf{w}, \mathbf{x} \rangle$ is then used to generate the label $y_i \in \{-1, +1\}$. We call this dataset SYN-SMALL.

SYN-SMALL-HIGH-DIM Another variant of the small dataset is generated by appending 19800 irrelevant features to every instances in the first variant with each value initialized with a uniform distribution $U(-1, 1)$ and a sparsity ratio of 15%. We call this dataset SYN-SMALL-HIGH-DIM. Ideally, a feature selection algorithm should be able to recover the relevant features from the first 200 features.

SYN-LARGE The second synthetic dataset is generated with the same setting as the first synthetic dataset except the second one is generated with 1,000,000 samples and 200 features. We call this dataset SYN-LARGE.

4.1.2 Real-World Datasets

The real-world datasets consist of two small datasets and two large datasets. All these four datasets are from the LIBSVM[20] website without further pre-processing. The two small datasets, **BREAST-CANCER** and **LEUKEMIA** are insufficient in the amount of instances. For the **BREAST-CANCER** dataset, we first merge the train set and the test set together and then randomly split them into two datasets of equal size and perform our experiment several times. For the **RCV1.BINARY** dataset, we also randomly split the whole dataset into two dataset of equal size. The **URL** dataset[21] consists of 120 days of URL data used for detecting malicious URL. Since the features are engineered from host-based features, lexical features and other real-valued features using the bag of words model, the number of features goes up to millions! We randomly shuffle the 120 days of data and split them into three datasets of **small**, **large** and **extra large** sizes respectively with details described in table1.

Table 1: Datasets used in the experiments

Dataset	# Training Points	# Testing Points	# Features
SYN-SMALL	10,000	10,000	200
SYN-SMALL-HIGH-DIM	10,000	10,000	20,000
BREAST-CANCER	21	21	7129
LEUKEMIA	38	34	7129
SYN-LARGE	1,000,000	10,000	10,000
RCV1.BINARY	338,700	338,699	47,236
URL-SMALL	159,745	159,740	3,231,961
URL-LARGE	958,460	479,220	3,231,961
URL-EXTRA	1,437,680	479,220	3,231,961

4.2 Experimental Set Up

In our experiments, comparisons are conducted among single machine FGM algorithm and distributed FGM algorithms. We denote distributed FGM run with K workers as DFGM-K. Dual Coordinate Descent SVM algorithm, denoted as DCD-SVM and Box constrained Quadratic Optimization SVM algorithm, denoted as BQO-SVM, from the LIBlinear library [22] serve as baselines, which select all features for classification. In all our experiments, the parameter C for SVM is set to 1 and the loss function for SVM is set to l2-loss.

4.3 Experiments on small datasets

To verify the correctness of our algorithms, two synthetic experiments are performed first. These two synthetic datasets share the same ground truth weight vector \mathbf{w} , with only 29 non-zero entries among the first 200 entries. As seen in Figure 2(a), as the number of kernels (number of selected features) increases, the testing accuracies of both FGM and DFGM increase and converge to the same testing accuracies as DCD-SVM and BQO-SVM when the number of kernels reaches 6 (number of selected features ≥ 29 because $B = 5$). From Figure 2(b), we can see that both FGM and DFGM successfully selects the relevant features among a pool of features mixed with a large number of irrelevant features when DCD-SVM and BQO-SVM suffers from giving weights to these irrelevant features and generalize poorly on the testing set. These two pictures show that our algorithm is effective as a feature selection algorithm in selecting useful features among noisy irrelevant features. We summarize the testing accuracies of the methods that we are interested in on both the synthetic small datasets and real-world small datasets in table 2. This table mainly serves to show the effectiveness of our algorithm in selecting useful features. As the primary purpose of this report is to propose a distributed algorithm that can scale in the presence of large scale high dimensional datasets, we focus on the next section, where we perform experiments on large datasets to show the scalability of our algorithm.

