

# Optimization Methods for Large-Scale Machine Learning\*

Léon Bottou<sup>†</sup>  
Frank E. Curtis<sup>‡</sup>  
Jorge Nocedal<sup>§</sup>

**Abstract.** This paper provides a review and commentary on the past, present, and future of numerical optimization algorithms in the context of machine learning applications. Through case studies on text classification and the training of deep neural networks, we discuss how optimization problems arise in machine learning and what makes them challenging. A major theme of our study is that large-scale machine learning represents a distinctive setting in which the stochastic gradient (SG) method has traditionally played a central role while conventional gradient-based nonlinear optimization techniques typically falter. Based on this viewpoint, we present a comprehensive theory of a straightforward, yet versatile SG algorithm, discuss its practical behavior, and highlight opportunities for designing algorithms with improved performance. This leads to a discussion about the next generation of optimization methods for large-scale machine learning, including an investigation of two main streams of research on techniques that diminish noise in the stochastic directions and methods that make use of second-order derivative approximations.

**Key words.** numerical optimization, machine learning, stochastic gradient methods, algorithm complexity analysis, noise reduction methods, second-order methods

**AMS subject classifications.** 65K05, 68Q25, 68T05, 90C06, 90C30, 90C90

**DOI.** 10.1137/16M1080173

## Contents

<b>1</b>	<b>Introduction</b>	<b>224</b>
<b>2</b>	<b>Machine Learning Case Studies</b>	<b>226</b>
2.1	Text Classification via Convex Optimization . . . . .	226
2.2	Perceptual Tasks via Deep Neural Networks . . . . .	228
2.3	Formal Machine Learning Procedure . . . . .	231
<b>3</b>	<b>Overview of Optimization Methods</b>	<b>235</b>
3.1	Formal Optimization Problem Statements . . . . .	235

\*Received by the editors June 16, 2016; accepted for publication (in revised form) April 19, 2017; published electronically May 8, 2018.

<http://www.siam.org/journals/sirev/60-2/M108017.html>

**Funding:** The work of the second author was supported by U.S. Department of Energy grant DE-SC0010615 and U.S. National Science Foundation grant DMS-1016291. The work of the third author was supported by Office of Naval Research grant N00014-14-1-0313 P00003 and Department of Energy grant DE-FG02-87ER25047s.

<sup>†</sup>Facebook AI Research, New York, NY 10003 (leon@bottou.org).

<sup>‡</sup>Department of Industrial and Systems Engineering, Lehigh University, Bethlehem, PA 18015 (frank.e.curtis@gmail.com).

<sup>§</sup>Department of Industrial Engineering and Management Sciences, Northwestern University, Evanston, IL 60201 (j-nocedal@northwestern.edu).

3.2	Stochastic vs. Batch Optimization Methods . . . . .	236
3.3	Motivation for Stochastic Methods . . . . .	237
3.4	Beyond SG: Noise Reduction and Second-Order Methods . . . . .	241
<b>4</b>	<b>Analyses of Stochastic Gradient Methods</b>	<b>242</b>
4.1	Two Fundamental Lemmas . . . . .	244
4.2	SG for Strongly Convex Objectives . . . . .	246
4.3	SG for General Objectives . . . . .	252
4.4	Work Complexity for Large-Scale Learning . . . . .	255
4.5	Commentary . . . . .	258
<b>5</b>	<b>Noise Reduction Methods</b>	<b>259</b>
5.1	Reducing Noise at a Geometric Rate . . . . .	261
5.2	Dynamic Sample Size Methods . . . . .	262
5.3	Gradient Aggregation . . . . .	265
5.4	Iterate Averaging Methods . . . . .	269
<b>6</b>	<b>Second-Order Methods</b>	<b>270</b>
6.1	Hessian-Free Inexact Newton Methods . . . . .	271
6.2	Stochastic Quasi-Newton Methods . . . . .	275
6.3	Gauss-Newton Methods . . . . .	279
6.4	Natural Gradient Method . . . . .	280
6.5	Methods that Employ Diagonal Scalings . . . . .	284
<b>7</b>	<b>Other Popular Methods</b>	<b>288</b>
7.1	Gradient Methods with Momentum . . . . .	288
7.2	Accelerated Gradient Methods . . . . .	289
7.3	Coordinate Descent Methods . . . . .	290
<b>8</b>	<b>Methods for Regularized Models</b>	<b>293</b>
8.1	First-Order Methods for Generic Convex Regularizers . . . . .	294
8.2	Second-Order Methods . . . . .	298
<b>9</b>	<b>Summary and Perspectives</b>	<b>301</b>
	<b>Appendix A. Convexity and Analyses of SG</b>	<b>302</b>
	<b>Appendix B. Proofs</b>	<b>303</b>
	<b>References</b>	<b>304</b>

**I. Introduction.** The promise of artificial intelligence has been a topic of both public and private interest for decades. Starting in the 1950s, there has been great hope that classical artificial intelligence techniques based on logic, knowledge representation, reasoning, and planning would result in revolutionary software that could, among other things, understand language, control robots, and provide expert advice. Although advances based on such techniques may be in store in the future, many researchers have started to doubt these classical approaches, choosing instead to focus their efforts on the design of systems based on statistical techniques, such as in the rapidly evolving and expanding field of *machine learning*.

