

# Recent Advances of Large-scale Linear Classification

Guo-Xun Yuan, Chia-Hua Ho, and Chih-Jen Lin

## Abstract

Linear classification is a useful tool in machine learning and data mining. For some data in a rich dimensional space, the performance (i.e., testing accuracy) of linear classifiers has shown to be close to that of nonlinear classifiers such as kernel methods, but training and testing speed is much faster. Recently, many research works have developed efficient optimization methods to construct linear classifiers and applied them to some large-scale applications. In this paper, we give a comprehensive survey on the recent development of this active research area.

## Index Terms

Large linear classification, Support vector machines, Logistic regression, Multi-class classification.

## I. INTRODUCTION

Linear classification is a useful tool in machine learning and data mining. In contrast to nonlinear classifiers such as kernel methods, which map data to a higher dimensional space, linear classifiers directly work on data in the original input space. While linear classifiers fail to handle some inseparable data, they may be sufficient for data in a rich dimensional space. For example, linear classifiers have shown to give competitive performances on document data with nonlinear classifiers. An important advantage of linear classification is that training and testing procedures are much more efficient. Therefore, linear classification can be very useful for some large-scale applications. Recently, the research on linear classification has been a very active topic. In this paper, we give a comprehensive survey on the recent advances.

We begin with explaining in Section II why linear classification is useful. The differences between linear and nonlinear classifiers are described. Through experiments, we demonstrate that for some data, a linear classifier achieves comparable accuracy to a nonlinear one, but both training and testing time is much shorter. Linear classifiers cover popular methods such as support vector machines (SVM) [1], [2], logistic regression (LR),<sup>1</sup> and others. In Section III, we show optimization problems of these methods and discuss their differences.

An important goal of the recent research on linear classification is to develop fast optimization algorithms for training (e.g., [4]–[6]). In Section IV, we discuss issues in finding a suitable algorithm and give details of some representative algorithms. Methods such as SVM and LR were originally proposed for two-class problems. Although past works have studied their extensions to multi-class problems, the focus was on nonlinear classification. In Section V, we systematically compare methods for multi-class linear classification.

Linear classification can be further applied to many other scenarios. We investigate some examples in Section VI. In particular, we show that linear classifiers can be effectively employed to either directly or indirectly approximate nonlinear classifiers. In Section VII, we discuss an ongoing research topic for data larger than memory or disk capacity. Existing algorithms often fail to handle such data because of assuming that data can be stored in a single computer's memory. We present some methods which try to reduce data reading or communication time. In Section VIII, we briefly discuss related topics such as structured learning and large-scale linear regression.

Finally, Section IX concludes this survey paper.

G.-X. Yuan, C.-H. Ho, and C.-J. Lin are with Department of Computer Science, National Taiwan University, Taipei 10617, Taiwan. Email: {r96042,b95082,cjlin}@csie.ntu.edu.tw

<sup>1</sup>It is difficult to trace the origin of logistic regression, which can be dated back to 19th century. Interested readers may check the investigation in [3].

## II. WHY IS LINEAR CLASSIFICATION USEFUL?

Given training data  $(y_i, \mathbf{x}_i) \in \{-1, +1\} \times \mathbf{R}^n, i = 1, \dots, l$ , where  $y_i$  is the label and  $\mathbf{x}_i$  is the feature vector, some classification methods construct the following decision function.

$$d(\mathbf{x}) \equiv \mathbf{w}^T \phi(\mathbf{x}) + b, \quad (1)$$

where  $\mathbf{w}$  is the weight vector and  $b$  is an intercept, or called the bias. A *nonlinear* classifier maps each instance  $\mathbf{x}$  to a higher dimensional vector  $\phi(\mathbf{x})$  if data are not linearly separable. If  $\phi(\mathbf{x}) = \mathbf{x}$  (i.e., data points are not mapped), we say (1) is a *linear* classifier. Because nonlinear classifiers use more features, generally they perform better than linear classifiers in terms of prediction accuracy.

For nonlinear classification, evaluating  $\mathbf{w}^T \phi(\mathbf{x})$  can be expensive because  $\phi(\mathbf{x})$  may be very high dimensional. Kernel methods (e.g., [2]) were introduced to handle such a difficulty. If  $\mathbf{w}$  is a linear combination of training data, i.e.,

$$\mathbf{w} \equiv \sum_{i=1}^l \alpha_i \phi(\mathbf{x}_i) \text{ for some } \boldsymbol{\alpha} \in \mathbf{R}^l, \quad (2)$$

and the following kernel function can be easily calculated

$$K(\mathbf{x}_i, \mathbf{x}_j) \equiv \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j),$$

then the decision function can be calculated by

$$d(\mathbf{x}) \equiv \sum_{i=1}^l \alpha_i K(\mathbf{x}_i, \mathbf{x}) + b, \quad (3)$$

regardless of the dimensionality of  $\phi(\mathbf{x})$ . For example,

$$K(\mathbf{x}_i, \mathbf{x}_j) \equiv (\mathbf{x}_i^T \mathbf{x}_j + 1)^2 \quad (4)$$

is the degree-2 polynomial kernel with

$$\begin{aligned} \phi(\mathbf{x}) = [1, \sqrt{2}x_1, \dots, \sqrt{2}x_n, \dots, x_1^2, \dots, x_n^2, \\ \sqrt{2}x_1x_2, \sqrt{2}x_1x_3, \dots, \sqrt{2}x_{n-1}x_n] \in \mathbf{R}^{(n+2)(n+1)/2}. \end{aligned} \quad (5)$$

This kernel trick makes methods such as SVM or kernel LR practical and popular; however, for large data, the training and testing processes are still time consuming. For a kernel like (4), the cost of predicting a testing instance  $\mathbf{x}$  via (3) can be up to  $O(ln)$ . In contrast, without using kernels,  $\mathbf{w}$  is available in an explicit form, so we can predict an instance by (1). With  $\phi(\mathbf{x}) = \mathbf{x}$ ,

$$\mathbf{w}^T \phi(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

costs only  $O(n)$ . It is also known that training a linear classifier is more efficient. Therefore, while a linear classifier may give inferior accuracy, it often enjoys faster training and testing.

We conduct an experiment to compare linear SVM and nonlinear SVM (with the RBF kernel). Table I shows the accuracy and training/testing time. Generally, nonlinear SVM has better accuracy, especially for problems `cod-RNA`,<sup>2</sup> `ijcnn1`, `covtype`, `webspam`, and `MNIST38`. This result is consistent with the theoretical proof that SVM with RBF kernel and suitable parameters gives at least as good accuracy as linear kernel [10]. However, for problems with large numbers of features, i.e. `real-sim`, `rcv1`, `astro-physic`, `yahoo-japan`, and `news20`, the accuracy values of linear and nonlinear SVMs are similar. Regarding training and testing time, Table I clearly indicates that linear classifiers are at least an order of magnitude faster.

In Table I, problems for which linear classifiers yield comparable accuracy to nonlinear are all document sets. In the area of document classification and natural language processing (NLP), a bag-of-word model is commonly used to generate feature vectors [11]. Each feature, corresponding to a word, indicates the existence of the word in a document. Because the number of features is the same as the number of possible words, the dimensionality is huge and the data set is often sparse. For this type of large sparse data, linear classifiers are very useful because of competitive accuracy and very fast training and testing.

<sup>2</sup>In this experiment, we scaled `cod-RNA` feature-wisely to  $[-1, 1]$  interval.

TABLE I

COMPARISON OF LINEAR AND NONLINEAR CLASSIFIERS. FOR LINEAR, WE USE THE SOFTWARE LIBLINEAR [7], WHILE FOR NONLINEAR WE USE LIBSVM [8] (RBF KERNEL). THE LAST COLUMN SHOWS THE ACCURACY DIFFERENCE BETWEEN LINEAR AND NONLINEAR CLASSIFIERS. TRAINING AND TESTING TIME IS IN SECONDS. THE EXPERIMENTAL SETTING FOLLOWS EXACTLY FROM [9, SECTION 4].

Data set	#instances		#features	Linear			Nonlinear (kernel)			Accuracy difference to nonlinear
	Training	Testing		Time (s) Training	Time (s) Testing	Testing accuracy	Time (s) Training	Time (s) Testing	Testing accuracy	
cod-RNA	59,535	271,617	8	3.1	0.05	70.71	80.2	126.02	96.67	-25.96
ijcnn1	49,990	91,701	22	1.7	0.01	92.21	26.8	20.29	98.69	-6.48
covtype	464,810	116,202	54	1.5	0.03	76.37	46,695.8	1,131.20	96.11	-19.74
webspam	280,000	70,000	254	26.8	0.04	93.35	15,681.8	853.34	99.26	-5.91
MNIST38	11,982	1,984	752	0.2	0.01	96.82	38.1	5.61	99.70	-2.88
real-sim	57,848	14,461	20,958	0.3	0.01	97.44	938.3	81.94	97.82	-0.38
rcv1	20,242	677,399	47,236	0.1	0.43	96.26	108.0	3,259.46	96.50	-0.24
astro-physic	49,896	12,473	99,757	0.3	0.01	97.09	735.7	111.59	97.31	-0.22
yahoo-japan	140,963	35,240	832,026	3.3	0.03	92.63	20,955.2	1,890.83	93.31	-0.68
news20	15,997	3,999	1,355,191	1.2	0.03	96.95	383.2	100.38	96.90	0.05

### III. BINARY LINEAR CLASSIFICATION METHODS

To generate a decision function (1), linear classification involves the following risk minimization problem.

$$\min_{\mathbf{w}, b} f(\mathbf{w}, b) \equiv r(\mathbf{w}) + C \sum_{i=1}^l \xi(\mathbf{w}, b; \mathbf{x}_i, y_i), \quad (6)$$

where  $r(\mathbf{w})$  is the regularization term and  $\xi(\mathbf{w}, b; \mathbf{x}, y)$  is the loss function associated with the observation  $(y, \mathbf{x})$ . Parameter  $C > 0$  is user-specified for balancing  $r(\mathbf{w})$  and the sum of losses.

Following the discussion in Section II, linear classification is often applied to data with many features, so the bias term  $b$  may not be needed in practice. Experiments in [12], [13] on document data sets showed similar performances with/without the bias term. In the rest of this paper, we omit the bias term  $b$ , so (6) is simplified to

$$\min_{\mathbf{w}} f(\mathbf{w}) \equiv r(\mathbf{w}) + C \sum_{i=1}^l \xi(\mathbf{w}; \mathbf{x}_i, y_i) \quad (7)$$

and the decision function becomes  $d(\mathbf{x}) \equiv \mathbf{w}^T \mathbf{x}$ .

#### A. Support Vector Machines and Logistic Regression

In (7), the loss function is used to penalize a wrongly-classified observation  $(\mathbf{x}, y)$ . There are three common loss functions considered in the literature of linear classification.

$$\xi_{L1}(\mathbf{w}; \mathbf{x}, y) \equiv \max(0, 1 - y\mathbf{w}^T \mathbf{x}), \quad (8)$$

$$\xi_{L2}(\mathbf{w}; \mathbf{x}, y) \equiv \max(0, 1 - y\mathbf{w}^T \mathbf{x})^2, \quad \text{and} \quad (9)$$

$$\xi_{LR}(\mathbf{w}; \mathbf{x}, y) \equiv \log(1 + e^{-y\mathbf{w}^T \mathbf{x}}). \quad (10)$$

Eqs. (8) and (9) are referred to as L1 and L2 losses, respectively. Problem (7) using (8) and (9) as the loss function is often called L1-loss and L2-loss SVM, while problem (7) using (10) is referred to as logistic regression (LR). Both SVM and LR are popular classification methods. The three loss functions in (8)–(10) are all convex and non-negative. L1 loss is not differentiable at the point  $y\mathbf{w}^T \mathbf{x} = 1$ , while L2 loss is differentiable, but not twice differentiable [14]. For logistic loss, it is twice differentiable. Figure 1 shows that these three losses are increasing functions of  $-y\mathbf{w}^T \mathbf{x}$ . They slightly differ in the amount of penalty imposed.

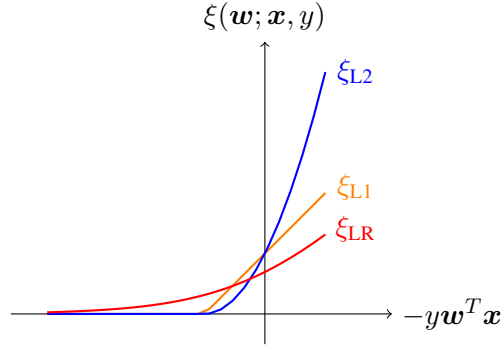


Fig. 1. Three loss functions:  $\xi_{L1}$ ,  $\xi_{L2}$ , and  $\xi_{LR}$ . The  $x$ -axis is  $-y\mathbf{w}^T \mathbf{x}$ .

### B. L1 and L2 Regularization

A classifier is used to predict the label  $y$  for a hidden (testing) instance  $\mathbf{x}$ . Overfitting training data to minimize the training loss may not imply that the classifier gives the best testing accuracy. The concept of regularization is introduced to prevent from overfitting observations. The following L2 and L1 regularization terms are commonly used.

$$r_{L2}(\mathbf{w}) \equiv \frac{1}{2} \|\mathbf{w}\|_2^2 = \frac{1}{2} \sum_{j=1}^n w_j^2 \quad \text{and} \quad (11)$$

$$r_{L1}(\mathbf{w}) \equiv \|\mathbf{w}\|_1 = \sum_{j=1}^n |w_j|. \quad (12)$$

Problem (7) with L2 regularization and L1 loss is the standard SVM proposed in [1]. Both (11) and (12) are convex and separable functions. The effect of regularization on a variable is to push it toward zero. Then, the search space of  $\mathbf{w}$  is more confined and overfitting may be avoided. It is known that an L1-regularized problem can generate a sparse model with few non-zero elements in  $\mathbf{w}$ . Note that  $w^2/2$  becomes more and more flat toward zero, but  $|w|$  is uniformly steep. Therefore, an L1-regularized variable is easier to be pushed to zero, but a caveat is that (12) is not differentiable. Because non-zero elements in  $\mathbf{w}$  may correspond to useful features [15], L1 regularization can be applied for feature selection. In addition, less memory is needed to store  $\mathbf{w}$  obtained by L1 regularization. Regarding testing accuracy, comparisons such as [13, Supplementary Materials Section D] show that L1 and L2 regularization generally give comparable performance.