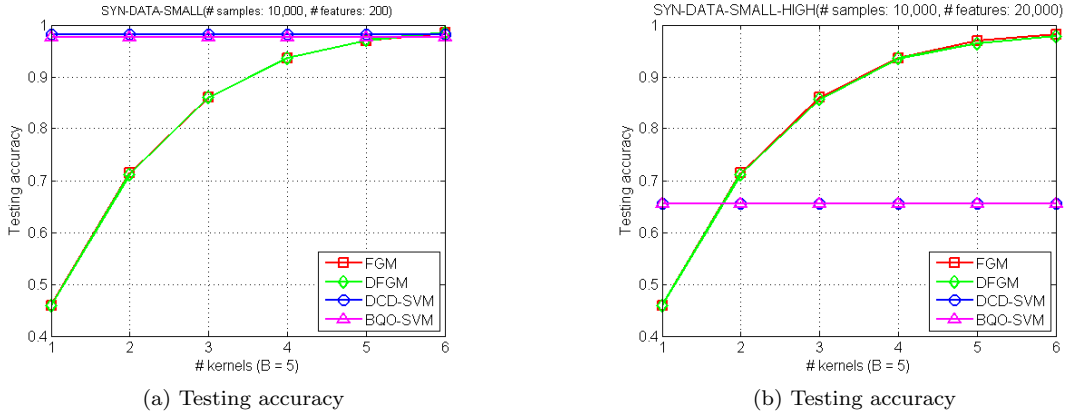


Figure 2: Testing accuracies on small synthetic datasets

4.4 Experiments on large datasets

In this subsection, we show the scalability of our algorithm first by comparing the performances of DFGM running on different number of worker machines with the performance of FGM on SYN-DATA-LARGE, RCV1.BINARY and URL-SMALL datasets. We then compare the

Table 2: Testing accuracies of various algorithms on small datasets

DATA SETS	DCD-SVM	BQO-SVM	FGM		
			B	DFGM	
SYN-DATA-SMALL	0.9809	0.9761	2	0.9618	0.9767
	0.9809	0.9761	5	0.9686	0.9831
	0.9809	0.9761	10	0.9850	0.9807
SYN-DATA-SMALL-HIGH	0.6550	0.6559	2	0.9618	0.9772
	0.6550	0.6559	5	0.9816	0.9784
	0.6550	0.6559	10	0.9704	0.9661
BREAST-CANCER	0.8095	0.8095	2	0.7619	0.8095
	0.8095	0.8095	5	0.8095	0.8095
	0.8095	0.8095	10	0.7619	0.7619
	0.8095	0.8095	20	0.8571	0.9048
LEUKEMIA	0.7941	0.7941	2	0.9412	0.9412
	0.7941	0.7941	5	0.9412	1.0000
	0.7941	0.7941	10	0.9706	0.9706
	0.7941	0.7941	20	0.9412	0.9706

performances among DFGM running on different number of worker machines on **URL-LARGE** and **URL-EXTRA** datasets. We separate the experiments into these two parts because FGM can only run on single machine, which makes it infeasible for very large datasets with millions of samples.

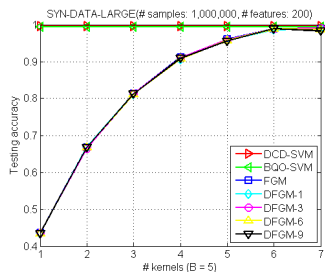


Figure 3: SYN-DATA-LARGE

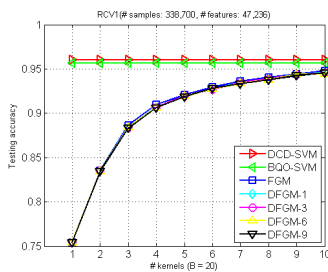


Figure 4: RCV1.BINARY

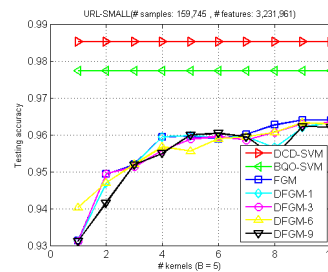


Figure 5: URL-SMALL

Figure 6: Testing accuracies on large datasets

4.4.1 FGM v.s. DFGM-K

Figure 3 summarizes the testing accuracies of various methods on the three large datasets respectively. Figure 7-9 summarizes the total running time and the speedup of DFGM compared to FGM. Finally, table 3 gives detailed information about the training time and communication time on the **URL** datasets. There are two simple observations that can be made from Figure 7-9. First, as expected, the training time reduces as we increase the number of worker machines in our cluster, which is what we expect of a distributed algorithm. Second, the speedup resulted by