Machine learning and the intelligent systems that have been borne out of it—such as search engines, recommendation platforms, and speech and image recognition software—have become an indispensable part of modern society. Rooted in statistics and relying heavily on the efficiency of numerical algorithms, machine learning techniques capitalize on the world’s increasingly powerful computing platforms and the availability of datasets of immense size. In addition, as the fruits of its efforts have become so easily accessible to the public through various modalities—such as *the cloud*—interest in machine learning is bound to continue its dramatic rise, yielding further societal, economic, and scientific impacts.

One of the pillars of machine learning is *mathematical optimization*, which, in this context, involves the numerical computation of parameters for a system designed to make decisions based on as yet unseen data. That is, based on currently available data, these parameters are chosen to be optimal with respect to a given learning problem. The success of certain optimization methods for machine learning has inspired great numbers of researchers in various communities to tackle even more challenging machine learning problems and to design new methods that are more widely applicable.

The purpose of this paper is to provide a review and commentary on the past, present, and future of the use of numerical optimization algorithms in the context of machine learning applications. A major theme of this work is that large-scale machine learning represents a distinctive setting in which traditional nonlinear optimization techniques typically falter, and so should be considered secondary to alternative classes of approaches that respect the statistical nature of the underlying problem of interest.

Overall, this paper attempts to provide answers for the following questions:

1. How do optimization problems arise in machine learning applications, and what makes them challenging?
2. What have been the most successful optimization methods for large-scale machine learning, and why?
3. What recent advances have been made in the design of algorithms, and what are open questions in this research area?

We answer the first question with the aid of two case studies. The first, a study of text classification, represents an application in which the success of machine learning has been widely recognized and celebrated. The second, a study of perceptual tasks such as speech or image recognition, represents an application in which machine learning still has had great success, but in a much more enigmatic manner that leaves many questions unanswered. These case studies also illustrate the variety of optimization problems that arise in machine learning: the first involves *convex* optimization problems—derived from the use of logistic regression or support vector machines—while the second typically involves *highly nonlinear and nonconvex* problems—derived from the use of deep neural networks.

With these case studies in hand, we turn our attention to the latter two questions on optimization algorithms, the discussions around which represent the bulk of the paper. Whereas traditional gradient-based methods may be effective for solving small-scale learning problems in which a *batch* approach may be used, in the context of large-scale machine learning it has been a *stochastic* algorithm—namely, the *stochastic gradient* method (SG) proposed by Robbins and Monro [130]—that has been the core strategy of interest. Due to this central role played by SG, we discuss its fundamental theoretical and practical properties within a few contexts of interest. We also discuss recent trends in the design of optimization methods for machine learning, organizing them according to their relationship to SG. We discuss (i) noise reduction methods that attempt to borrow from the strengths of batch methods, such

as their fast convergence rates and ability to exploit parallelism; (ii) methods that incorporate approximate second-order derivative information with the goal of dealing with nonlinearity and ill-conditioning; and (iii) methods for solving regularized problems designed to avoid overfitting and allow for the use of high-dimensional models. Rather than contrast SG and other methods based on the results of numerical experiments—which might bias our review toward a limited test set and implementation details—we focus our attention on fundamental computational trade-offs and theoretical properties of optimization methods.

We close the paper with a look back at what has been learned, as well as additional thoughts about the future of optimization methods for machine learning.

**2. Machine Learning Case Studies.** Optimization problems arise throughout machine learning. We provide two case studies that illustrate their role in the selection of prediction functions in state-of-the-art machine learning systems. We focus on cases that involve very large datasets and for which the number of model parameters to be optimized is also large. By remarking on the structure and scale of such problems, we provide a glimpse into the challenges that make them difficult to solve.

**2.1. Text Classification via Convex Optimization.** The assignment of natural language text to predefined classes based on their contents is one of the fundamental tasks of information management [56]. Consider, for example, the task of determining whether a text document is one that discusses politics. Educated humans can make a determination of this type unambiguously, say, by observing that the document contains the names of well-known politicians. Early text classification systems attempted to consolidate knowledge from human experts by building on such observations to formally characterize the word sequences that signify a discussion about a topic of interest (e.g., politics). Unfortunately, however, concise characterizations of this type are difficult to formulate. Rules about which word sequences do or do not signify a topic need to be refined when new documents arise that cannot be classified accurately based on previously established rules. The need to coordinate such a growing collection of possibly contradictory rules limits the applicability of such systems to relatively simple tasks.