In statistics literature, a model related to L1 regularization is LASSO [16].

$$\begin{aligned} \min_{\mathbf{w}} \quad & \sum_{i=1}^l \xi(\mathbf{w}; \mathbf{x}_i, y_i) \\ \text{subject to} \quad & \|\mathbf{w}\|_1 \leq K, \end{aligned} \quad (13)$$

where  $K > 0$  is a parameter. This optimization problem is equivalent to (7) with L1 regularization. That is, for a given  $C$  in (7), there exists  $K$  such that (13) gives the same solution as (7). The explanation for this relationship can be found in, for example, [17].

Any combination of the above-mentioned two regularizations and three loss functions has been well studied in linear classification. Of them, L2-regularized L1/L2-loss SVM can be geometrically interpreted as maximum margin classifiers. L1/L2-regularized LR can be interpreted in a Bayesian view by maximizing the posterior probability with Laplacian/Gaussian prior of  $\mathbf{w}$ .

A convex combination of L1 and L2 regularizations forms the elastic net [18].

$$r_e(\mathbf{w}) \equiv \lambda \|\mathbf{w}\|_2^2 + (1 - \lambda) \|\mathbf{w}\|_1, \quad (14)$$

where  $\lambda \in [0, 1)$ . The elastic net is used to break the following limitations of L1 regularization. First, L1 regularization term is not strictly convex, so the solution may not be unique. Second, for two highly correlated

features, the solution obtained by L1-regularization may select only one of these features. Consequently, L1 regularization may discard the group effect of variables with high correlation [18].

#### IV. TRAINING TECHNIQUES

To obtain the model  $w$ , in the training phase we need to solve the convex optimization problem (7). Although many convex optimization methods are available, for large linear classification, we must carefully consider some factors in designing a suitable algorithm. In this section, we first discuss these design issues and follow by showing details of some representative algorithms.

##### A. Issues in Finding Suitable Algorithms

- **Data property** Algorithms that are efficient for some data sets may be slow for others. We must take data properties into account in selecting algorithms. For example, we can check if the number of instances is much larger than features, or vice versa. Other useful properties include the number of non-zero feature values, feature distribution, and feature correlation, etc.
- **Optimization formulation** Algorithm design is strongly related to the problem formulation. For example, most unconstrained optimization techniques can be applied to L2-regularized logistic regression, while specialized algorithms may be needed for the non-differentiable L1-regularized problems.

In some situations, by reformulation, we are able to transform a non-differentiable problem to be differentiable. For example, by letting  $w = w^+ - w^-$  ( $w^+, w^- \geq 0$ ), L1-regularized classifiers can be written as

$$\begin{aligned} \min_{w^+, w^-} \quad & \sum_{j=1}^n w_j^+ + \sum_{j=1}^n w_j^- + \sum_{i=1}^l \xi(w^+ - w^-; \mathbf{x}_i, y_i) \\ \text{subject to} \quad & w_j^+, w_j^- \geq 0, \quad j = 1, \dots, n. \end{aligned} \quad (15)$$

However, there is no guarantee that solving a differentiable form is faster. Recent comparisons [13] show that for L1-regularized classifiers, methods directly minimizing the non-differentiable form are often more efficient than those solving (15).

- **Solving primal or dual problems** Problem (7) has  $n$  variables. In some applications, the number of instances  $l$  is much smaller than the number of features  $n$ . By Lagrangian duality, a dual problem of (7) has  $l$  variables. If  $l \ll n$ , solving the dual form may be easier due to the smaller number of variables. Further, in some situations, the dual problem possesses nice properties not in the primal form. For example, the dual problem of the standard SVM (L2-regularized L1-loss SVM) is the following quadratic program.<sup>3</sup>

$$\begin{aligned} \min_{\alpha} \quad & f^D(\alpha) \equiv \frac{1}{2} \alpha^T Q \alpha - e^T \alpha \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C, \quad \forall i = 1, \dots, l, \end{aligned} \quad (16)$$

where  $Q_{ij} \equiv y_i y_j \mathbf{x}_i^T \mathbf{x}_j$ . Although the primal objective function is non-differentiable because of the L1 loss, in (16), the dual objective function is smooth (i.e., derivatives of all orders are available). Hence, solving the dual problem may be easier than primal because we can apply differentiable optimization techniques. Note that the primal optimal  $w$  and the dual optimal  $\alpha$  satisfy the relationship (2),<sup>4</sup> so solving primal and dual problems leads to the same decision function.

Dual problems come with another nice property that each variable  $\alpha_i$  corresponds to a training instance  $(y_i, \mathbf{x}_i)$ . In contrast, for primal problems, each variable  $w_i$  corresponds to a feature. Optimization methods which update some variables at a time often need to access the corresponding instances (if solving dual) or the corresponding features (if solving primal). In practical applications, instance-wise data storage is more common than feature-wise storage. Therefore, a dual-based algorithm can directly work on the input data without any transformation.

<sup>3</sup>Because the bias term  $b$  is not considered, therefore, different from the dual problem considered in SVM literature, an inequality constraint  $\sum y_i \alpha_i = 0$  is absent from (16).

<sup>4</sup>However, we do not necessarily need the dual problem to get (2). For example, the reduced SVM [19] directly assumes that  $w$  is the linear combination of a subset of data.

Unfortunately, the dual form may not be always easier to solve. For example, the dual form of L1-regularized problems involves general linear constraints rather than bound constraints in (16), so solving primal may be easier.

- **Using low-order or high-order information** Low-order methods, such as gradient or sub-gradient methods, have been widely considered in large-scale training. They characterize low-cost update, low-memory requirement, and slow convergence. In classification tasks, slow convergence may not be a serious concern because a loose solution of (7) may already give similar testing performances to that by an accurate solution. High-order methods such as Newton methods often require the smoothness of the optimization problems. Further, the cost per step is more expensive; sometimes a linear system must be solved. However, their convergence rate is superior. These high-order methods are useful for applications needing an accurate solution of problem (7). Some (e.g., [20]) have tried a hybrid setting by using low-order methods in the beginning and switching to higher-order methods in the end.
- **Cost of different types of operations** In a real-world computer, not all types of operations cost equally. For example, exponential and logarithmic operations are much more expensive than multiplication and division. For training large-scale LR, because exp/log operations are required, the cost of this type of operations may accumulate faster than that of other types. An optimization method which can avoid intensive exp/log evaluations is potentially efficient; see more discussion in, for example, [12], [21], [22].
- **Parallelization** Most existing training algorithms are inherently sequential, but a parallel algorithm can make good use of the computational power in a multi-core machine or a distributed system. However, the communication cost between different cores or nodes may become a new bottleneck. See more discussion in Section VII.

Earlier developments of optimization methods for linear classification tend to focus on data with few features. By taking this property, they are able to easily train millions of instances [23]. However, these algorithms may not be suitable for sparse data with both large numbers of instances and features, for which we show in Section II that linear classifiers often give competitive accuracy with nonlinear classifiers. Many recent studies have proposed algorithms for such data. We list some of them (and their software name if any) according to regularization and loss functions used.

- **L2-regularized L1-loss SVM:** Available approaches include, for example, cutting plane methods for the primal form (SVM<sup>perf</sup> [4], OCAS [24], and BMRM [25]), a stochastic (sub-)gradient descent method for the primal form (Pegasos [5] and SGD [26]), and a coordinate descent method for the dual form (LIBLINEAR [6]).
- **L2-regularized L2-loss SVM:** Existing methods for the primal form include a coordinate descent method [21], a Newton method [27], and a trust region Newton method (LIBLINEAR [28]). For the dual problem, a coordinate descent method is in the software LIBLINEAR [6].
- **L2-regularized LR:** Most unconstrained optimization methods can be applied to solve the primal problem. An early comparison on small-scale data is [29]. Existing studies for large sparse data include iterative scaling methods [12], [30], [31], a truncated Newton method [32], and a trust region Newton method (LIBLINEAR [28]). Few works solve the dual problem. One example is a coordinate descent method (LIBLINEAR [33]).
- **L1-regularized L1-loss SVM:** It seems no studies have applied L1-regularized L1-loss SVM on large sparse data although some early works for data with either few features or few instances are available [34]–[36].
- **L1-regularized L2-loss SVM:** Some proposed methods include a coordinate descent method (LIBLINEAR [13]) and a Newton-type method [22].
- **L1-regularized LR:** Most methods solve the primal form, for example, an interior-point method (l1\_logreg [37]), (block) coordinate descent methods (BBR [38] and CGD [39]), a quasi-Newton method (OWL-QN [40]), Newton-type methods (GLMNET [41] and LIBLINEAR [22]), and a Nesterov’s method (SLEP [42]). Recently, an augmented Lagrangian method (DAL [43]) is proposed for solving the dual problem. Comparisons of methods for L1-regularized LR include [13], [44].

In the rest of this section, we show details of some optimization algorithms. We select them not only because they are popular but also because many design issues discussed earlier can be covered.

### B. Example: A Sub-gradient Method (Pegasos with Deterministic Settings)

Shalev-Shwartz et al. [5] proposed a method **Pegasos** for solving the primal form of L2-regularized L1-loss SVM. It can be used for batch and online learning. Here we discuss only the deterministic setting and leave the

---

**ALGORITHM 1: Pegasos for L2-regularized L1-loss SVM (deterministic setting for batch learning) [5]**


---

1. Given  $\mathbf{w}$  such that  $\|\mathbf{w}\|_2 \leq \sqrt{Cl}$ .
  2. For  $k = 1, 2, 3, \dots$ 
    - (a) Let  $B = \{(y_i, \mathbf{x}_i)\}_{i=1}^l$ .
    - (b) Compute the learning rate  $\eta = (Cl)/k$ .
    - (c) Compute  $\nabla^S f(\mathbf{w}; B)$  by (17).
    - (d)  $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla^S f(\mathbf{w}; B)$ .
    - (e) Project  $\mathbf{w}$  by (19) to ensure  $\|\mathbf{w}\|_2 \leq \sqrt{Cl}$ .
- 

stochastic setting in Section VII-A.

Given a training subset  $B$ , at each iteration, Pegasos approximately solves the following problem.

$$\min_{\mathbf{w}} f(\mathbf{w}; B) \equiv \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i \in B} \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i).$$

Here, for the deterministic setting,  $B$  is the whole training set. Because L1 loss is not differentiable, Pegasos takes the following sub-gradient direction of  $f(\mathbf{w}; B)$

$$\nabla^S f(\mathbf{w}; B) \equiv \mathbf{w} - C \sum_{i \in B^+} y_i \mathbf{x}_i, \quad (17)$$

where  $B^+ \equiv \{i \mid i \in B, 1 - y_i \mathbf{w}^T \mathbf{x}_i > 0\}$ , and updates  $\mathbf{w}$  by

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla^S f(\mathbf{w}; B), \quad (18)$$

where  $\eta = (Cl)/k$  is the learning rate and  $k$  is the iteration index. Different from earlier sub-gradient descent methods, after the update by (18), Pegasos further projects  $\mathbf{w}$  onto the ball set  $\{\mathbf{w} \mid \|\mathbf{w}\|_2 \leq \sqrt{Cl}\}$ .<sup>5</sup> That is,

$$\mathbf{w} \leftarrow \min(1, \frac{\sqrt{Cl}}{\|\mathbf{w}\|_2}) \mathbf{w}. \quad (19)$$

We show the overall procedure of Pegasos in Algorithm 1.

For convergence, it is proved that in  $O(1/\epsilon)$  iterations, Pegasos achieves an average  $\epsilon$ -accurate solution. That is,

$$f\left(\frac{\sum_{k=1}^T \mathbf{w}^k}{T}\right) - f(\mathbf{w}^*) \leq \epsilon,$$

where  $\mathbf{w}^k$  is the  $k$ th iterate and  $\mathbf{w}^*$  is the optimal solution.

Pegasos has been applied in many studies. One implementation issue is that information obtained in the algorithm cannot be directly used for designing a suitable stopping condition.

### C. Example: Trust Region Newton Method (TRON)

Trust region Newton method (TRON) is an effective approach for unconstrained and bound-constrained optimization. In [28], it applies the setting in [45] to solve (7) with L2 regularization and differentiable losses.

At each iteration, given an iterate  $\mathbf{w}$ , a trust region interval  $\Delta$ , and a quadratic model

$$q(\mathbf{d}) \equiv \nabla f(\mathbf{w})^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \nabla^2 f(\mathbf{w}) \mathbf{d} \quad (20)$$

as an approximation of  $f(\mathbf{w} + \mathbf{d}) - f(\mathbf{w})$ , TRON finds a truncated Newton step confined in the trust region by approximately solving the following sub-problem.

$$\min_{\mathbf{d}} q(\mathbf{d}) \quad \text{subject to} \quad \|\mathbf{d}\|_2 \leq \Delta. \quad (21)$$

<sup>5</sup> The optimal solution of  $f(\mathbf{w})$  is proven to be in the ball set  $\{\mathbf{w} \mid \|\mathbf{w}\|_2 \leq \sqrt{Cl}\}$ ; see Theorem 1 in [5].

---

**ALGORITHM 2:** TRON for L2-regularized LR and L2-loss SVM [28]

---

1. Given  $\mathbf{w}$ ,  $\Delta$ , and  $\sigma_0$ .
  2. For  $k = 1, 2, 3, \dots$ 
    - (a) Find an approximate solution  $\mathbf{d}$  of (21) by the conjugate gradient method.
    - (b) Check the ratio  $\sigma$  in (22).
    - (c) If  $\sigma > \sigma_0$   
 $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{d}$ .
    - (d) Adjust  $\Delta$  according to  $\sigma$ .
- 

Then, by checking the ratio

$$\sigma \equiv \frac{f(\mathbf{w} + \mathbf{d}) - f(\mathbf{w})}{q(\mathbf{d})} \quad (22)$$

of actual function reduction to estimated function reduction, TRON decides if  $\mathbf{w}$  should be updated and then adjusts  $\Delta$ . A large enough  $\sigma$  indicates that the quadratic model  $q(\mathbf{d})$  is close to  $f(\mathbf{w} + \mathbf{d}) - f(\mathbf{w})$ , so TRON updates  $\mathbf{w}$  to be  $\mathbf{w} + \mathbf{d}$  and slightly enlarges the trust region interval  $\Delta$  for the next iteration. Otherwise, the current iterate  $\mathbf{w}$  is unchanged and the trust region interval  $\Delta$  shrinks by multiplying a factor less than one. The overall procedure of TRON is presented in Algorithm 2. If the loss function is not twice differentiable (e.g., L2 loss), we can use generalized Hessian [14] as  $\nabla^2 f(\mathbf{w})$  in (20).