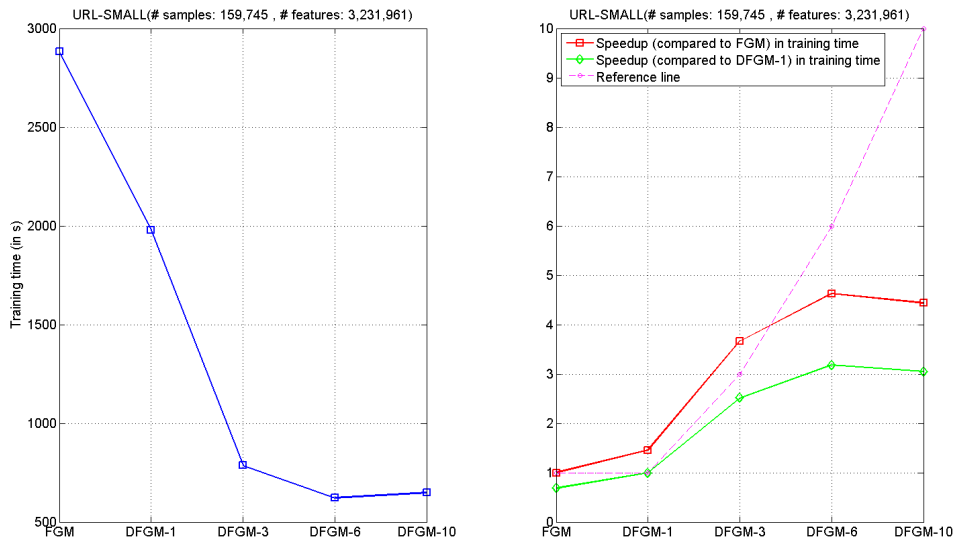


Figure 7: URL-SMALL

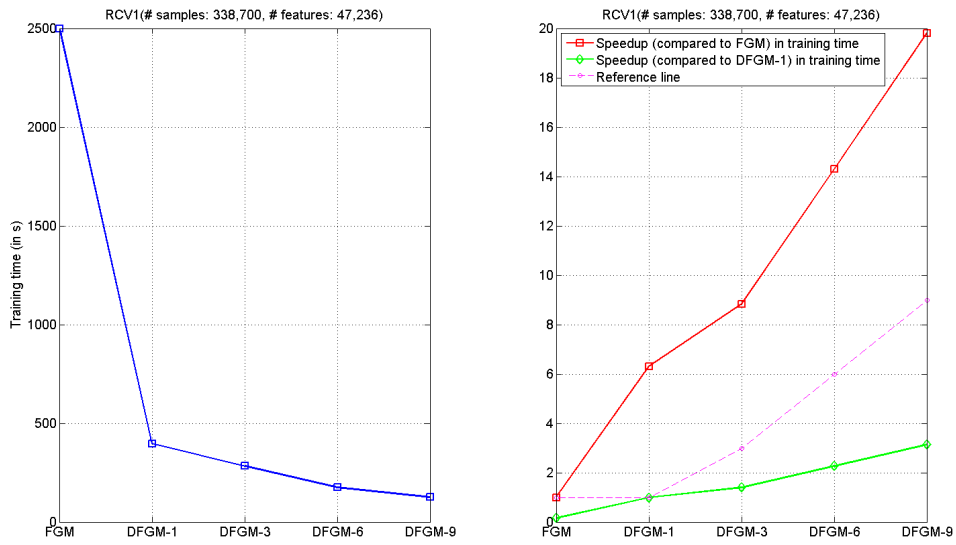


Figure 8: RCV1.BINARY

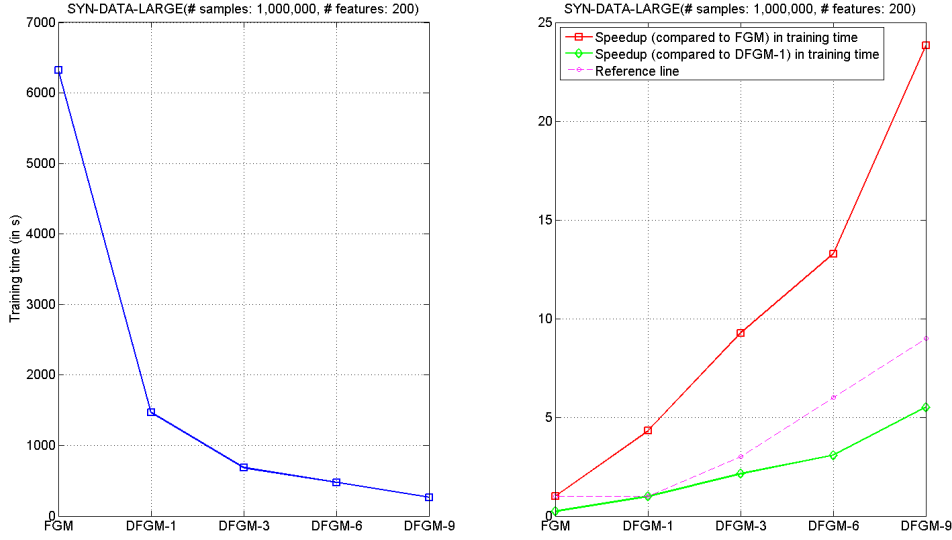


Figure 9: SYN-DATA-LARGE

increasing the number of worker machine is not linear but sub-linear.

Table 3: Training time and communication time on **URL** datasets

DATA SETS	DFGM-K	training time	communication time	$\frac{\text{communication time}}{\text{training time}}$
URL-SMALL	DFGM-1	1978.68	0	0%
	DFGM-3	786.40	96.35	12.25%
	DFGM-6	622.62	167.00	26.66%
	DFGM-10	649.45	256.58	39.51%
URL-LARGE	DFGM-5	1046.45	33.91	3.24%
	DFGM-10	707.76	55.14	7.79%
	DFGM-15	564.70	79.06	14.00%
URL-EXTRA	DFGM-5	2044.01	40.42	1.98%
	DFGM-10	1302.82	67.15	5.15%
	DFGM-15	969.79	92.68	9.57%

4.4.2 Comparisons among DFGM-K

Figure 10 shows our experiments of running 5, 10, 15 and 20 worker machines on **URL-LARGE** and **URL-EXTRA** respectively. This experiment is conducted mainly to investigate the scalability of our distributed algorithms in the presence of large amount of data. As shown in the figure, having a larger dataset (**URL-EXTRA** approximately 1.5 times larger than **URL-LARGE**) that is still within the capacity of the cluster does make the speed up curve steeper, indicating that the algorithm can scale in the presence of larger amount of data and more worker machines.