By contrast, the statistical machine learning approach begins with the collection of a sizable set of examples  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ , where for each  $i \in \{1, \dots, n\}$  the vector  $x_i$  represents the *features* of a text document (e.g., the words it includes) and the scalar  $y_i$  is a *label* indicating whether the document belongs ( $y_i = 1$ ) or not ( $y_i = -1$ ) to a particular class (i.e., topic of interest). With such a set of examples, one can construct a classification program, defined by a *prediction function*  $h$ , and measure its performance by counting how often the program prediction  $h(x_i)$  differs from the correct prediction  $y_i$ . In this manner, it seems judicious to search for a prediction function that minimizes the frequency of observed misclassifications, otherwise known as the *empirical risk* of misclassification:

$$(2.1) \quad R_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}[h(x_i) \neq y_i], \quad \text{where } \mathbb{1}[A] = \begin{cases} 1 & \text{if } A \text{ is true,} \\ 0 & \text{otherwise.} \end{cases}$$

The idea of minimizing such a function gives rise to interesting conceptual issues. Consider, for example, a function that simply memorizes the examples, such as

$$(2.2) \quad h^{\text{rote}}(x) = \begin{cases} y_i & \text{if } x = x_i \text{ for some } i \in \{1, \dots, n\}, \\ \pm 1 & \text{(arbitrarily) otherwise.} \end{cases}$$

This prediction function clearly minimizes (2.1), but it offers no performance guarantees on documents that do not appear in the examples. To avoid such rote memorization, one should aim to find a prediction function that *generalizes* the concepts that may be learned from the examples.

One way to achieve good generalized performance is to choose among a carefully selected class of prediction functions, perhaps satisfying certain smoothness conditions and enjoying a convenient parametric representation. How such a class of functions may be selected is not straightforward; we discuss this issue in more detail in section 2.3. For now, we only mention that common practice for choosing between prediction functions belonging to a given class is to compare them using *cross-validation* procedures that involve splitting the examples into three disjoint subsets: a *training set*, a *validation set*, and a *testing set*. The process of optimizing the choice of  $h$  by minimizing  $R_n$  in (2.1) is carried out on the training set over the set of candidate prediction functions, the goal of which is to pinpoint a small subset of viable candidates. The generalized performance of each of these remaining candidates is then estimated using the validation set, and the best performing is chosen as the selected function. The testing set is only used to estimate the generalized performance of this selected function.

Such experimental investigations have, for instance, shown that the *bag of words* approach works very well for text classification [56, 81]. In such an approach, a text document is represented by a feature vector  $x \in \mathbb{R}^d$  whose components are associated with a prescribed set of vocabulary words; i.e., each nonzero component indicates that the associated word appears in the document. (The capabilities of such a representation can also be increased by augmenting the set of words with entries representing short sequences of words.) This encoding scheme leads to very sparse vectors. For instance, the canonical encoding for the standard RCV1 dataset [94] uses a vocabulary of  $d = 47,152$  words to represent news stories that typically contain fewer than 1,000 words. Scaling techniques can be used to give more weight to distinctive words, while scaling to ensure that each document has  $\|x\| = 1$  can be used to compensate for differences in document lengths [94].

Thanks to such a high-dimensional sparse representation of documents, it has been deemed empirically sufficient to consider prediction functions of the form  $h(x; w, \tau) = w^T x - \tau$ . Here,  $w^T x$  is a linear discriminant parameterized by  $w \in \mathbb{R}^d$  and  $\tau \in \mathbb{R}$  is a bias that provides a way to compromise between precision and recall.<sup>1</sup> The accuracy of the predictions could be determined by counting the number of times that  $\text{sign}(h(x; w, \tau))$  matches the correct label, i.e., 1 or  $-1$ . However, while such a prediction function may be appropriate for classifying new documents, formulating an optimization problem around it to choose the parameters  $(w, \tau)$  is impractical in large-scale settings due to the combinatorial structure introduced by the sign function, which is discontinuous. Instead, one typically employs a continuous approximation through a *loss* function that measures a cost for predicting  $h$  when the true label is  $y$ ; e.g., one may choose a *log-loss* function of the form  $\ell(h, y) = \log(1 + \exp(-hy))$ . Moreover, one obtains a class of prediction functions via the addition of a regularization term parameterized by a scalar  $\lambda > 0$ , leading to a convex optimization problem:

$$(2.3) \quad \min_{(w, \tau) \in \mathbb{R}^d \times \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \ell(h(x_i; w, \tau), y_i) + \frac{\lambda}{2} \|w\|_2^2.$$

<sup>1</sup>Precision and recall, defined as the probabilities  $\mathbb{P}[y = 1|h(x) = 1]$  and  $\mathbb{P}[h(x) = 1|y = 1]$ , respectively, are convenient measures of classifier performance when the class of interest represents only a small fraction of all documents.