Some difficulties of applying Newton methods to linear classification include that  $\nabla^2 f(\mathbf{w})$  may be a huge  $n$  by  $n$  matrix and solving (21) is expensive. Fortunately,  $\nabla^2 f(\mathbf{w})$  of linear classification problems takes the following special form

$$\nabla^2 f(\mathbf{w}) = \mathcal{I} + CX^T D_w X,$$

where  $\mathcal{I}$  is an identity matrix,  $X \equiv [\mathbf{x}_1, \dots, \mathbf{x}_l]^T$ , and  $D_w$  is a diagonal matrix. In [28], a conjugate gradient method is applied to solve (21), where the main operation is the product between  $\nabla^2 f(\mathbf{w})$  and a vector  $\mathbf{v}$ . By

$$\nabla^2 f(\mathbf{w})\mathbf{v} = \mathbf{v} + C(X^T(D_w(X\mathbf{v}))), \quad (23)$$

the Hessian matrix  $\nabla^2 f(\mathbf{w})$  need not be stored.

Because of using high-order information (Newton directions), TRON gives fast quadratic local convergence. It has been extended to solve L1-regularized LR and L2-loss SVM in [13] by reformulating (7) to a bound-constrained optimization problem in (15).

#### D. Example: Solving Dual SVM by Coordinate Descent Methods (Dual-CD)

Hsieh et al. [6] proposed a coordinate descent method for the dual L2-regularized linear SVM in (16). We call this algorithm Dual-CD. Here, we focus on L1-loss SVM, although the same method has been applied to L2-loss SVM in [6].

A coordinate descent method sequentially selects one variable for update and fixes others. To update the  $i$ th variable, the following one-variable problem is solved.

$$\begin{aligned} \min_d \quad & f^D(\boldsymbol{\alpha} + d\mathbf{e}_i) - f^D(\boldsymbol{\alpha}) \\ \text{subject to} \quad & 0 \leq \alpha_i + d \leq C, \end{aligned}$$

where  $f(\boldsymbol{\alpha})$  is defined in (16),  $\mathbf{e}_i = \underbrace{[0, \dots, 0]}_{i-1}, 1, 0, \dots, 0]^T$ , and

$$f^D(\boldsymbol{\alpha} + d\mathbf{e}_i) - f^D(\boldsymbol{\alpha}) = \frac{1}{2}Q_{ii}d^2 + \nabla_i f^D(\boldsymbol{\alpha})d.$$

This simple quadratic function can be easily minimized. After considering the constraint, a simple update rule for  $\alpha_i$  is

$$\alpha_i \leftarrow \min(\max(\alpha_i - \frac{\nabla_i f^D(\boldsymbol{\alpha})}{Q_{ii}}, 0), C). \quad (24)$$



---

**ALGORITHM 3:** A coordinate descent method for L2-regularized L1-loss SVM [6]

---

1. Given  $\alpha$  and the corresponding  $\mathbf{u} = \sum_{i=1}^l y_i \alpha_i \mathbf{x}_i$ .
  2. Compute  $Q_{ii}, \forall i = 1, \dots, l$ .
  3. For  $k = 1, 2, 3, \dots$ 
    - For  $i = 1, \dots, l$ 
      - (a) Compute  $G = y_i \mathbf{u}^T \mathbf{x}_i - 1$  in (27).
      - (b)  $\bar{\alpha}_i \leftarrow \alpha_i$ .
      - (c)  $\alpha_i \leftarrow \min(\max(\alpha_i - G/Q_{ii}, 0), C)$ .
      - (d)  $\mathbf{u} \leftarrow \mathbf{u} + y_i(\alpha_i - \bar{\alpha}_i)\mathbf{x}_i$ .
- 

From (24),  $Q_{ii}$  and  $\nabla_i f^D(\alpha)$  are our needs. The diagonal entries of  $Q$ ,  $Q_{ii}, \forall i$ , are computed only once initially, but

$$\nabla_i f^D(\alpha) = (Q\alpha)_i - 1 = \sum_{t=1}^l (y_i y_t \mathbf{x}_i^T \mathbf{x}_t) \alpha_t - 1 \quad (25)$$

requires  $O(nl)$  cost for  $l$  inner products  $\mathbf{x}_i^T \mathbf{x}_t, \forall t = 1, \dots, l$ . To make coordinate descent methods viable for large linear classification, a crucial step is to maintain

$$\mathbf{u} \equiv \sum_{t=1}^l y_t \alpha_t \mathbf{x}_t, \quad (26)$$

so that (25) becomes

$$\nabla_i f^D(\alpha) = (Q\alpha)_i - 1 = y_i \mathbf{u}^T \mathbf{x}_i - 1. \quad (27)$$

If  $\mathbf{u}$  is available through the training process, then the cost  $O(nl)$  in (25) is significantly reduced to  $O(n)$ . The remaining task is to maintain  $\mathbf{u}$ . Following (26), if  $\bar{\alpha}_i$  and  $\alpha_i$  are values before and after the update (24), respectively, then we can easily maintain  $\mathbf{u}$  by the following  $O(n)$  operation.

$$\mathbf{u} \leftarrow \mathbf{u} + y_i(\alpha_i - \bar{\alpha}_i)\mathbf{x}_i. \quad (28)$$

Therefore, the total cost for updating an  $\alpha_i$  is  $O(n)$ . The overall procedure of the coordinate descent method is in Algorithm 3.

The vector  $\mathbf{u}$  defined in (26) is in the same form as  $\mathbf{w}$  in (2). In fact, as  $\alpha$  approaches a dual optimal solution,  $\mathbf{u}$  will converge to the primal optimal  $\mathbf{w}$  following the primal-dual relationship.

The linear convergence of Algorithm 3 is established in [6] using techniques in [46]. They propose two implementation tricks to speed up the convergence. First, instead of a sequential update, they repeatedly permute  $\{1, \dots, l\}$  to decide the order. Second, similar to the shrinking technique used in training nonlinear SVM [47], they identify some bounded variables which may already be optimal and remove them during the optimization procedure. Experiments in [6] show that for large sparse data, Algorithm 3 is much faster than TRON in the early stage. However, it is less competitive if the parameter  $C$  is large.

Algorithm 3 is very related to popular decomposition methods used in training nonlinear SVM (e.g., [8], [47]). These decomposition methods also update very few variables at each step, but use more sophisticated schemes for selecting variables. The main difference is that for linear SVM, we can define  $\mathbf{u}$  in (26) because  $\mathbf{x}_i, \forall i$  are available. For nonlinear SVM,  $\nabla_i f^D(\mathbf{w})$  in (25) needs  $O(nl)$  cost for calculating  $l$  kernel elements. This difference between  $O(n)$  and  $O(nl)$  is similar to that in the testing phase discussed in Section II.

### E. Example: Solving L1-regularized Problems by Combining Newton and Coordinate Descent Methods (newGLMNET)

GLMNET proposed by Friedman et al. [41] is a Newton method for L1-regularized minimization. An improved version newGLMNET [22] is proposed for large-scale training.

---

**ALGORITHM 4: newGLMNET for L1-regularized minimization [22]**


---

1. Given  $\mathbf{w}$ . Given  $0 < \beta, \sigma < 1$ .
  2. For  $k = 1, 2, 3, \dots$ 
    - (a) Find an approximate solution  $\mathbf{d}$  of (29) by a coordinate descent method.
    - (b) Find  $\lambda = \max\{1, \beta, \beta^2, \dots\}$  such that (31) holds.
    - (c)  $\mathbf{w} \leftarrow \mathbf{w} + \lambda \mathbf{d}$ .
- 

Because the 1-norm term is not differentiable, we represent  $f(\mathbf{w})$  as the sum of two terms  $\|\mathbf{w}\|_1 + L(\mathbf{w})$ , where

$$L(\mathbf{w}) \equiv C \sum_{i=1}^l \xi(\mathbf{w}; \mathbf{x}_i, y_i).$$

At each iteration, newGLMNET considers the second-order approximation of  $L(\mathbf{w})$  and solves the following problem.

$$\min_{\mathbf{d}} q(\mathbf{d}) \equiv \|\mathbf{w} + \mathbf{d}\|_1 - \|\mathbf{w}\|_1 + \nabla L(\mathbf{w})^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T H \mathbf{d}, \quad (29)$$

where  $H \equiv \nabla^2 L(\mathbf{w}) + \nu \mathcal{I}$  and  $\nu$  is a small number to ensure  $H$  to be positive definite.

Although (29) is similar to (21), its optimization is more difficult because of the 1-norm term. Thus, newGLMNET further breaks (29) to sub-problems by a coordinate descent procedure. In a setting similar to the method in Section IV-D, each time a one-variable function is minimized.

$$\begin{aligned} & q(\mathbf{d} + z \mathbf{e}_j) - q(\mathbf{d}) \\ &= |w_j + d_j + z| - |w_j + d_j| + G_j z + \frac{1}{2} H_{jj} z^2, \end{aligned} \quad (30)$$

where  $G \equiv \nabla L(\mathbf{w}) + H \mathbf{d}$ . This one-variable function (30) has a simple closed-form minimizer (see [48], [49], and Appendix B of [13]).

$$z = \begin{cases} -\frac{G_j+1}{H_{jj}} & \text{if } G_j + 1 \leq H_{jj}(w_j + d_j), \\ -\frac{G_j-1}{H_{jj}} & \text{if } G_j - 1 \geq H_{jj}(w_j + d_j), \\ -(w_j + d_j) & \text{otherwise.} \end{cases}$$

At each iteration of newGLMNET, the coordinate descent method does not solve problem (29) exactly. Instead, newGLMNET designs an adaptive stopping condition so that initially problem (29) is solved loosely and in the final iterations, (29) is more accurately solved.

After an approximate solution  $\mathbf{d}$  of (29) is obtained, we need a line search procedure to ensure the sufficient function decrease. It finds  $\lambda \in (0, 1]$  such that

$$f(\mathbf{w} + \lambda \mathbf{d}) - f(\mathbf{w}) \leq \sigma \lambda (\|\mathbf{w} + \mathbf{d}\|_1 - \|\mathbf{w}\|_1 + \nabla L(\mathbf{w})^T \mathbf{d}), \quad (31)$$

where  $\sigma \in (0, 1)$ . The overall procedure of newGLMNET is in Algorithm 4.

Due to the adaptive setting, in the beginning, newGLMNET behaves like a coordinate descent method, which is able to quickly obtain an approximate  $\mathbf{w}$ ; however, in the final stage, the iterate  $\mathbf{w}$  converges quickly because a Newton step is taken. Recall in Section IV-A, we mentioned that exp/log operations are more expensive than basic operations such as multiplication/division. Because (30) does not involve any exp/log operation, we successfully achieve that time spent on exp/log operations is only a small portion of the whole procedure. In addition, newGLMNET is an example of accessing data feature-wisely; see details in [22] about how  $G_j$  in (30) is updated.

#### F. A Comparison of the Four Examples

The four methods discussed in Sections IV-B–IV-E differ in various aspects. By considering design issues mentioned in Section IV-A, we compare these methods in Table II. We point out that three methods are primal-based, but one is dual-based. Next, both Pegasos and Dual-CD use only low-order information (sub-gradient and gradient), but TRON and newGLMNET employ high-order information by Newton directions. Also, we check how

TABLE II  
A COMPARISON OF THE FOUR METHODS IN SECTIONS IV-B–IV-E.

	Pegasos	TRON	Dual-CD	newGLMNET
primal-/dual-based	primal	primal	dual	primal
low/high order	low	high	low	high
data access	instance-wisely	both	instance-wisely	feature-wisely

data instances are accessed. Clearly, Pegasos and Dual-CD instance-wisely access data, but we have mentioned in Section IV-E that newGLMNET must employ a feature-wisely setting. Interestingly, TRON can use both because in Eq. (23), matrix-vector products can be conducted by accessing data instance-wisely or feature-wisely.

We analyze the complexity of the four methods by showing the cost at the  $k$ th iteration:

**Pegasos:**  $O(|B^+|n)$

**TRON:**  $\#CG\_iter \times O(ln)$

**Dual-CD:**  $O(ln)$

**newGLMNET:**  $\#CD\_iter \times O(ln)$ .

The cost of Pegasos and TRON easily follows from (17) and (23), respectively. For Dual-CD, both (27) and (28) cost  $O(n)$ , so one iteration of going through all variables is  $O(nl)$ . For newGLMNET, please see details in [22]. We can clearly see that each iteration of Pegasos and Dual-CD is cheaper because of using low-order information. However, they need more iterations than high-order methods in order to accurately solve the optimization problem.

## V. MULTI-CLASS LINEAR CLASSIFICATION

Most classification methods are originally proposed to solve a two-class problem; however, extensions of these methods to multi-class classification have been studied. For non-linear SVM, some works (e.g., [50], [51]) have comprehensively compared different multi-class solutions. In contrast, few studies have focused on multi-class linear classification. This section introduces and compares some commonly used methods.

### A. Solving Several Binary Problems

Multi-class classification can be decomposed to several binary classification problems. One-against-rest and one-against-one methods are two of the most common decomposition approaches. Studies that broadly discussed various approaches of decomposition include, for example, [52], [53].

- **One-against-rest method** If there are  $k$  classes in the training data, the one-against-rest method [54] constructs  $k$  binary classification models. To obtain the  $m$ th model, instances from the  $m$ th class of the training set are treated as positive, and all other instances are negative. Then the weight vector  $\mathbf{w}_m$  for the  $m$ th model can be generated by any linear classifier.