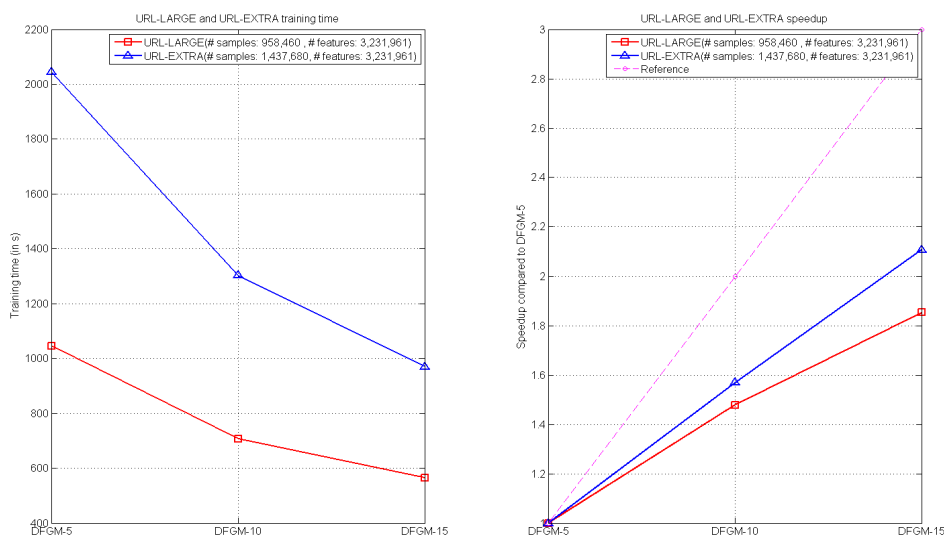


Figure 10: URL-LARGE and URL-EXTRA

5 Discussions

The simple observations and figures above show that our algorithm can scale in the presence of more data and more worker machines. However, the scaling factor, which is at the best in the **SYN-DATA-LARGE** dataset, achieving 0.6 (5.5 speedup with 9 worker machines), is still relatively small. In this section we discuss some of the issues that might have caused this problem and possible improvements. In the mean time, we further discuss in details about the results obtained in the experiments section.

5.1 Load Balancing and Data Heterogeneity

In a distributed computing setting, it is usually assumed that computing nodes are homogeneous in their capability. But this assumption is easily violated in a real-world setting where the users only communicate with the resource manager to negotiate computing resources for their jobs. It is possible that the some machines are running several different tasks submitted by different users while others are mostly idle, thus violating the nodes homogeneity assumptions. This heterogeneity not only affects the available computing performances at run-time but also the amount of data handled by each worker machine. Since Husky adopts a Master-Slave cluster architecture, during run time the Slave machine sends a request to the Master requesting locations of input data

and the Master responds with the block location in HDFS[9]. Thus if a worker machine is slow, then it may acquire far less data compared to other worker machines. For examples, in our experiments, it is observed that some slow workers can have as high as 25% less data than the worker machines holding the most data. Therefore, the difference of amount of data processed by each worker machine and the loading status of each worker machine causes the "stragglers" [23] in our experiments, which in turn hurts the scalability performance of our algorithm. Another issues that affect the scaling factor of our algorithm is the inherent heterogeneity of the input data. Intuitively, for a SVM problem, some samples, which at the optimum correspond to the "Support Vectors" are harder to optimize and some samples, which at the optimum have zero-weight in the α , are easier to optimize. This is evident from the fact that the scaling factor is at the best on the **SYN-DATA-LARGE** dataset. Intuitive this is so because the **SYN-DATA-LARGE** dataset is generated by a uniform distribution and each sample is relatively identical to each other and **URL** and **RCV1** are real world datasets, which contain more realistic and extreme cases. While the second issue is related to the input data and may not be easily solved without further preprocessing of the data, the first issue can be alleviated by doing a load balancing between different workers before running the DFGM algorithm. In particular, we can make use of the *Migrate Channel* in Husky to allow highly loaded worker machines to offload some of the samples to less-loaded worker machines. By doing load balancing, we can better utilize the computing resources available and bridging the gap between normal workers and the "stragglers", thereby reducing the average waiting time during the communication procedure.

5.2 Number of worker machines

From the right picture of Figure 7-9 we can see that approximately it takes two times or even three times amount of worker machines to reduce the training time to half on **RCV1** and **URL-SMALL** datasets when compared to DFGM-1. In particular, we can observe that DFGM-10 is worse than DFGM-6 on the **URL-SMALL** dataset. To understand why this is so, we refer the readers to table 3 where it can be seen from the first row of the table that DFGM-10 spends 40% of the time in communication, which is a clear indication that with only 159,745 samples in the **URL-SMALL** dataset, having 10 worker machines in the cluster is an overkill and hurts the performance. Furthermore, we can see from **RCV1**, **SYN-DATA-LARGE**, **URL-LARGE** and **URL-EXTRA**, which has more samples than **URL-SMALL**, that the speedup curve gets steeper as we increase the number of samples in the data given the same number of worker

machines. Thus the lesson learned is that when running a distributed machine learning algorithm where communication through network is frequent and unavoidable, we need to tune the number of worker machines according to the size of the input data rather than blindly increasing the number of worker machines in the cluster.