This problem may be solved multiple times for a given training set with various values of  $\lambda > 0$ , with the ultimate solution  $(w_*, \tau_*)$  being the one that yields the best performance on a validation set.

Many variants of problem (2.3) have also appeared in the literature. For example, the well-known support vector machine problem [38] amounts to using the *hinge loss*  $\ell(h, y) = \max(0, 1 - hy)$ . Another popular feature is to use an  $\ell_1$ -norm regularizer, namely,  $\lambda \|w\|_1$ , which favors sparse solutions and therefore restricts attention to a subset of vocabulary words. Other more sophisticated losses can also target a specific precision/recall trade-off or deal with hierarchies of classes; e.g., see [65, 49]. The choice of loss function is usually guided by experimentation.

In summary, both theoretical arguments [38] and experimental evidence [81] indicate that a carefully selected family of prediction functions and such high-dimensional representations of documents lead to good performance while avoiding *overfitting*—recall (2.2)—of the training set. The use of a simple surrogate, such as (2.3), facilitates experimentation and has been very successful for a variety of problems even beyond text classification. Simple surrogate models are, however, not the most effective in all applications. For example, for certain perceptual tasks, the use of deep neural networks—which lead to large-scale, highly nonlinear, and *nonconvex* optimization problems—has produced great advances that are unmatched by approaches involving simpler, convex models. We discuss such problems next.

**2.2. Perceptual Tasks via Deep Neural Networks.** Just like text classification tasks, perceptual tasks such as speech and image recognition are not well performed in an automated manner using computer programs based on sets of prescribed rules. For instance, the infinite diversity of writing styles easily defeats attempts to concisely specify which pixel combinations represent the digit *four*; see Figure 2.1. One may attempt to design heuristic techniques to handle such eclectic instantiations of the same object, but as previous attempts to design such techniques have often failed, computer vision researchers are increasingly embracing machine learning techniques.



**Fig. 2.1** No known prescribed rules express all pixel combinations that represent the digit four.

During the past five years, spectacular applicative successes on perceptual problems have been achieved by machine learning techniques through the use of *deep neural networks* (DNNs). Although there are many kinds of DNNs, most recent advances have been made using essentially the same types that were popular in the 1990s and almost forgotten in the 2000s [111]. What have made recent successes possible are the availability of much larger datasets and greater computational resources.

Because DNNs were initially inspired by simplified models of biological neurons [134, 135], they are often described with jargon borrowed from neuroscience. For the most part, this jargon turns into a rather effective language for describing a prediction function  $h$  whose value is computed by applying successive transformations to a given input vector  $x_i \in \mathbb{R}^{d_0}$ . These transformations are made in *layers*. For example, a canonical fully connected layer performs the computation

$$(2.4) \quad x_i^{(j)} = s(W_j x_i^{(j-1)} + b_j) \in \mathbb{R}^{d_j},$$

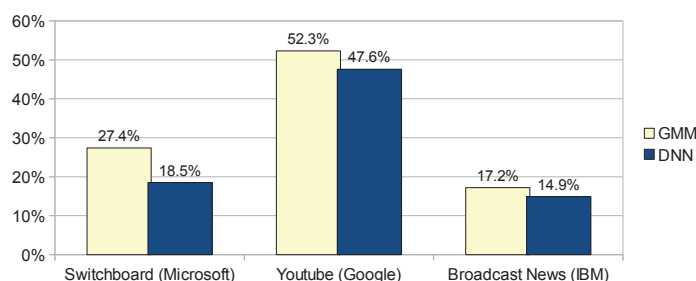
where  $x_i^{(0)} = x_i$ , the matrix  $W_j \in \mathbb{R}^{d_j \times d_{j-1}}$  and vector  $b_j \in \mathbb{R}^{d_j}$  contain the  $j$ th layer parameters, and  $s$  is a component-wise nonlinear *activation* function. Popular choices for the activation function include the sigmoid function  $s(x) = 1/(1 + \exp(-x))$  and the hinge function  $s(x) = \max\{0, x\}$  (often called a rectified linear unit (ReLU) in this context). In this manner, the ultimate output vector  $x_i^{(J)}$  leads to the prediction function value  $h(x_i; w)$ , where the parameter vector  $w$  collects all the parameters  $\{(W_1, b_1), \dots, (W_J, b_J)\}$  of the successive layers.