After obtaining all  $k$  models, we say an instance  $\mathbf{x}$  is in the  $m$ th class if the decision value (1) of the  $m$ th model is the largest, i.e.,

$$\text{class of } \mathbf{x} \equiv \arg \max_{m=1,\dots,k} \mathbf{w}_m^T \mathbf{x}. \quad (32)$$

The cost for testing an instance is  $O(nk)$ .

- **One-against-one method** One-against-one method [55] solves  $k(k-1)/2$  binary problems. Each binary classifier constructs a model with data from one class as positive and another class as negative. Since there are  $k(k-1)/2$  combination of two classes,  $k(k-1)/2$  weight vectors are constructed:  $\mathbf{w}_{1,2}, \mathbf{w}_{1,3}, \dots, \mathbf{w}_{1,k}, \mathbf{w}_{2,3}, \dots, \mathbf{w}_{(k-1),k}$ . There are different methods for testing. One approach is by voting [56]. For a testing instance  $\mathbf{x}$ , if model  $(i, j)$  predicts  $\mathbf{x}$  as in the  $i$ th class, then a counter for the  $i$ th class is added by one; otherwise, the counter for the  $j$ th class is added. Then we say  $\mathbf{x}$  is in the  $i$ th class if the  $i$ th counter has the largest value. Other prediction methods are similar though they differ in how to use the  $k(k-1)/2$  decision values; see some examples in [52], [53]. For linear classifiers, one-against-one method is shown to give better testing accuracy than one-against-rest [57]. However, it requires  $O(k^2n)$  spaces for storing models and  $O(k^2n)$  cost for testing an instance; both are more expensive than the one-against-rest method. Interestingly, for nonlinear classifiers via kernels, one-against-one method does not have such disadvantages [50].

DAGSVM [58] is the same as one-against-one but it attempts to reduce the testing cost. Starting with a candidate set of all classes, this method sequentially selects a pair of classes for prediction and removes one of the two. That is, if a binary classifier of class  $i$  and  $j$  predicts  $i$ , then  $j$  is removed from the candidate set. Alternatively, a prediction of class  $j$  will cause  $i$  to be removed. Finally, the only remained class is the predicted result. For any pair  $(i, j)$  considered, the true class may be neither  $i$  nor  $j$ . However, it does not matter which one is removed because all we need is that if the true class is involved in a binary prediction, it is the winner. Because classes are sequentially removed, only  $k - 1$  models are used. The testing time complexity of DAGSVM is thus  $O(nk)$ .

### B. Considering All Data at Once

In contrast to using many binary models, some have proposed solving a single optimization problem for multi-class classification [59]–[61]. Here we discuss details of Crammer and Singer’s approach [60]. Assume class labels are  $1, \dots, k$ . They consider the following optimization problem:

$$\min_{\mathbf{w}_1, \dots, \mathbf{w}_k} \frac{1}{2} \sum_{m=1}^k \|\mathbf{w}_m\|_2^2 + C \sum_{i=1}^l \xi_{\text{CS}}(\{\mathbf{w}_m\}_{m=1}^k; \mathbf{x}_i, y_i), \quad (33)$$

where

$$\xi_{\text{CS}}(\{\mathbf{w}_m\}_{m=1}^k; \mathbf{x}, y) \equiv \max_{m \neq y} \max(0, 1 - (\mathbf{w}_y - \mathbf{w}_m)^T \mathbf{x}). \quad (34)$$

The setting is like to combine all binary models of the one-against-rest method. There are  $k$  weight vectors  $\mathbf{w}_1, \dots, \mathbf{w}_k$  for  $k$  classes. In the loss function (34), for each  $m$ ,  $\max(0, 1 - (\mathbf{w}_{y_i} - \mathbf{w}_m)^T \mathbf{x}_i)$  is similar to the L1 loss in (8) for binary classification. Overall, we hope that the decision value of  $\mathbf{x}_i$  by the model  $\mathbf{w}_{y_i}$  is at least one larger than the values by other models. For testing, the decision function is also (32).

Early works of this method focus on the nonlinear (i.e., kernel) case [50], [60], [62]. A study for linear classification is in [63], which applies a coordinate descent method to solve the dual problem of (33). The idea is similar to the method in Section IV-D; however, at each step, a larger sub-problem of  $k$  variables is solved. A nice property of this  $k$ -variable sub-problem is that it has a closed-form solution. Experiments in [63] show that solving (33) gives slightly better accuracy than one-against-rest, but the training time is competitive. This result is different from the nonlinear case, where the longer training time than one-against-rest and one-against-one has made the approach of solving one single optimization problem less practical [50]. A careful implementation of the approach in [63] is given in [7, Appendix E].

### C. Maximum Entropy

Maximum Entropy (ME) [64] is a generalization of logistic regression for multi-class problems<sup>6</sup> and a special case of conditional random fields [65] (see Section VIII-A). It is widely applied by NLP applications. We still assume class labels  $1, \dots, k$  for an easy comparison to (33) in our subsequent discussion. ME models the following conditional probability function of label  $y$  given data  $\mathbf{x}$ .

$$P(y|\mathbf{x}) \equiv \frac{\exp(\mathbf{w}_y^T \mathbf{x})}{\sum_{m=1}^k \exp(\mathbf{w}_m^T \mathbf{x})}, \quad (35)$$

where  $\mathbf{w}_m, \forall m$  are weight vectors like those in (32) and (33). This model is also called multinomial logistic regression.

ME minimizes the following regularized negative log-likelihood.

$$\min_{\mathbf{w}_1, \dots, \mathbf{w}_k} \frac{1}{2} \sum_{m=1}^k \|\mathbf{w}_m\|_2^2 + C \sum_{i=1}^l \xi_{\text{ME}}(\{\mathbf{w}_m\}_{m=1}^k; \mathbf{x}_i, y_i), \quad (36)$$

where

$$\xi_{\text{ME}}(\{\mathbf{w}_m\}_{m=1}^k; \mathbf{x}, y) \equiv -\log P(y|\mathbf{x}).$$

<sup>6</sup>Details of the connection between logistic regression and maximum entropy can be found in, for example, [12, Section 5.2].

TABLE III

COMPARISON OF METHODS FOR MULTI-CLASS LINEAR CLASSIFICATION IN STORAGE (MODEL SIZE) AND TESTING TIME.  $n$  IS THE NUMBER OF FEATURES AND  $k$  IS THE NUMBER OF CLASSES.

Method	Storage	Testing
one-against-rest	$O(kn)$	$O(kn)$
one-against-one	$O(k^2n)$	$O(k^2n)$
DAGSVM	$O(k^2n)$	$O(kn)$
Crammer and Singer's	$O(kn)$	$O(kn)$
maximum entropy	$O(kn)$	$O(kn)$

Clearly, (36) is similar to (33) and  $\xi_{\text{ME}}(\cdot)$  can be considered as a loss function. If  $\mathbf{w}_{y_i}^T \mathbf{x}_i \gg \mathbf{w}_m^T \mathbf{x}_i, \forall m \neq y_i$ , then  $\xi_{\text{ME}}(\{\mathbf{w}_m\}_{m=1}^k; \mathbf{x}_i, y_i)$  is close to zero (i.e., no loss). On the other hand, if  $\mathbf{w}_{y_i}^T \mathbf{x}_i$  is smaller than other  $\mathbf{w}_m^T \mathbf{x}_i, m \neq y_i$ , then  $P(y_i|\mathbf{x}_i) \ll 1$  and the loss is large. For prediction, the decision function is also (32).

NLP applications often consider a more general ME model by using a function  $\mathbf{f}(\mathbf{x}, y)$  to generate the feature vector.

$$P(y|\mathbf{x}) \equiv \frac{\exp(\mathbf{w}^T \mathbf{f}(\mathbf{x}, y))}{\sum_{y'} \exp(\mathbf{w}^T \mathbf{f}(\mathbf{x}, y'))}. \quad (37)$$

Eq. (35) is a special case of (37) by

$$\mathbf{f}(\mathbf{x}_i, y) = \left[ \begin{array}{c} 0 \\ \vdots \\ 0 \\ \mathbf{x}_i \\ 0 \\ \vdots \\ 0 \end{array} \right] \left. \vphantom{\begin{array}{c} 0 \\ \vdots \\ 0 \\ \mathbf{x}_i \\ 0 \\ \vdots \\ 0 \end{array}} \right\} y - 1 \in \mathbf{R}^{nk} \text{ and } \mathbf{w} = \left[ \begin{array}{c} \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_k \end{array} \right]. \quad (38)$$

Many studies have investigated optimization methods for L2-regularized ME. For example, Malouf [66] compares iterative scaling methods [67], gradient descent, nonlinear conjugate gradient, and L-BFGS (quasi Newton) method [68] to solve (36). Experiments show that quasi Newton performs better. In [12], a framework is proposed to explain variants of iterative scaling methods [30], [67], [69] and make a connection to coordinate descent methods. For L1-regularized ME, [40] proposes an extension of L-BFGS.

Recently, instead of solving the primal problem (36), some works solve the dual problem. A detailed derivation of the dual ME is in [33, Appendix A.7]. Memisevic [70] proposed a two-level decomposition method. Similar to the coordinate descent method [63] for (33) in Section V-B, in [70], a sub-problem of  $k$  variables is considered at a time. However, the sub-problem does not have a closed-form solution, so a second-level coordinate descent method is applied. Collin et al. [71] proposed an exponential gradient method to solve ME dual. They also decompose the problem into  $k$ -variable sub-problems, but only approximately solve each sub-problem. The work in [33] follows [70] to apply a two-level coordinate descent method, but uses a different method in the second level to decide variables for update.

#### D. Comparison

We summarize storage (model size) and testing time of each method in Table III. Clearly, one-against-one and DAGSVM are less practical because of the much higher storage, although the comparison in [57] indicates that one-against-one gives slightly better testing accuracy. Note that the situation is very different for the kernel case [50], where one-against-one and DAGSVM are very useful methods.

## VI. LINEAR-CLASSIFICATION TECHNIQUES FOR NONLINEAR CLASSIFICATION

Many recent developments of linear classification can be extended to handle non-standard scenarios. Interestingly, most of them are related to training nonlinear classifiers.

TABLE IV

RESULTS OF TRAINING/TESTING DEGREE-2 POLYNOMIAL MAPPINGS BY THE COORDINATE DESCENT METHOD IN SECTION IV-D. THE DEGREE-2 POLYNOMIAL MAPPING IS DYNAMICALLY COMPUTED DURING TRAINING, INSTEAD OF EXPANDED BEFOREHAND. THE LAST COLUMN SHOWS THE ACCURACY DIFFERENCE BETWEEN DEGREE-2 POLYNOMIAL MAPPINGS AND RBF SVM.

Data set	Degree-2 polynomial		Testing accuracy	Accuracy difference to RBF kernel
	Training Time (s)	Testing Time (s)		
cod-RNA	1.27	0.12	96.35	-0.3210
ijcnn1	10.38	0.08	97.84	-0.8517
covtype	5,166.30	0.07	80.21	-15.8999

### A. Training and Testing Explicit Data Mappings via Linear Classifiers

In some problems, training a linear classifier in the original feature space may not lead to competitive performances. For example, on `ijcnn1` in Table I, the testing accuracy (92.21%) of a linear classifier is inferior to 98.69% of a nonlinear one with the RBF kernel. However, the higher accuracy comes with longer training and testing time. Taking the advantage of linear classifiers' fast training, some studies have proposed using the explicit nonlinear data mappings. That is, we consider  $\phi(\mathbf{x}_i)$ ,  $i = 1, \dots, l$  as the new training set and employ a linear classifier. In some problems, this type of approaches may still enjoy fast training/testing, but achieve accuracy close to that of using highly nonlinear kernels.

Some early works, e.g., [72]–[74], have directly trained nonlinearly mapped data in their experiments. Chang et al. [9] analyze when this approach leads to faster training and testing. Assume that the coordinate descent method in Section IV-D is used for training linear/kernelized classifiers<sup>7</sup> and  $\phi(\mathbf{x}) \in \mathbf{R}^d$ . From Section IV-D, each coordinate descent step takes  $O(d)$  and  $O(nl)$  operations for linear and kernelized settings, respectively. Thus, if  $d \ll nl$ , the approach of training explicit mappings may be faster than using kernels. In [9], they particularly study degree-2 polynomial mappings such as (5). The dimensionality is  $d = O(n^2)$ , but for sparse data, the  $O(n^2)$  versus  $O(nl)$  comparison is changed to  $O(\bar{n}^2)$  versus  $O(\bar{n}l)$ , where  $\bar{n}$  is the average number of non-zero values per instance. For large sparse data sets,  $\bar{n} \ll l$ , so their approach can be very efficient. Table IV shows results of training/testing degree-2 polynomial mappings using three data sets in Table I with significant lower linear-SVM accuracy than RBF. We apply the same setting as [9, Section 4]. From Tables I and IV, we observed that training  $\phi(\mathbf{x}_i), \forall i$  by a linear classifier may give accuracy close to RBF kernel, but is faster in training/testing.

A general framework was proposed in [75] for various nonlinear mappings of data. They noticed that to perform the coordinate descent method in Section IV-D, one only needs that  $\mathbf{u}^T \phi(\mathbf{x})$  in (27) and  $\mathbf{u} \leftarrow \mathbf{u} + y(\alpha_i - \bar{\alpha}_i)\phi(\mathbf{x})$  in (28) can be performed. Thus, even if  $\phi(\mathbf{x})$  cannot be explicitly represented, as long as these two operations can be performed, Algorithm 3 is applicable.

Studies in [76], [77] designed linear classifiers to train explicit mappings of sequence data, where features correspond to subsequences. Using the relation between subsequences, they are able to design efficient training methods for very high dimensional mappings.

### B. Approximation of Kernel Methods via Linear Classification

Methods in Section VI-A train  $\phi(\mathbf{x}_i), \forall i$  explicitly, so they obtain the same model as a kernel method using  $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ . However, they have limitations when the dimensionality of  $\phi(\mathbf{x})$  is very high. To resolve the slow training/testing of kernel methods, approximation is sometimes unavoidable. Among the many available methods to approximate the kernel, some of them lead to training a linear classifier. Following [78], we categorize these methods to the following two types.