5.3 SVM-Solver

From the comparisons between FGM and DFGM-K, we can see that our distributed algorithm running on only one machine, performs extremely well when compared to FGM. We investigate into this problem and find out that the difference of performance is due to the differences of the SVM solver that these two methods use. In particular, FGM uses DCD-SVM, which circularly goes through the training dataset and use projected gradient information to update the Lagrange multiplier α_i until all the α_i s are at optimum, at which point the projected gradient vanishes to a zero vector. In comparison, DFGM uses BQO-SVM, which also circularly goes through the training dataset but considers second order approximation. For convergence criterion, BQO-SVM stops when the first order term of the Taylor series expansion of the objective function with respect to α becomes greater than 0. Although each iteration of BQO-SVM is slower due to the use of second order approximation, we find out in our experiments that BQO-SVM usually takes far less iterations to converge when compared to DCD-SVM. For example, DCD-SVM takes 72 iterations, totaling 50 seconds to converge and achieve a classification accuracy of 97.74% on the testing set. In comparison, BQO-SVM takes only 2 iterations, totaling 8 seconds to converge on a solution that yields classification accuracy of 97.41% on the testing set. From this example, we can see that DCD-SVM actually spends a majority of the training time achieving a solution that is not significantly better than the one found by BQO-SVM. Since the SimpleMKL algorithm performs several steps of searching of the parameter using the SVM solver in each iteration, sometimes the solution returned from the solver may simply be discarded due to some other constraints. Thus the problem of the DCD-SVM solver gets magnified during the many iterations of the SimpleMKL algorithm. Therefore, it may not be a good idea to spend too much computational resources in finding a solution that brings very little or maybe no returns and this explains why DFGM-1 performs much better than FGM in our experiment. Although it is not done in this report, this problem can be avoided if we modify the solver to use relative changes in the objective values as a stopping criterion.

5.4 Trade-off between Communication and Local Computation

As mentioned in section 3.2.1, each iteration of BQO-SVM involves solving a local SVM problem using the DCD-SVM solver and a communication procedure to synchronize with other worker machines. Thus there is a trade-off between communication cost and local computational cost. If we allow the DCD-SVM more time to find a better solution, then the communication cost will be smaller. On the other hand, if we stop DCD-SVM prematurely, there may be more noises/perturbations in our descent direction and the BQO-SVM algorithm will need more iteration to converge, thus increasing the communication cost. Since for the datasets used in our experiments, BQO-SVM typically takes a few iterations to converge, we set the maximum iteration for DCD-SVM to 1 like the implementation given by [12], which yields faster convergence for BQO-SVM but incurs more communication cost. A possible improvement should be to adopt a more general stopping criterion for DCD-SVM by taking into consideration the maximum number of iterations, the duality gap and the relative changes in the objective values. For tuning of the parameters for the general stopping criterion, we can sample a portion of the input data and list out a similar table as table 3 to determine whether the training time is dominated by communication or by local computation. If the training time is dominated by communication, we can impose a more stringent stopping criterion for DCD-SVM and vice versa.

6 Conclusion

In this report, we propose a distributed algorithm to learn a Sparse SVM for feature selection on high dimensional data. By introducing a 0-1 control variable, our method achieves the goal of sparsity while being computationally efficient. The use of 0-1 control variable allows our algorithm to be unaware of the number of features in the input data. By taking advantages of the linear Kernel and cleverly caching important data, we design an efficient and scalable distributed solver for the SVM problem with multiple linear kernels. Finally, by leveraging an efficient distributed computing framework - Husky, and making an implementation choice that reduces the amount of communication between workers, we only need to pay a communication cost that is linear in the number of desired features to scale up the performance in the presence of large scale high dimensional data.