Similar to (2.3), an optimization problem in this setting involves the collection of a training set  $\{(x_1, y_1) \dots (x_n, y_n)\}$  and the choice of a loss function  $\ell$  leading to

$$(2.5) \quad \min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(h(x_i; w), y_i).$$

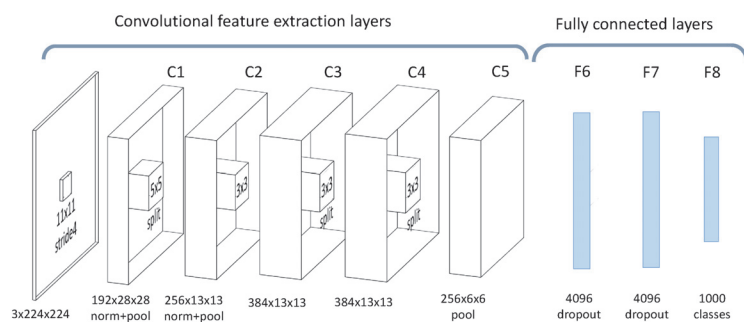
However, in contrast to (2.3), this optimization problem is highly nonlinear and non-convex, making it intractable to solve to global optimality. That being said, machine learning experts have made great strides in the use of DNNs by computing approximate solutions by gradient-based methods. This has been made possible by the conceptually straightforward, yet crucial observation that the gradient of the objective in (2.5) with respect to the parameter vector  $w$  can be computed by the chain rule using algorithmic differentiation [71]. This differentiation technique is known in the machine learning community as *back propagation* [134, 135].

The number of layers in a DNN and the size of each layer are usually determined by performing comparative experiments and evaluating the system performance on a validation set, as in the procedure in section 2.1. A contemporary fully connected neural network for speech recognition typically has five to seven layers. This amounts to tens of millions of parameters to be optimized, the training of which may require up to thousands of hours of speech data (representing hundreds of millions of training examples) and weeks of computation on a supercomputer. Figure 2.2 illustrates the word error rate gains achieved by using DNNs for acoustic modeling in three state-of-the-art speech recognition systems. These gains in accuracy are so significant that DNNs are now used in all the main commercial speech recognition products.

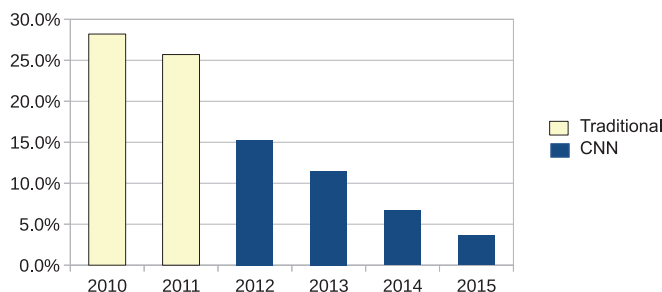


**Fig. 2.2** Word error rates reported by three different research groups on three standard speech recognition benchmarks. For all three groups, deep neural networks (DNNs) significantly outperform the traditional Gaussian mixture models (GMMs) [50]. These experiments were performed between 2010 and 2012 and were instrumental in the recent revival of DNNs.

At the same time, convolutional neural networks (CNNs) have proved to be very effective for computer vision and signal processing tasks [87, 24, 88, 85]. Such a network is composed of convolutional layers, wherein the parameter matrix  $W_j$  is a



**Fig. 2.3** Architecture for image recognition. The 2012 ILSVRC winner consists of eight layers [85]. Each layer performs a linear transformation (specifically, convolutions in layers C1–C5 and matrix multiplication in layers F6–F8) followed by nonlinear transformations (rectification in all layers, contrast normalization in C1–C2, and pooling in C1–C2 and C5). Regularization with dropout noise is used in layers F6–F7.



**Fig. 2.4** Historical top5 error rate of the annual winner of the ImageNet image classification challenge. A convolutional neural network (CNN) achieved a significant performance improvement over all traditional methods in 2012. The following years have cemented CNNs as the current state of the art in visual object recognition [85, 129].

circulant matrix and the input  $x_i^{(j-1)}$  is interpreted as a multichannel image. The product  $W_j x_i^{(j-1)}$  then computes the convolution of the image by a trainable filter while the activation functions—which are piecewise linear functions as opposed to sigmoids—can perform more complex operations that may be interpreted as image rectification, contrast normalization, or subsampling. Figure 2.3 represents the architecture of the winner of the landmark 2012 ImageNet Large Scale Visual Recognition Competition (ILSVRC) [148]. The figure illustrates a CNN with five convolutional layers and three fully connected layers [85]. The input vector represents the pixel values of a  $224 \times 224$  image, while the output scores represent the odds that the image belongs to each of 1,000 categories. This network contains about 60 million parameters, the training of which on a few million labeled images takes a few days on a dual GPU workstation.