- **Kernel matrix approximation** This type of approaches finds a low-rank matrix  $\bar{\Phi} \in \mathbf{R}^{d \times l}$  with  $d \ll l$  such that  $\bar{\Phi}^T \bar{\Phi}$  can approximate the kernel matrix  $Q$ .

$$\bar{Q} = \bar{\Phi}^T \bar{\Phi} \approx Q. \quad (39)$$

<sup>7</sup>See the discussion in the end of Section IV-D about the connection between Algorithm 3 and the popular decomposition methods for nonlinear SVMs.

Assume  $\bar{\Phi} \equiv [\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_l]$ . If we replace  $Q$  in (16) with  $\bar{Q}$ , then (16) becomes the dual problem of training a linear SVM on the new set  $(y_i, \bar{\mathbf{x}}_i)$ ,  $i = 1, \dots, l$ . Thus, optimization methods discussed in Section IV can be directly applied. An advantage of this approach is that we do not need to know an explicit mapping function corresponding to a kernel of our interest (see the other type of approaches discussed below). However, this property causes a complicated testing procedure. That is, the approximation in (39) does not directly reveal how to adjust the decision function (3).

Early developments focused on finding a good approximation matrix  $\bar{\Phi}$ . Some examples include Nyström method [79], [80] and incomplete Cholesky factorization [81], [82]. Some works (e.g., [19]) consider approximations other than (39), but also lead to linear classification problems.

A recent study [78] addresses more on training and testing linear SVM after obtaining the low-rank approximation. In particular, details of the testing procedures can be found in Section 2.4 of [78]. Note that linear SVM problems obtained after kernel approximations are often dense and have more instances than features. Thus, training algorithms suitable for such problems may be different from those for sparse document data.

- **Feature mapping approximation** This type of approaches finds a mapping function  $\bar{\phi} : \mathbf{R}^n \rightarrow \mathbf{R}^d$  such that

$$\bar{\phi}(\mathbf{x})^T \bar{\phi}(\mathbf{t}) \approx K(\mathbf{x}, \mathbf{t}).$$

Then, linear classifiers can be applied to new data  $\bar{\phi}(\mathbf{x}_1), \dots, \bar{\phi}(\mathbf{x}_l)$ . The testing phase is straightforward because the mapping  $\bar{\phi}(\cdot)$  is available.

Many mappings have been proposed. Examples include random Fourier projection [83], random projections [84]–[87], polynomial approximation [88], and hashing [89]–[92]. They differ in various aspects, which are beyond the scope of this paper. An issue related to the subsequent linear classification is that some methods (e.g., [93]) generate dense  $\bar{\phi}(\mathbf{x})$  vectors, while others give sparse vectors (e.g., [85]). A recent study focusing on the linear classification after obtaining  $\bar{\phi}(\mathbf{x}_i)$ ,  $\forall i$  is in [94].

## VII. TRAINING LARGE DATA BEYOND THE MEMORY OR THE DISK CAPACITY

Recall that we described some binary linear classification algorithms in Section IV. Those algorithms can work well under the assumption that the training set is stored in the computer memory. However, as the training size goes beyond the memory capacity, traditional algorithms may become very slow because of frequent disk access. Indeed, even if the memory is enough, loading data to memory may take more time than subsequent computation [95]. Therefore, the design of algorithms for data larger than memory is very different from that of traditional algorithms.

If the data set is beyond the disk capacity of a single computer, then it must be stored distributively. Internet companies now routinely handle such large data sets in data centers. In such a situation, linear classification faces even more challenges because of expensive communication cost between different computing nodes. In some recent works [96], [97], parallel SVM on distributed environments has been studied but they investigated only kernel SVM. The communication overhead is less serious because of expensive kernel computation. For distributed linear classification, the research is still in its infancy. The current trend is to design algorithms so that computing nodes access data locally and the communication between nodes is minimized. The implementation is often conducted using distributed computing environments such as Hadoop [98]. In this section, we will discuss some ongoing research results.

Among the existing developments, some can be easily categorized as online methods. We describe them in Section VII-A. Batch methods are discussed in Section VII-B, while other approaches are in Section VII-C

### A. Online Methods

An online learning method receives a sequence of training samples and processes some instances at a time. Because training instances may be used only once, not only can online methods handle data larger than memory, they are also suitable for streaming data. An online optimization method can also be applied to a batch setting, where a fixed set of training instances is given; it iteratively updates the model  $\mathbf{w}$  using some instances at a time.

One popular online algorithm is the stochastic gradient descent method (SGD), which can be traced back to stochastic approximation method [99], [100]. Take the primal L2-regularized L1-loss SVM in (7) as an example.

At each step, a training instance  $\mathbf{x}_i$  is chosen and  $\mathbf{w}$  is updated by

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla^S \left( \frac{1}{2} \|\mathbf{w}\|_2^2 + C \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i) \right), \quad (40)$$

where  $\nabla^S$  is a sub-gradient operator and  $\eta$  is the learning rate. Specifically, (40) becomes the following update rule.

$$\begin{aligned} \text{If } 1 - y_i \mathbf{w}^T \mathbf{x}_i > 0, \quad \text{then} \\ \mathbf{w} \leftarrow (1 - \eta) \mathbf{w} + \eta C y_i \mathbf{x}_i. \end{aligned} \quad (41)$$

The learning rate  $\eta$  is gradually reduced along iterations.

It is well known that stochastic gradient descent methods have slow convergence. However, they are suitable for large data because of accessing only one instance at a time. Early studies which have applied SGD to linear classification include, for example, [101], [102]. For data with many features, recent studies [5], [26] show that SGD is effective. They allow more flexible settings such as using more than one training instance at a time. We briefly discuss the online setting of Pegasos [5]. In Algorithm 1, at each Step (a), a small random subset  $B$  is used instead of the full set. Similar convergence properties to that described in Section IV-B still hold but in expectation (see Theorem 2 in [5]).

Instead of solving the primal problem, we can design an online algorithm to solve the dual problem [6], [103]. For example, the coordinate descent method in Algorithm 3 can be easily extended to an online setting by replacing the sequential selection of variables with a random selection. Notice that the update rule (28) is similar to (41), but has the advantage of not needing to decide the learning rate  $\eta$ . This online setting falls into the general framework of randomized coordinate descent methods in [104], [105]. Using the proof in [104], the linear convergence in expectation is obtained in Appendix 7.5 of [6].

To improve the convergence of SGD, some [106], [107] have proposed using higher-order information. The rule in (40) is replaced by

$$\mathbf{w} \leftarrow \mathbf{w} - \eta H \nabla^S(\cdot), \quad (42)$$

where  $H$  is an approximation of the inverse Hessian  $\nabla^2 f(\mathbf{w})^{-1}$ . To save the cost at each update, practically  $H$  is a diagonal scaling matrix. Experiments [106], [107] show that using (42) is faster than (40).

The update rule in (40) assumes L2 regularization. While SGD is applicable for other regularization, it may not perform as well because of not taking special properties of the regularization term into consideration. For example, if L1 regularization is used, a standard SGD may face difficulties to generate a sparse  $\mathbf{w}$ . To address this problem, recently several approaches have been proposed [108]–[113]. The stochastic coordinate descent method in [109] has been extended to a parallel version [114].

Unfortunately, most existing studies of online algorithms conduct experiments by assuming enough memory and reporting the number of times to access data. To apply them in a real scenario without sufficient memory, many practical issues must be checked. Vowpal-Wabbit [115] is one of the very few implementations which can handle data larger than memory. Because the same data may be accessed several times and the disk reading time is expensive, at the first pass, Vowpal-Wabbit stores data to a compressed cache file. This is similar to the compression strategy in [95], which will be discussed in Section VII-B. Currently, Vowpal-Wabbit supports unregularized linear classification and regression. It is extended to solve L1-regularized problems in [108].

For data in a distributed environment, some studies have proposed parallel SGD. Each node only computes the sub-gradient corresponding to locally stored data instances. In [116], a delayed SGD is proposed. Instead of computing the sub-gradient of the current iterate  $\mathbf{w}^k$ , in delayed SGD, each node computes the sub-gradient of a previous iterator  $\mathbf{w}^{\tau(k)}$ , where  $\tau(k) \leq k$ . Delayed SGD is useful to reduce the synchronization delay because of communication overheads or uneven computational time at various nodes. Recent works [117], [118] show that delayed SGD is efficient when the number of nodes is large, and the delay is asymptotically negligible. Other online distributed learning studies include, for example, [119]–[121]. Unfortunately, so far few existing software packages are available. Recently, Vowpal-Wabbit (after version 6.0) has supported distributed online learning using the Hadoop [98] framework. We are aware that other Internet companies have constructed online linear classifiers on distributed environments, although details have not been fully available. One example is the system SETI at Google [122].



## B. Batch Methods

In some situations, we still would like to consider the whole training set and solve a corresponding optimization problem. While this task is very challenging, some (e.g., [95], [123]) have checked the situation that data are larger than memory but smaller than disk. Because of expensive disk I/O, they design algorithms by reading a continuous chunk of data at a time and minimizing the number of disk accesses. The method in [95] extends the coordinate descent method in Section IV-D for linear SVM. The major change is to update more variables at a time so that a block of data is used together. Specifically, in the beginning the training set is randomly partitioned to  $m$  files  $B_1, \dots, B_m$ . The available memory space needs to be able to accommodate one block of data and the working space of a training algorithm. To solve (16), sequentially one block of data  $B$  is read and the following function of  $\mathbf{d}$  is minimized under the condition  $0 \leq \alpha_i + d_i \leq C, \forall i \in B$  and  $d_i = 0, \forall i \notin B$ .

$$\begin{aligned} & f^D(\boldsymbol{\alpha} + \mathbf{d}) - f^D(\boldsymbol{\alpha}) \\ &= \frac{1}{2} \mathbf{d}_B^T Q_{BB} \mathbf{d}_B + \mathbf{d}_B^T (Q\boldsymbol{\alpha} - \mathbf{e})_B \\ &= \frac{1}{2} \mathbf{d}_B^T Q_{BB} \mathbf{d}_B + \sum_{i \in B} y_i d_i (\mathbf{u}^T \mathbf{x}_i) - \mathbf{d}_B^T \mathbf{e}_B, \end{aligned} \quad (43)$$

where  $Q_{BB}$  is a sub-matrix of  $Q$  and  $\mathbf{u}$  is defined in (26). By maintaining  $\mathbf{u}$  in a way similar to (28), Eq. (43) involves only data in the block  $B$ , which can be stored in memory. Eq. (43) can be minimized by any traditional algorithm. Experiments in [95] demonstrate that they can train data 20 times larger than the memory capacity. This method is extended in [124] to cache informative data points in the computer memory. That is, at each iteration, not only the selected block but also the cached points are used for updating corresponding variables. Their way to select informative points is inspired by the shrinking techniques used in training nonlinear SVM [8], [47].

For distributed batch learning, all existing parallel optimization methods [125] can possibly be applied. We discuss some existing approaches.

1. ADMM: Recently, [126] considers the ADMM (Alternating Direction Method of Multiplier) method [127] for distributed learning. Take SVM as an example and assume data points are partitioned to  $m$  distributively stored sets  $B_1, \dots, B_m$ . This method solves the following approximation of the original optimization problem.

$$\begin{aligned} \min_{\mathbf{w}_1, \dots, \mathbf{w}_m, \mathbf{z}} \quad & \frac{1}{2} \mathbf{z}^T \mathbf{z} + C \sum_{j=1}^m \sum_{i \in B_j} \xi_{L_1}(\mathbf{w}_j; \mathbf{x}_i, y_i) \\ & + \frac{\rho}{2} \sum_{j=1}^m \|\mathbf{w}_j - \mathbf{z}\|^2 \\ \text{subject to} \quad & \mathbf{w}_j - \mathbf{z} = \mathbf{0}, \forall j, \end{aligned}$$

where  $\rho$  is a pre-specific parameter. It then employs an optimization method of multipliers by alternatively minimizing the Lagrangian function over  $\mathbf{w}_1, \dots, \mathbf{w}_m$ , minimizing the Lagrangian over  $\mathbf{z}$ , and updating dual multipliers. The minimization of Lagrangian over  $\mathbf{w}_1, \dots, \mathbf{w}_m$  can be decomposed to  $m$  independent problems. Other steps do not involve data at all. Therefore, data points are locally accessed and the communication cost is kept minimum. Examples of using ADMM for distributed training include [128]. Some known problems of this approaches are first, the convergence rate is not very fast, and second, it is unclear how to choose the parameter  $\rho$ .

2. Parallel coordinate descent: In Section IV-D, the coordinate descent method is a sequential process that updates one variable at a time. To make it parallel, a possible approach is to select a subset of variables, construct independent sub-problems, and solve them in parallel. The main difficulty is that these independent sub-problems may be obtained by only a loose approximation of the function value, so the convergence can be slow. Existing works include [114], [129]–[131].
3. Quasi Newton method: Currently, a special quasi Newton method called LBFGS [68] has been commonly used in computational advertising [132], [133]. It has also been implemented in the software Vowpal-Wabbit [115] and Spark MLlib.<sup>8</sup>

<sup>8</sup>Spark is an cluster-computing platform. Its machine learning library MLlib is an integrated component of the framework.

4. Newton: the Newton method discussed in Section IV-C has been parallelized in, for example, [134], [135]. Their main idea is to conduct parallel matrix-vector products for the Hessian-vector product in (23).

### C. Other Approaches

We briefly discuss some other approaches which cannot be clearly categorized as batch or online methods.

The most straightforward method to handle large data is probably to randomly select a subset that can fit in memory. This approach works well if the data quality is good; however, sometimes using more data gives higher accuracy. To improve the performance of using only a subset, some have proposed techniques to include important data points into the subset. For example, the approach in [136] selects a subset by reading data from disk only once. For data in a distributed environment, sub-sampling can be a complicated operation. Moreover, a subset fitting the memory of one single computer may be too small to give good accuracy.

Bagging [137] is a popular classification method to split a learning task to several easier ones. It selects several random subsets, trains each of them, and ensembles (e.g., averaging) the results during testing. This method may be particularly useful for distributively stored data because we can directly consider data in each node as a subset. However, if data quality in each node is not good (e.g., all instances with the same class label), the model generated by each node may be poor. Thus, ensuring data quality of each subset is a concern. Some studies have applied the bagging approach on a distributed system [138], [139]. For example, in the application of web advertising, [138] trains a set of individual classifiers in a distributed way. Then, a final model is obtained by averaging the separate classifiers. In the linguistic applications, [140] extends the simple model average to the weighted average and achieves better performance. An advantage of the bagging-like approach is the easy implementation using distributed computing techniques such as MapReduce [141].<sup>9</sup>

## VIII. RELATED TOPICS

In this section, we discuss some other linear models. They are related to linear classification models discussed in earlier sections.