References

- [1] Andrew McAfee, Erik Brynjolfsson, Thomas H Davenport, DJ Patil, and Dominic Barton. Big data. *The management revolution. Harvard Bus Rev*, 90(10):61–67, 2012.
- [2] Isabelle Guyon and André Elisseeff. An introduction to variable and feature selection. *Journal of machine learning research*, 3(Mar):1157–1182, 2003.
- [3] Jason Weston, André Elisseeff, Bernhard Schölkopf, and Mike Tipping. Use of the zero-norm with linear models and kernel methods. *Journal of machine learning research*, 3(Mar):1439–1461, 2003.
- [4] Paul S Bradley and Olvi L Mangasarian. Feature selection via concave minimization and support vector machines. In *ICML*, volume 98, pages 82–90, 1998.
- [5] Isabelle Guyon, Jason Weston, Stephen Barnhill, and Vladimir Vapnik. Gene selection for cancer classification using support vector machines. *Machine learning*, 46(1-3):389–422, 2002.
- [6] Mingkui Tan, Ivor W Tsang, and Li Wang. Towards ultrahigh dimensional feature selection for big data. *Journal of Machine Learning Research*, 15(1):1371–1429, 2014.
- [7] Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J Franklin, Scott Shenker, and Ion Stoica. Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing. In *Proceedings of the 9th USENIX conference on Networked Systems Design and Implementation*, pages 2–2. USENIX Association, 2012.
- [8] Michael Isard, Mihai Budiu, Yuan Yu, Andrew Birrell, and Dennis Fetterly. Dryad: distributed data-parallel programs from sequential building blocks. In *ACM SIGOPS Operating Systems Review*, volume 41, pages 59–72. ACM, 2007.
- [9] Apache Hadoop. <http://hadoop.apache.org/>.
- [10] Fan Yang, Jinfeng Li, and James Cheng. Husky: Towards a more efficient and expressive distributed computing framework. *Proceedings of the VLDB Endowment*, 9(5):420–431, 2016.
- [11] Mingkui Tan, Li Wang, and Ivor W Tsang. Learning sparse svm for feature selection on very high dimensional datasets. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 1047–1054, 2010.

- [12] Ching-Pei Lee and Dan Roth. Distributed box-constrained quadratic optimization for dual linear svm. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 987–996, 2015.
- [13] Cho-Jui Hsieh, Kai-Wei Chang, Chih-Jen Lin, S Sathiya Keerthi, and Sellamanickam Sundararajan. A dual coordinate descent method for large-scale linear svm. In *Proceedings of the 25th international conference on Machine learning*, pages 408–415. ACM, 2008.
- [14] Bernhard E Boser, Isabelle M Guyon, and Vladimir N Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the fifth annual workshop on Computational learning theory*, pages 144–152. ACM, 1992.
- [15] Bernhard Schölkopf. The kernel trick for distances. *Advances in neural information processing systems*, 13:301–307, 2001.
- [16] Alain Rakotomamonjy, Francis R Bach, Stéphane Canu, and Yves Grandvalet. Simplekl. *Journal of Machine Learning Research*, 9(Nov):2491–2521, 2008.
- [17] Matei Zaharia, Reynold S. Xin, Patrick Wendell, Tathagata Das, Michael Armbrust, Ankur Dave, Xiangrui Meng, Josh Rosen, Shivaram Venkataraman, Michael J. Franklin, Ali Ghodsi, Joseph Gonzalez, Scott Shenker, and Ion Stoica. Apache spark: A unified engine for big data processing. *Commun. ACM*, 59(11):56–65, October 2016.
- [18] Ashish Thusoo, Joydeep Sen Sarma, Namit Jain, Zheng Shao, Prasad Chakka, Suresh Anthony, Hao Liu, Pete Wyckoff, and Raghotham Murthy. Hive: A warehousing solution over a map-reduce framework. *Proc. VLDB Endow.*, 2(2):1626–1629, August 2009.
- [19] John Platt et al. Sequential minimal optimization: A fast algorithm for training support vector machines. 1998.
- [20] Chih-Chung Chang and Chih-Jen Lin. Libsvm: A library for support vector machines. *ACM Trans. Intell. Syst. Technol.*, 2(3):27:1–27:27, May 2011.
- [21] Justin Ma, Lawrence K. Saul, Stefan Savage, and Geoffrey M. Voelker. Identifying suspicious urls: An application of large-scale online learning. In *Proceedings of the 26th Annual International Conference on Machine Learning, ICML '09*, pages 681–688, New York, NY, USA, 2009. ACM.

- [22] Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. Liblinear: A library for large linear classification. *Journal of machine learning research*, 9(Aug):1871–1874, 2008.
- [23] Jeffrey Dean and Sanjay Ghemawat. Mapreduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1):107–113, 2008.