Figure 2.4 illustrates the historical error rates of the winner of the 2012 ILSVRC. In this competition, a classification is deemed successful if the correct category appeared among the top five categories returned by the system. The large performance gain achieved in 2012 was confirmed in the following years, and today CNNs are



considered the tool of choice for visual object recognition [129]. They are currently deployed by numerous Internet companies for image search and face recognition.

The successes of DNNs in modern machine learning applications are undeniable. Although the training process requires extreme skill and care—e.g., it is crucial to initialize the optimization process with a good starting point and to monitor its progress while correcting conditioning issues as they appear [89]—the mere fact that one can do anything useful with such large, highly nonlinear, and nonconvex models is remarkable.

**2.3. Formal Machine Learning Procedure.** Through our case studies, we have illustrated how a process of machine learning leads to the selection of a prediction function  $h$  through solving an optimization problem. Moving forward, it is necessary to formalize our presentation by discussing in greater detail the principles behind the selection process, stressing the theoretical importance of *uniform laws of large numbers* as well as the practical importance of *structural risk minimization*.

For simplicity, we continue to focus on the problems that arise in the context of *supervised classification*; i.e., we focus on the optimization of prediction functions for labeling unseen data based on information contained in a set of labeled training data. Such a focus is reasonable as many unsupervised and other learning techniques reduce to optimization problems of comparable form; see, e.g., [155].

**Fundamentals.** Our goal is to determine a prediction function  $h : \mathcal{X} \rightarrow \mathcal{Y}$  from an input space  $\mathcal{X}$  to an output space  $\mathcal{Y}$  such that, given  $x \in \mathcal{X}$ , the value  $h(x)$  offers an accurate prediction about the true output  $y$ . That is, our goal is to choose a prediction function that avoids rote memorization and instead generalizes the concepts that can be learned from a given set of examples. To do this, one should choose the prediction function  $h$  by attempting to minimize a risk measure over an adequately selected family of prediction functions [157], call it  $\mathcal{H}$ .

To formalize this idea, suppose that the examples are sampled from a joint probability distribution function  $P(x, y)$  that simultaneously represents the distribution  $P(x)$  of inputs as well as the conditional probability  $P(y|x)$  of the label  $y$  being appropriate for an input  $x$ . (With this view, one often refers to the examples as *samples*; we use both terms throughout the rest of the paper.) Rather than one that merely minimizes the empirical risk (2.1), one should seek to find  $h$  that yields a small *expected risk* of misclassification *over all possible inputs*, i.e., an  $h$  that minimizes

$$(2.6) \quad R(h) = \mathbb{P}[h(x) \neq y] = \mathbb{E}[\mathbb{1}[h(x) \neq y]],$$

where  $\mathbb{P}[A]$  and  $\mathbb{E}[A]$ , respectively, denote the probability and expected value of  $A$ . Such a framework is *variational* since we are optimizing over a set of functions, and it is *stochastic* since the objective function involves an expectation.

While one may desire to minimize the expected risk (2.6), in practice one must attempt to do so without explicit knowledge of  $P$ . Instead, the only tractable option is to construct a surrogate problem that relies solely on the examples  $\{(x_i, y_i)\}_{i=1}^n$ . Overall, there are two main issues that must be addressed: (i) how to choose the parameterized family of prediction functions  $\mathcal{H}$ , and (ii) how to determine (and find) the particular prediction function  $h \in \mathcal{H}$  that is optimal.

**Choice of Prediction Function Family.** The family of functions  $\mathcal{H}$  should be determined with three *potentially competing* goals in mind. First,  $\mathcal{H}$  should contain prediction functions that are able to achieve a low empirical risk over the training

set, so as to avoid bias or underfitting the data. This can be achieved by selecting a rich family of functions or by using a priori knowledge to select a well-targeted family. Second, the gap between expected risk and empirical risk, namely,  $R(h) - R_n(h)$ , should be small over all  $h \in \mathcal{H}$ . Generally, this gap decreases when one uses more training examples, but, due to potential overfitting, it increases when one uses richer families of functions (see below). This latter fact puts the second goal at odds with the first. Third,  $\mathcal{H}$  should be selected so that one can efficiently solve the corresponding optimization problem, the difficulty of which may increase when one employs a richer family of functions and/or a larger training set.