### A. Structured Learning

In the discussion so far, we assume that the label  $y_i$  is a single value. For binary classification, it is  $+1$  or  $-1$ , while for multi-class classification, it is one of the  $k$  class labels. However, in some applications, the label may be a more sophisticated object. For example, in part-of-speech (POS) tagging applications, the training instances are sentences and the labels are sequences of POS tags of words. If there are  $l$  sentences, we can write the training instances as  $(\mathbf{y}_i, \mathbf{x}_i) \in Y^{n_i} \times X^{n_i}, \forall i = 1, \dots, l$ , where  $\mathbf{x}_i$  is the  $i$ th sentence,  $\mathbf{y}_i$  is a sequence of tags,  $X$  is a set of unique words in the context,  $Y$  is a set of candidate tags for each word, and  $n_i$  is the number of words in the  $i$ th sentence. Note that we may not be able to split the problem to several independent ones by treating each value  $y_{ij}$  of  $\mathbf{y}_i$  as the label, because  $y_{ij}$  not only depends on the sentence  $\mathbf{x}_i$  but also other tags  $(y_{i1}, \dots, y_{i(j-1)}, y_{i(j+1)}, \dots, y_{in_i})$ . To handle these problems, we could use structured learning models like conditional random fields [65] and structured SVM [142], [143].

- **Conditional Random Fields** Conditional random fields (CRF) [65] is a linear structured model commonly used in NLP. Using notation mentioned above and a feature function  $\mathbf{f}(\mathbf{x}, \mathbf{y})$  like ME, CRF solves the following problem:

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^l \xi_{\text{CRF}}(\mathbf{w}; \mathbf{x}_i, \mathbf{y}_i), \quad (44)$$

where

$$\begin{aligned} \xi_{\text{CRF}}(\mathbf{w}; \mathbf{x}_i, \mathbf{y}_i) &\equiv -\log P(\mathbf{y}_i | \mathbf{x}_i), \text{ and} \\ P(\mathbf{y} | \mathbf{x}) &\equiv \frac{\exp(\mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y}'} \exp(\mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y}'))}. \end{aligned} \quad (45)$$

If elements in  $\mathbf{y}_i$  are independent to each other, then CRF reduces to ME.

<sup>9</sup>We mentioned earlier the Hadoop system, which includes a MapReduce implementation.

The optimization of (44) is challenging because in the probability model (45), the number of possible  $\mathbf{y}$ 's is exponentially many. An important property to make CRF practical is that the gradient of the objective function in (44) can be efficiently evaluated by dynamic programming [65]. Some available optimization methods include L-BFGS (quasi Newton) and conjugate gradient [144], stochastic gradient descent [145], stochastic quasi Newton [106], [146], and trust region Newton method [147]. It is shown in [147] that the Hessian-vector product (23) of the Newton method can also be evaluated by dynamic programming.

- **Structured SVM** Structured SVM solves the following optimization problem generalized form multi-class SVM in [59], [60].

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^l \xi_{\text{ss}}(\mathbf{w}; \mathbf{x}_i, \mathbf{y}_i), \quad (46)$$

where

$$\begin{aligned} & \xi_{\text{ss}}(\mathbf{w}; \mathbf{x}_i, \mathbf{y}_i) \\ & \equiv \max_{\mathbf{y} \neq \mathbf{y}_i} \left( \max \left( 0, \Delta(\mathbf{y}_i, \mathbf{y}) - \mathbf{w}^T (\mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{f}(\mathbf{x}_i, \mathbf{y})) \right) \right), \end{aligned}$$

and  $\Delta(\cdot)$  is a distance function with  $\Delta(\mathbf{y}_i, \mathbf{y}_i) = 0$  and  $\Delta(\mathbf{y}_i, \mathbf{y}_j) = \Delta(\mathbf{y}_j, \mathbf{y}_i)$ . Similar to the relation between conditional random fields and maximum entropy, if

$$\Delta(\mathbf{y}_i, \mathbf{y}_j) = \begin{cases} 0 & \text{if } \mathbf{y}_i = \mathbf{y}_j \\ 1 & \text{otherwise,} \end{cases}$$

and  $\mathbf{y}_i \in \{1, \dots, k\}, \forall i$ , then structured SVM becomes Crammer and Singer's problem in (33) following the definition of  $\mathbf{f}(\mathbf{x}, \mathbf{y})$  and  $\mathbf{w}$  in (38).

Like CRF, the main difficulty to solve (46) is on handling an exponential number of  $\mathbf{y}$  values. Some works (e.g., [25], [142], [148]) use a cutting plane method [149] to solve (46). In [150], a stochastic subgradient descent method is applied for both online and batch settings.

## B. Regression

Given training data  $\{(z_i, \mathbf{x}_i)\}_{i=1}^l \subset \mathbf{R} \times \mathbf{R}^n$ , a regression problem finds a weight vector  $\mathbf{w}$  such that  $\mathbf{w}^T \mathbf{x}_i \approx z_i, \forall i$ . Like classification, a regression task solves a risk minimization problem involving regularization and loss terms. While L1 and L2 regularization is still used, loss functions are different, where two popular ones are

$$\xi_{\text{LS}}(\mathbf{w}; \mathbf{x}, z) \equiv \frac{1}{2} (z - \mathbf{w}^T \mathbf{x})^2 \quad \text{and} \quad (47)$$

$$\xi_{\epsilon}(\mathbf{w}; \mathbf{x}, z) \equiv \max(0, |z - \mathbf{w}^T \mathbf{x}| - \epsilon). \quad (48)$$

The least square loss in (47) is widely used in many places, while the  $\epsilon$ -insensitive loss in (48) is extended from the L1 loss in (8), where there is a user-specified parameter  $\epsilon$  as the error tolerance. Problem (7) with L2 regularization and  $\epsilon$ -insensitive loss is called support vector regression (SVR) [151]. Contrary to the success of linear classification, so far not many applications of linear regression on large sparse data have been reported. We believe that this topic has not been fully explored yet.

Regarding the minimization of (7), if L2 regularization is used, many optimization methods mentioned in Section IV can be easily modified for linear regression.

We then particularly discuss L1-regularized least-square regression, which has recently drawn much attention for signal processing and image applications. This research area is so active that many optimization methods (e.g., [49], [152]–[156]) have been proposed. However, as pointed out in [13], optimization methods most suitable for signal/image applications via L1-regularized regression may be very different from those in Section IV for classifying large sparse data. One reason is that data from signal/image problems tend to be dense. Another is that  $\mathbf{x}_i, \forall i$  may be not directly available in some signal/image problems. Instead, we can only evaluate the product between the data matrix and a vector through certain operators. Thus, optimization methods that can take this property into their design may be more efficient.

## IX. CONCLUSIONS

In this article, we have comprehensively reviewed recent advances of large linear classification. For some applications, linear classifiers can give comparable accuracy to nonlinear classifiers, but enjoy much faster training and testing speed. However, these results do not imply that nonlinear classifiers should no longer be considered. Both linear and nonlinear classifiers are useful under different circumstances.

Without mapping data to another space, for linear classification we can easily prepare, select, and manipulate features. We have clearly shown that linear classification is not limited to standard scenarios like document classification. It can be applied in many other places such as efficiently approximating nonlinear classifiers. We are confident that future research works will make linear classification a useful technique for more large-scale applications.

## X. ACKNOWLEDGMENTS

This work was supported in part by the National Science Council of Taiwan via the grant 98-2221-E-002-136-MY3.

## REFERENCES

- [1] B. E. Boser, I. Guyon, and V. Vapnik, "A training algorithm for optimal margin classifiers," in *Proceedings of the Fifth Annual Workshop on Computational Learning Theory*. ACM Press, 1992, pp. 144–152.
- [2] C. Cortes and V. Vapnik, "Support-vector network," *Machine Learning*, vol. 20, pp. 273–297, 1995.
- [3] J. S. Cramer, "The origins of logistic regression," Tinbergen Institute, Tech. Rep., 2002. [Online]. Available: <http://ideas.repec.org/p/dgr/uvatin/20020119.html>
- [4] T. Joachims, "Training linear SVMs in linear time," in *Proceedings of the Twelfth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2006.
- [5] S. Shalev-Shwartz, Y. Singer, and N. Srebro, "Pegasos: primal estimated sub-gradient solver for SVM," in *Proceedings of the Twenty Fourth International Conference on Machine Learning (ICML)*, 2007.
- [6] C.-J. Hsieh, K.-W. Chang, C.-J. Lin, S. S. Keerthi, and S. Sundararajan, "A dual coordinate descent method for large-scale linear SVM," in *Proceedings of the Twenty Fifth International Conference on Machine Learning (ICML)*, 2008. [Online]. Available: <http://www.csie.ntu.edu.tw/~cjlin/papers/cddual.pdf>
- [7] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin, "LIBLINEAR: A library for large linear classification," *Journal of Machine Learning Research*, vol. 9, pp. 1871–1874, 2008. [Online]. Available: <http://www.csie.ntu.edu.tw/~cjlin/papers/liblinear.pdf>
- [8] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," *ACM Transactions on Intelligent Systems and Technology*, vol. 2, no. 3, pp. 27:1–27:27, 2011, software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [9] Y.-W. Chang, C.-J. Hsieh, K.-W. Chang, M. Ringgaard, and C.-J. Lin, "Training and testing low-degree polynomial data mappings via linear SVM," *Journal of Machine Learning Research*, vol. 11, pp. 1471–1490, 2010. [Online]. Available: [http://www.csie.ntu.edu.tw/~cjlin/papers/lowpoly\\_journal.pdf](http://www.csie.ntu.edu.tw/~cjlin/papers/lowpoly_journal.pdf)
- [10] S. S. Keerthi and C.-J. Lin, "Asymptotic behaviors of support vector machines with Gaussian kernel," *Neural Computation*, vol. 15, no. 7, pp. 1667–1689, 2003.
- [11] Z. S. Harris, "Distributional structure," *Word*, vol. 10, pp. 146–162, 1954.
- [12] F.-L. Huang, C.-J. Hsieh, K.-W. Chang, and C.-J. Lin, "Iterative scaling and coordinate descent methods for maximum entropy," *Journal of Machine Learning Research*, vol. 11, pp. 815–848, 2010. [Online]. Available: [http://www.csie.ntu.edu.tw/~cjlin/papers/maxent\\_journal.pdf](http://www.csie.ntu.edu.tw/~cjlin/papers/maxent_journal.pdf)
- [13] G.-X. Yuan, K.-W. Chang, C.-J. Hsieh, and C.-J. Lin, "A comparison of optimization methods and software for large-scale l1-regularized linear classification," *Journal of Machine Learning Research*, vol. 11, pp. 3183–3234, 2010. [Online]. Available: <http://www.csie.ntu.edu.tw/~cjlin/papers/l1.pdf>
- [14] O. L. Mangasarian, "A finite Newton method for classification," *Optimization Methods and Software*, vol. 17, no. 5, pp. 913–929, 2002.
- [15] A. Y. Ng, "Feature selection, L1 vs. L2 regularization, and rotational invariance," in *Proceedings of the Twenty First International Conference on Machine Learning (ICML)*, 2004.
- [16] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society Series B*, vol. 58, pp. 267–288, 1996.
- [17] D. L. Donoho and Y. Tsaig, "Fast solution of l1 minimization problems when the solution may be sparse," *IEEE Transactions on Information Theory*, vol. 54, pp. 4789–4812, 2008.
- [18] H. Zou and T. Hastie, "Regularization and variable selection via the elastic net," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 67, no. 2, pp. 301–320, 2005.
- [19] Y.-J. Lee and O. L. Mangasarian, "RSVM: Reduced support vector machines," in *Proceedings of the First SIAM International Conference on Data Mining*, 2001.
- [20] J. Shi, W. Yin, S. Osher, and P. Sajda, "A fast hybrid algorithm for large scale l1-regularized logistic regression," *Journal of Machine Learning Research*, vol. 11, pp. 713–741, 2010.
- [21] K.-W. Chang, C.-J. Hsieh, and C.-J. Lin, "Coordinate descent method for large-scale L2-loss linear SVM," *Journal of Machine Learning Research*, vol. 9, pp. 1369–1398, 2008. [Online]. Available: <http://www.csie.ntu.edu.tw/~cjlin/papers/cdl2.pdf>

- [22] G.-X. Yuan, C.-H. Ho, and C.-J. Lin, “An improved GLMNET for l1-regularized logistic regression,” *Journal of Machine Learning Research*, vol. 13, pp. 1999–2030, 2012. [Online]. Available: [http://www.csie.ntu.edu.tw/~cjlin/papers/l1\\_glmnet/long-glmnet.pdf](http://www.csie.ntu.edu.tw/~cjlin/papers/l1_glmnet/long-glmnet.pdf)
- [23] P. S. Bradley and O. L. Mangasarian, “Massive data discrimination via linear support vector machines,” *Optimization Methods and Software*, vol. 13, no. 1, pp. 1–10, 2000.
- [24] V. Franc and S. Sonnenburg, “Optimized cutting plane algorithm for support vector machines,” in *Proceedings of the 25th International Conference on Machine Learning*. ACM, 2008, pp. 320–327.
- [25] C. H. Teo, S. Vishwanathan, A. Smola, and Q. V. Le, “Bundle methods for regularized risk minimization,” *Journal of Machine Learning Research*, vol. 11, pp. 311–365, 2010.
- [26] L. Bottou, “Stochastic gradient descent examples,” 2007, <http://leon.bottou.org/projects/sgd>.
- [27] S. S. Keerthi and D. DeCoste, “A modified finite Newton method for fast solution of large scale linear SVMs,” *Journal of Machine Learning Research*, vol. 6, pp. 341–361, 2005.
- [28] C.-J. Lin, R. C. Weng, and S. S. Keerthi, “Trust region Newton method for large-scale logistic regression,” *Journal of Machine Learning Research*, vol. 9, pp. 627–650, 2008. [Online]. Available: <http://www.csie.ntu.edu.tw/~cjlin/papers/logistic.pdf>
- [29] T. P. Minka, “A comparison of numerical optimizers for logistic regression,” 2003, <http://research.microsoft.com/en-us/um/people/minka/papers/logreg/minka-logreg.pdf>.
- [30] J. Goodman, “Sequential conditional generalized iterative scaling,” in *Proceedings of the 40th Annual Meeting of the Association of Computational Linguistics (ACL)*, 2002, pp. 9–16.
- [31] R. Jin, R. Yan, J. Zhang, and A. G. Hauptmann, “A faster iterative scaling algorithm for conditional exponential model,” in *Proceedings of the Twentieth International Conference on Machine Learning (ICML)*, 2003, pp. 282–289.
- [32] P. Komarek and A. W. Moore, “Making logistic regression a core data mining tool: A practical investigation of accuracy, speed, and simplicity,” Robotics Institute, Carnegie Mellon University, Tech. Rep. TR-05-27, 2005.
- [33] H.-F. Yu, F.-L. Huang, and C.-J. Lin, “Dual coordinate descent methods for logistic regression and maximum entropy models,” *Machine Learning*, vol. 85, no. 1-2, pp. 41–75, October 2011. [Online]. Available: [http://www.csie.ntu.edu.tw/~cjlin/papers/maxent\\_dual.pdf](http://www.csie.ntu.edu.tw/~cjlin/papers/maxent_dual.pdf)
- [34] J. Zhu, S. Rosset, T. Hastie, and R. Tibshirani, “1-norm support vector machines,” in *Advances in Neural Information Processing Systems 16*, S. Thrun, L. Saul, and B. Schölkopf, Eds. Cambridge, MA: MIT Press, 2004.
- [35] G. M. Fung and O. L. Mangasarian, “A feature selection Newton method for support vector machine classification,” *Computational optimization and applications*, vol. 28, pp. 185–202, 2004.
- [36] O. L. Mangasarian, “Exact 1-norm support vector machines via unconstrained convex differentiable minimization,” *Journal of Machine Learning Research*, vol. 7, pp. 1517–1530, 2006.
- [37] K. Koh, S.-J. Kim, and S. Boyd, “An interior-point method for large-scale l1-regularized logistic regression,” *Journal of Machine Learning Research*, vol. 8, pp. 1519–1555, 2007. [Online]. Available: [http://www.stanford.edu/~boyd/l1\\_logistic\\_reg.html](http://www.stanford.edu/~boyd/l1_logistic_reg.html)
- [38] A. Genkin, D. D. Lewis, and D. Madigan, “Large-scale Bayesian logistic regression for text categorization,” *Technometrics*, vol. 49, no. 3, pp. 291–304, 2007.
- [39] S. Yun and K.-C. Toh, “A coordinate gradient descent method for l1-regularized convex minimization,” *Computational Optimizations and Applications*, vol. 48, no. 2, pp. 273–307, 2011.
- [40] G. Andrew and J. Gao, “Scalable training of L1-regularized log-linear models,” in *Proceedings of the Twenty Fourth International Conference on Machine Learning (ICML)*, 2007.
- [41] J. H. Friedman, T. Hastie, and R. Tibshirani, “Regularization paths for generalized linear models via coordinate descent,” *Journal of Statistical Software*, vol. 33, no. 1, pp. 1–22, 2010.
- [42] J. Liu, J. Chen, and J. Ye, “Large-scale sparse logistic regression,” in *Proceedings of The 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2009, pp. 547–556.
- [43] R. Tomioka, T. Suzuki, and M. Sugiyama, “Super-linear convergence of dual augmented Lagrangian algorithm for sparse learning,” *Journal of Machine Learning Research*, vol. 12, pp. 1537–1586, 2011.
- [44] M. Schmidt, G. Fung, and R. Rosales, “Optimization methods for l1-regularization,” University of British Columbia, Technical Report TR-2009-19, 2009.
- [45] C.-J. Lin and J. J. Moré, “Newton’s method for large-scale bound constrained problems,” *SIAM Journal on Optimization*, vol. 9, pp. 1100–1127, 1999.
- [46] Z.-Q. Luo and P. Tseng, “On the convergence of coordinate descent method for convex differentiable minimization,” *Journal of Optimization Theory and Applications*, vol. 72, no. 1, pp. 7–35, 1992.
- [47] T. Joachims, “Making large-scale SVM learning practical,” in *Advances in Kernel Methods – Support Vector Learning*, B. Schölkopf, C. J. C. Burges, and A. J. Smola, Eds. Cambridge, MA: MIT Press, 1998, pp. 169–184.
- [48] J. H. Friedman, T. Hastie, H. Höfling, and R. Tibshirani, “Pathwise coordinate optimization,” *Annals of Applied Statistics*, vol. 1, no. 2, pp. 302–332, 2007.
- [49] S. J. Wright, R. D. Nowak, and M. A. Figueiredo, “Sparse reconstruction by separable approximation,” *IEEE Transactions on Signal Processing*, vol. 57, pp. 2479–2493, 2009.
- [50] C.-W. Hsu and C.-J. Lin, “A comparison of methods for multi-class support vector machines,” *IEEE Transactions on Neural Networks*, vol. 13, no. 2, pp. 415–425, 2002.
- [51] R. Rifkin and A. Klautau, “In defense of one-vs-all classification,” *Journal of Machine Learning Research*, vol. 5, pp. 101–141, 2004.
- [52] E. L. Allwein, R. E. Schapire, and Y. Singer, “Reducing multiclass to binary: a unifying approach for margin classifiers,” *Journal of Machine Learning Research*, vol. 1, pp. 113–141, 2001.
- [53] T.-K. Huang, R. C. Weng, and C.-J. Lin, “Generalized Bradley-Terry models and multi-class probability estimates,” *Journal of Machine Learning Research*, vol. 7, pp. 85–115, 2006. [Online]. Available: <http://www.csie.ntu.edu.tw/~cjlin/papers/generalBT.pdf>
- [54] L. Bottou, C. Cortes, J. S. Denker, H. Drucker, I. Guyon, L. Jackel, Y. LeCun, U. A. Müller, E. Säckinger, P. Simard, and V. Vapnik, “Comparison of classifier methods: a case study in handwriting digit recognition,” in *International Conference on Pattern Recognition*, 1994, pp. 77–87.

- [55] S. Knerr, L. Personnaz, and G. Dreyfus, "Single-layer learning revisited: a stepwise procedure for building and training a neural network," in *Neurocomputing: Algorithms, Architectures and Applications*, J. Fogelman, Ed. Springer-Verlag, 1990.
- [56] J. H. Friedman, "Another approach to polychotomous classification," Department of Statistics, Stanford University, Tech. Rep., 1996. [Online]. Available: <http://www-stat.stanford.edu/~jhf/ftp/poly.pdf>
- [57] T.-L. Huang, "Comparison of L2-regularized multi-class linear classifiers," Master's thesis, Department of Computer Science and Information Engineering, National Taiwan University, 2010.
- [58] J. C. Platt, N. Cristianini, and J. Shawe-Taylor, "Large margin DAGs for multiclass classification," in *Advances in Neural Information Processing Systems*, vol. 12. MIT Press, 2000, pp. 547–553.
- [59] J. Weston and C. Watkins, "Multi-class support vector machines," in *Proceedings of ESANN99*, M. Verleysen, Ed. Brussels: D. Facto Press, 1999, pp. 219–224.
- [60] K. Crammer and Y. Singer, "On the algorithmic implementation of multiclass kernel-based vector machines," *Journal of Machine Learning Research*, vol. 2, pp. 265–292, 2001.
- [61] Y. Lee, Y. Lin, and G. Wahba, "Multicategory support vector machines," *Journal of the American Statistical Association*, vol. 99, no. 465, pp. 67–81, 2004.
- [62] C.-J. Lin, "A formal analysis of stopping criteria of decomposition methods for support vector machines," *IEEE Transactions on Neural Networks*, vol. 13, no. 5, pp. 1045–1052, 2002. [Online]. Available: <http://www.csie.ntu.edu.tw/~cjlin/papers/stop.ps.gz>
- [63] S. S. Keerthi, S. Sundararajan, K.-W. Chang, C.-J. Hsieh, and C.-J. Lin, "A sequential dual method for large scale multi-class linear SVMs," in *Proceedings of the Forteenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2008, pp. 408–416. [Online]. Available: [http://www.csie.ntu.edu.tw/~cjlin/papers/sdm\\_kdd.pdf](http://www.csie.ntu.edu.tw/~cjlin/papers/sdm_kdd.pdf)
- [64] A. L. Berger, V. J. Della Pietra, and S. A. Della Pietra, "A maximum entropy approach to natural language processing," *Computational Linguistics*, vol. 22, no. 1, pp. 39–71, 1996.
- [65] J. Lafferty, A. McCallum, and F. Pereira, "Conditional random fields: Probabilistic models for segmenting and labeling sequence data," in *Proceedings of the 18th International Conference on Machine Learning (ICML)*, 2001, pp. 282–289.
- [66] R. Malouf, "A comparison of algorithms for maximum entropy parameter estimation," in *Proceedings of the 6th conference on Natural language learning*. Association for Computational Linguistics, 2002, pp. 1–7.
- [67] J. N. Darroch and D. Ratcliff, "Generalized iterative scaling for log-linear models," *The Annals of Mathematical Statistics*, vol. 43, no. 5, pp. 1470–1480, 1972.
- [68] D. C. Liu and J. Nocedal, "On the limited memory BFGS method for large scale optimization," *Mathematical Programming*, vol. 45, no. 1, pp. 503–528, 1989.
- [69] S. Della Pietra, V. Della Pietra, and J. Lafferty, "Inducing features of random fields," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 19, no. 4, pp. 380–393, 1997.
- [70] R. Memisevic, "Dual optimization of conditional probability models," Department of Computer Science, University of Toronto, Tech. Rep., 2006.
- [71] M. Collins, A. Globerson, T. Koo, X. Carreras, and P. Bartlett, "Exponentiated gradient algorithms for conditional random fields and max-margin Markov networks," *Journal of Machine Learning Research*, vol. 9, pp. 1775–1822, 2008.
- [72] E. M. Gertz and J. D. Griffin, "Support vector machine classifiers for large data sets," Argonne National Laboratory, Tech. Rep. ANL/MCS-TM-289, 2005.
- [73] J. H. Jung, D. P. O'Leary, and A. L. Tits, "Adaptive constraint reduction for training support vector machines," *Electronic Transactions on Numerical Analysis*, vol. 31, pp. 156–177, 2008.
- [74] Y. Moh and J. M. Buhmann, "Kernel expansion for online preference tracking," in *Proceedings of The International Society for Music Information Retrieval (ISMIR)*, 2008, pp. 167–172.
- [75] S. Sonnenburg and V. Franc, "COFFIN : A computational framework for linear SVMs," in *Proceedings of the Twenty Seventh International Conference on Machine Learning (ICML)*, 2010, pp. 999–1006.
- [76] G. Ifrim, G. Bakir, and G. Weikum, "Fast logistic regression for text categorization with variable-length n-grams," in *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2008, pp. 354–362.
- [77] G. Ifrim and C. Wiuf, "Bounded coordinate-descent for biological sequence classification in high dimensional predictor space," in *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2011.
- [78] S. Lee and S. J. Wright, "ASSET: Approximate stochastic subgradient estimation training for support vector machines," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2010, submitted.
- [79] C. K. I. Williams and M. Seeger, "Using the Nyström method to speed up kernel machines," in *Neural Information Processing Systems 13*, T. Leen, T. Dietterich, and V. Tresp, Eds. MIT Press, 2001, pp. 682–688.
- [80] P. Drineas and M. W. Mahoney, "On the Nyström method for approximating a gram matrix for improved kernel-based learning," *Journal of Machine Learning Research*, vol. 6, pp. 2153–2175, 2005.
- [81] S. Fine and K. Scheinberg, "Efficient SVM training using low-rank kernel representations," *Journal of Machine Learning Research*, vol. 2, pp. 243–264, 2001.
- [82] F. R. Bach and M. I. Jordan, "Predictive low-rank decomposition for kernel methods," in *Proceedings of the Twenty Second International Conference on Machine Learning*, 2005, pp. 33–40.
- [83] A. Rahimi and B. Recht, "Random features for large-scale kernel machines," *Advances in Neural Information Processing Systems*, pp. 1177–1184, 2008.
- [84] D. Achlioptas, "Database-friendly random projections: Johnson-Lindenstrauss with binary coins," *Journal of Computer and System Sciences*, vol. 66, pp. 671–687, 2003.
- [85] P. Li, T. J. Hastie, and K. W. Church, "Very sparse random projections," in *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2006, pp. 287–296.
- [86] P. Kar and H. Karnick, "Random feature maps for dot product kernels," in *Proceedings of the 15th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2012, pp. 583–591.