Our observation about the gap between expected and empirical risk can be understood by recalling certain *laws of large numbers*. For instance, when the expected risk represents a misclassification probability as in (2.6), the Hoeffding inequality [75] guarantees that for any fixed prediction function  $h \in \mathcal{H}$ , the probability of drawing  $n$  examples  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  such that

$$|R(h) - R_n(h)| \leq \sqrt{\frac{1}{2n} \log \left( \frac{2}{\eta} \right)}$$

is greater than  $1 - \eta$ . This bound offers the intuitive explanation that the gap decreases as one uses more training examples. However, this view is insufficient for our purposes since, in the context of machine learning, the prediction function  $h$  is not fixed but depends on the training sample! Rather,  $h$  is the variable over which one is optimizing.

For this reason, one often turns to *uniform laws of large numbers* and the concept of the Vapnik–Chervonenkis (VC) dimension of  $\mathcal{H}$ , a measure of the *capacity* of such a family of functions [157]. For the intuition behind this concept, consider, e.g., a binary classification scheme in  $\mathbb{R}^2$  where one assigns a label of 1 for points above a polynomial and  $-1$  for points below. The set of linear polynomials has a low capacity in the sense that it is only capable of accurately classifying training points that can be separated by a line; e.g., in two variables, a linear classifier has a VC dimension of three. A set of high-degree polynomials, on the other hand, has a high capacity since it can accurately separate training points that are interspersed; the VC dimension of a polynomial of degree  $D$  in  $d$  variables is  $\binom{d+D}{d}$ . That being said, the gap between empirical and expected risk can be larger for a set of high-degree polynomials since the high capacity allows them to overfit a given set of training data.

Mathematically, with the VC dimension measuring capacity, one can establish one of the most important results in learning theory: with  $d_{\mathcal{H}}$  defined as the VC dimension of  $\mathcal{H}$ , one has with probability at least  $1 - \eta$  that

$$(2.7) \quad \sup_{h \in \mathcal{H}} |R(h) - R_n(h)| \leq \mathcal{O} \left( \sqrt{\frac{1}{2n} \log \left( \frac{2}{\eta} \right)} + \frac{d_{\mathcal{H}}}{n} \log \left( \frac{n}{d_{\mathcal{H}}} \right) \right).$$

This bound gives a more accurate picture of the dependence of the gap on the choice of  $\mathcal{H}$ . For example, it shows that for a fixed  $d_{\mathcal{H}}$ , uniform convergence is obtained by increasing the number of training points  $n$ . However, it also shows that, for a fixed  $n$ , the gap can widen for larger  $d_{\mathcal{H}}$ . Indeed, to maintain the same gap, one must increase  $n$  at the same rate if  $d_{\mathcal{H}}$  is increased. The uniform convergence embodied in this result is crucial in machine learning since one wants to ensure that the prediction system performs well with any data provided to it. In section 4.4, we employ a

slight variant of this result to discuss computational trade-offs that arise in large-scale learning.<sup>2</sup>

Interestingly, one quantity that does not enter in (2.7) is the number of parameters that distinguish a particular member function  $h$  of the family  $\mathcal{H}$ . In some settings such as logistic regression, this number is essentially the same as  $d_{\mathcal{H}}$ , which might suggest that the task of optimizing over  $h \in \mathcal{H}$  is more cumbersome as  $d_{\mathcal{H}}$  increases. However, this is not always the case. Certain families of functions are amenable to minimization despite having a very large or even infinite number of parameters [156, section 4.11]. For example, support vector machines [38] were designed to take advantage of this fact [156, Theorem 10.3].

All in all, while bounds such as (2.7) are theoretically interesting and provide useful insight, they are rarely used directly in practice since, as we have suggested in sections 2.1 and 2.2, it is typically easier to estimate the gap between empirical and expected risk with *cross-validation* experiments. We now present ideas underlying a practical framework that respects the trade-offs mentioned above.

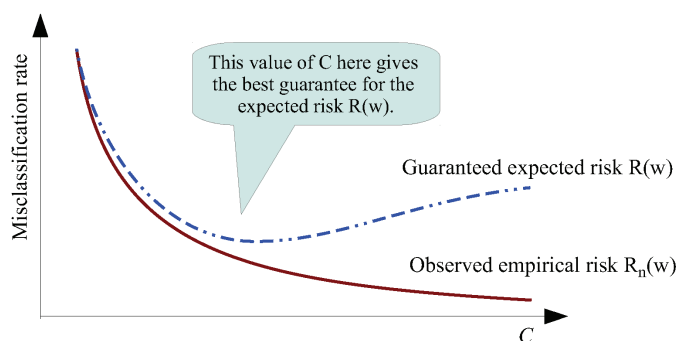
**Structural Risk Minimization.** An approach for choosing a prediction function that has proved to be widely successful in practice is *structural risk minimization* [158, 156]. Rather than choose a generic family of prediction functions—over which it would be difficult both to optimize and to estimate the gap between empirical and expected risks—one chooses a *structure*, i.e., a collection of nested function families. For instance, such a structure can be formed as a collection of subsets of a given family  $\mathcal{H}$  in the following manner: given a preference function  $\Omega$ , choose various values of a *hyperparameter*  $C$ , according to each of which one obtains the subset  $\mathcal{H}_C := \{h \in \mathcal{H} : \Omega(h) \leq C\}$ . Given a fixed number of examples, increasing  $C$  reduces the empirical risk (i.e., the minimum of  $R_n(h)$  over  $h \in \mathcal{H}_C$ ), but, after some point, it typically increases the gap between expected and empirical risks. This phenomenon is illustrated in Figure 2.5.