- [87] N. Pham and R. Pagh, “Fast and scalable polynomial kernels via explicit feature maps,” in *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2013, pp. 239–247.
- [88] K.-P. Lin and M.-S. Chen, “Efficient kernel approximation for large-scale support vector machine classification,” in *Proceedings of the Eleventh SIAM International Conference on Data Mining*, 2011, pp. 211–222.
- [89] Q. Shi, J. Petterson, G. Dror, J. Langford, A. Smola, A. Strehl, and S. Vishwanathan, “Hash kernels,” in *JMLR Workshop and Conference Proceedings: Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics*, vol. 5, 2009, pp. 496–503.
- [90] K. Weinberger, A. Dasgupta, J. Langford, A. Smola, and J. Attenberg, “Feature hashing for large scale multitask learning,” in *Proceedings of the Twenty Sixth International Conference on Machine Learning (ICML)*, 2009, pp. 1113–1120.
- [91] P. Li and A. C. König, “b-Bit minwise hashing,” in *Proceedings of the Nineteenth International Conference on World Wide Web*, 2010, pp. 671–680.
- [92] —, “Theory and applications of b-Bit minwise hashing,” *Communications of the ACM*, vol. 54, no. 8, pp. 101–109, 2011.
- [93] A. Rakotomamonjy, F. Bach, S. Canu, and Y. Grandvalet, “SimpleMKL,” *Journal of Machine Learning Research*, vol. 9, pp. 2491–2521, 2008.
- [94] P. Li, A. Shrivastava, J. Moore, and A. C. König, “Hashing algorithms for large-scale learning,” Cornell University, Tech. Rep., 2011. [Online]. Available: <http://www.stat.cornell.edu/~li/reports/HashLearning.pdf>
- [95] H.-F. Yu, C.-J. Hsieh, K.-W. Chang, and C.-J. Lin, “Large linear classification when data cannot fit in memory,” in *Proceedings of the Sixteenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2010, pp. 833–842. [Online]. Available: [http://www.csie.ntu.edu.tw/~cjlin/papers/kdd\\_disk\\_decomposition.pdf](http://www.csie.ntu.edu.tw/~cjlin/papers/kdd_disk_decomposition.pdf)
- [96] E. Chang, K. Zhu, H. Wang, H. Bai, J. Li, Z. Qiu, and H. Cui, “Parallelizing support vector machines on distributed computers,” in *Advances in Neural Information Processing Systems 20*, J. Platt, D. Koller, Y. Singer, and S. Roweis, Eds. Cambridge, MA: MIT Press, 2008, pp. 257–264.
- [97] Z. A. Zhu, W. Chen, G. Wang, C. Zhu, and Z. Chen, “P-packSVM: Parallel primal gradient descent kernel SVM,” in *Proceedings of the IEEE International Conference on Data Mining*, 2009.
- [98] T. White, *Hadoop: The definitive guide*, 2nd ed. O’Reilly Media, 2010.
- [99] H. Robbins and S. Monro, “A stochastic approximation method,” *The Annals of Mathematical Statistics*, vol. 22, no. 3, pp. 400–407, 1951.
- [100] J. Kiefer and J. Wolfowitz, “Stochastic estimation of the maximum of a regression function,” *The Annals of Mathematical Statistics*, vol. 23, no. 3, pp. 462–466, 1952.
- [101] T. Zhang, “Solving large scale linear prediction problems using stochastic gradient descent algorithms,” in *Proceedings of the 21th International Conference on Machine Learning (ICML)*, 2004.
- [102] L. Bottou and Y. LeCun, “Large scale online learning,” *Advances in Neural Information Processing Systems*, vol. 16, pp. 217–224, 2004.
- [103] A. Bordes, S. Ertekin, J. Weston, and L. Bottou, “Fast kernel classifiers with online and active learning,” *Journal of Machine Learning Research*, vol. 6, pp. 1579–1619, 2005.
- [104] Y. E. Nesterov, “Efficiency of coordinate descent methods on huge-scale optimization problems,” *SIAM Journal on Optimization*, vol. 22, no. 2, pp. 341–362, 2012.
- [105] P. Richtárik and M. Takáč, “Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function,” *Mathematical Programming*, vol. 144, pp. 1–38, 2014.
- [106] A. Bordes, L. Bottou, and P. Gallinari, “SGD-QN: Careful quasi-Newton stochastic gradient descent,” *Journal of Machine Learning Research*, vol. 10, pp. 1737–1754, 2009.
- [107] A. Bordes, L. Bottou, P. Gallinari, J. Chang, and S. A. Smith, “Erratum: SGD-QN is less careful than expected,” *Journal of Machine Learning Research*, vol. 11, pp. 2229–2240, 2010.
- [108] J. Langford, L. Li, and T. Zhang, “Sparse online learning via truncated gradient,” *Journal of Machine Learning Research*, vol. 10, pp. 771–801, 2009.
- [109] S. Shalev-Shwartz and A. Tewari, “Stochastic methods for  $l_1$ -regularized loss minimization,” *Journal of Machine Learning Research*, vol. 12, pp. 1865–1892, 2011.
- [110] Y. E. Nesterov, “Primal-dual subgradient methods for convex problems,” *Mathematical Programming*, vol. 120, no. 1, pp. 221–259, 2009.
- [111] J. Duchi and Y. Singer, “Efficient online and batch learning using forward backward splitting,” *Journal of Machine Learning Research*, vol. 10, pp. 2899–2934, 2009.
- [112] J. Duchi, E. Hazan, and Y. Singer, “Adaptive subgradient methods for online learning and stochastic optimization,” *Journal of Machine Learning Research*, vol. 12, pp. 2121–2159, 2011.
- [113] L. Xiao, “Dual averaging methods for regularized stochastic learning and online optimization,” *Journal of Machine Learning Research*, vol. 11, pp. 2543–2596, 2010.
- [114] J. K. Bradley, A. Kyrola, D. Bickson, and C. Guestrin, “Parallel coordinate descent for  $l_1$ -regularized loss minimization,” in *Proceedings of the Twenty Eighth International Conference on Machine Learning (ICML)*, 2011, pp. 321–328.
- [115] J. Langford, L. Li, and A. Strehl, “Vowpal Wabbit,” 2007, [https://github.com/JohnLangford/vowpal\\_wabbit/wiki](https://github.com/JohnLangford/vowpal_wabbit/wiki).
- [116] A. Nedić, D. P. Bertsekas, and V. S. Borkar, “Distributed asynchronous incremental subgradient methods,” *Studies in Computational Mathematics*, vol. 8, pp. 381–407, 2001.
- [117] J. Langford, A. Smola, and M. Zinkevich, “Slow learners are fast,” in *Advances in Neural Information Processing Systems 22*, Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, and A. Culotta, Eds., 2009, pp. 2331–2339.
- [118] A. Agarwal and J. Duchi, “Distributed delayed stochastic optimization,” in *Advances in Neural Information Processing Systems 24*, J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K. Weinberger, Eds., 2011, pp. 873–881.
- [119] O. Dekel, R. Gilad-Bachrach, O. Shamir, and L. Xiao, “Optimal distributed online prediction using mini-batches,” *Journal of Machine Learning Research*, vol. 13, pp. 165–202, 2012.

- [120] D. Hsu, N. Karampatziakis, and J. Langford, "Parallel online learning," in *Scaling Up Machine Learning*, R. Bekkerman, M. Bilenko, and J. Langford, Eds. Cambridge University Press, 2011, ch. 14.
- [121] M. Li, T. Zhang, Y. Chen, and A. J. Smola, "Efficient mini-batch training for stochastic optimization," in *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2014, pp. 661–670.
- [122] S. Tong, "Lessons learned developing a practical large scale machine learning system," Google Research Blog, 2010, <http://googleresearch.blogspot.com/2010/04/lessons-learned-developing-practical.html>.
- [123] M. Ferris and T. Munson, "Interior point methods for massive support vector machines," *SIAM Journal on Optimization*, vol. 13, no. 3, pp. 783–804, 2003.
- [124] K.-W. Chang and D. Roth, "Selective block minimization for faster convergence of limited memory large-scale linear models," in *Proceedings of the Seventeenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2011.
- [125] Y. Censor and S. A. Zenios, *Parallel Optimization: Theory, Algorithms, and Applications*. Oxford University Press, 1998.
- [126] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [127] D. Gabay and B. Mercier, "A dual algorithm for the solution of nonlinear variational problems via finite element approximation," *Computers and Mathematics with Applications*, vol. 2, pp. 17–40, 1976.
- [128] P. A. Forero, A. Cano, and G. B. Giannakis, "Consensus-based distributed support vector machines," *Journal of Machine Learning*, vol. 11, pp. 1663–1707, 2010.
- [129] P. Richtárik and M. Takáč, "Parallel coordinate descent methods for big data optimization," *Mathematical Programming*, 2012, under revision.
- [130] Y. Bian, X. Li, M. Cao, and Y. Liu, "Bundle CDN: A highly parallelized approach for large-scale  $\ell_1$ -regularized logistic regression," in *Proceedings of European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML/PKDD)*, 2013.
- [131] M. Jaggi, V. Smith, M. Takáč, J. Terhorst, T. Hofmann, and M. I. Jordan, "Communication-efficient distributed dual coordinate ascent," in *Advances in Neural Information Processing Systems 27*, 2014.
- [132] A. Agarwal, O. Chapelle, M. Dudík, and J. Langford, "A reliable effective terascale linear learning system," *Journal of Machine Learning Research*, vol. 15, pp. 1111–1133, 2014.
- [133] O. Chapelle, E. Manavoglu, and R. Rosales, "Simple and scalable response prediction for display advertising," *ACM Transactions on Intelligent Systems and Technology*, 2014, to appear.
- [134] Y. Zhuang, W.-S. Chin, Y.-C. Juan, and C.-J. Lin, "Distributed Newton method for regularized logistic regression," Department of Computer Science and Information Engineering, National Taiwan University, Tech. Rep., 2014.
- [135] C.-Y. Lin, C.-H. Tsai, C.-P. Lee, and C.-J. Lin, "Large-scale logistic regression and linear support vector machines using Spark," in *Proceedings of the IEEE International Conference on Big Data*, 2014. [Online]. Available: <http://www.csie.ntu.edu.tw/~cjlin/papers/spark-liblinear/spark-liblinear.pdf>
- [136] H. Yu, J. Yang, and J. Han, "Classifying large data sets using SVMs with hierarchical clusters," in *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. New York, NY, USA: ACM Press, 2003, pp. 306–315.
- [137] L. Breiman, "Bagging predictors," *Machine Learning*, vol. 24, no. 2, pp. 123–140, August 1996.
- [138] D. Chakrabarti, D. Agarwal, and V. Josifovski, "Contextual advertising by combining relevance with click feedback," in *Proceeding of the 17th international conference on World Wide Web*, 2008, pp. 417–426.
- [139] M. Zinkevich, M. Weimer, A. Smola, and L. Li, "Parallelized stochastic gradient descent," in *Advances in Neural Information Processing Systems 23*, J. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta, Eds., 2010, pp. 2595–2603.
- [140] R. McDonald, K. Hall, and G. Mann, "Distributed training strategies for the structured perceptron," in *Proceedings of the 48th Annual Meeting of the Association of Computational Linguistics (ACL)*, 2010, pp. 456–464.
- [141] J. Dean and S. Ghemawat, "MapReduce: simplified data processing on large clusters," *Communications of the ACM*, vol. 51, no. 1, pp. 107–113, 2008.
- [142] I. Tsochantaridis, T. Joachims, T. Hofmann, and Y. Altun, "Large margin methods for structured and interdependent output variables," *Journal of Machine Learning Research*, vol. 6, pp. 1453–1484, 2005.
- [143] B. Taskar, C. Guestrin, and D. Koller, "Max-margin markov networks," in *Advances in Neural Information Processing Systems 16*. Cambridge, MA: MIT Press, 2004.
- [144] F. Sha and F. C. N. Pereira, "Shallow parsing with conditional random fields," in *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology (NAACL)*, 2003, pp. 134–141.
- [145] S. Vishwanathan, N. N. Schraudolph, M. W. Schmidt, and K. Murphy, "Accelerated training of conditional random fields with stochastic gradient methods," in *Proceedings of the Twenty Third International Conference on Machine Learning (ICML)*, 2006, pp. 969–976.
- [146] N. N. Schraudolph, J. Yu, and S. Gunter, "A stochastic quasi-Newton method for online convex optimization," in *Proceedings of the 11th International Conference Artificial Intelligence and Statistics (AISTATS)*, 2007, pp. 433–440.
- [147] P.-J. Chen, "Newton methods for conditional random fields," Master's thesis, Department of Computer Science and Information Engineering, National Taiwan University, 2009.
- [148] T. Joachims, T. Finley, and C.-N. J. Yu, "Cutting-plane training of structural SVMs," *Machine Learning Journal*, vol. 77, no. 1, pp. 27–59, 2009.
- [149] J. E. Kelley, "The cutting-plane method for solving convex programs," *Journal of the Society for Industrial and Applied Mathematics*, 1960.
- [150] N. D. Ratliff, J. A. Bagnell, and M. A. Zinkevich, "(Online) subgradient methods for structured prediction," in *Proceedings of the Eleventh International Conference Artificial Intelligence and Statistics (AISTATS)*, vol. 2, 2007, pp. 380–387.
- [151] V. Vapnik, *Statistical Learning Theory*. New York, NY: Wiley, 1998.
- [152] I. Daubechies, M. Defrise, and C. De Mol, "An iterative thresholding algorithm for linear inverse problems with a sparsity constraint," *Communications on Pure and Applied Mathematics*, vol. 57, pp. 1413–1457, 2004.



- [153] M. A. T. Figueiredo, R. Nowak, and S. Wright, "Gradient projection for sparse reconstruction: Applications to compressed sensing and other inverse problems," *IEEE Journal on Selected Topics in Signal Processing*, vol. 1, pp. 586–598, 2007.
- [154] S.-J. Kim, K. Koh, M. Lustig, S. Boyd, and D. Gorinevsky, "An interior point method for large-scale  $l_1$ -regularized least squares," *IEEE Journal on Selected Topics in Signal Processing*, vol. 1, pp. 606–617, 2007.
- [155] J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra, "Efficient projections onto the  $L_1$ -ball for learning in high dimensions," in *Proceedings of the Twenty Fifth International Conference on Machine Learning (ICML)*, 2008.
- [156] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM Journal on Imaging Sciences*, vol. 2, no. 1, pp. 183–202, 2009.



**Guo-Xun Yuan** is currently a PhD student at University of California Davis. He has BS and MS degrees at National Tsinghua University and National Taiwan University, respectively. His research interest is large-scale data classification.



**Chia-Hua Ho** is a master student in the Department of Computer Science at National Taiwan University. His research interests are machine learning and data mining. He obtained his B.S. degree from National Taiwan University in 2010.



**Chih-Jen Lin** is currently a distinguished professor at the Department of Computer Science, National Taiwan University. He obtained his B.S. degree from National Taiwan University in 1993 and Ph.D. degree from University of Michigan in 1998. His major research areas include machine learning, data mining, and numerical optimization. He is best known for his work on support vector machines (SVM) for data classification. His software LIBSVM is one of the most widely used and cited SVM packages. Nearly all major companies apply his software for classification and regression applications. He has received many awards for his research work. A recent one is the ACM KDD 2010 best paper award. He is an IEEE fellow and an ACM distinguished scientist for his contribution to machine learning algorithms and software design. More information about him and his software tools can be found at <http://www.csie.ntu.edu.tw/~cjlin>.