Other ways to introduce structures are to consider a regularized empirical risk  $R_n(h) + \lambda\Omega(h)$  (an idea introduced in problem (2.3), which may be viewed as the Lagrangian for minimizing  $R_n(h)$  subject to  $\Omega(h) \leq C$ ), to enlarge the dictionary in a bag-of-words representation, to increase the degree of a polynomial model function, or to add to the dimension of an inner layer of a DNN.

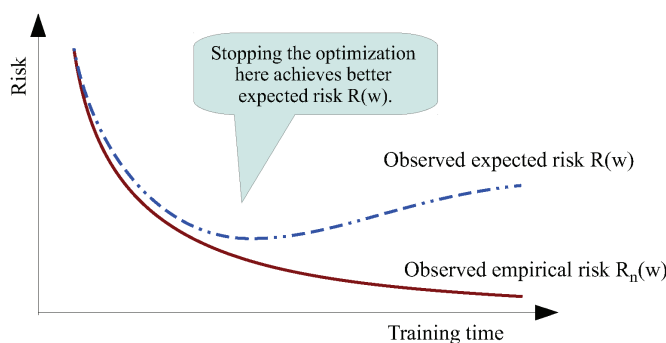
Given such a setup, one can avoid estimating the gap between empirical and expected risk by splitting the available data into subsets: a *training set* used to produce a subset of candidate solutions, a *validation set* used to estimate the expected risk for each such candidate, and a *testing set* used to estimate the expected risk for the candidate that is ultimately chosen. Specifically, over the training set, one minimizes an empirical risk measure  $R_n$  over  $\mathcal{H}_C$  for various values of  $C$ . This results in a handful of candidate functions. The validation set is then used to estimate the expected risk corresponding to each candidate solution, after which one chooses the function yielding the lowest estimated risk value. Assuming a large enough range for  $C$  has been used, one often finds that the best solution does not correspond to the largest value of  $C$  considered; again, see Figure 2.5.

Another, albeit indirect, avenue toward risk minimization is to employ an algorithm for minimizing  $R_n$ , but terminate the algorithm *early*, i.e., before an actual

<sup>2</sup>We also note that considerably better bounds hold when one can collect statistics on actual examples, e.g., by determining gaps dependent on an observed variance of the risk or by considering uniform bounds restricted to families of prediction functions that achieve a risk within a certain threshold of the optimum [55, 102, 26].



**Fig. 2.5** Illustration of structural risk minimization. Given a set of  $n$  examples, a decision function family  $\mathcal{H}$ , and a relative preference function  $\Omega$ , the figure illustrates a typical relationship between the expected and empirical risks corresponding to a prediction function obtained by an optimization algorithm that minimizes an empirical risk  $R_n(h)$  subject to  $\Omega(h) \leq C$ . The optimal empirical risk decreases when  $C$  increases. Meanwhile, the deviation between empirical and expected risk is bounded above by a quantity—which depends on  $\mathcal{H}$  and  $\Omega$ —that increases with  $C$ . While not shown in the figure, the value of  $C$  that offers the best guarantee on the expected risk increases with  $n$ , i.e., the number of examples; recall (2.7).



**Fig. 2.6** Illustration of early stopping. Prematurely stopping the optimization of the empirical risk  $R_n$  often results in a better expected risk  $R$ . In this manner, the stopping time plays a similar role as the hyperparameter  $C$  in the illustration of structural risk minimization in Figure 2.5.

minimizer of  $R_n$  is found. In this manner, the role of the hyperparameter is played by the training time allowed, according to which one typically finds the relationships illustrated in Figure 2.6. Theoretical analyses related to the idea of early stopping are much more challenging than those for other forms of structural risk minimization. However, it is worthwhile to mention these effects since early stopping is a popular technique in practice, and is often *essential* due to computational budget limitations.

Overall, the structural risk minimization principle has proved useful for many applications, and it can be viewed as an alternative to the approach of employing expert human knowledge mentioned in section 2.1. Rather than encoding knowledge as formal classification rules, one encodes it via preferences for certain prediction functions over others, then explores the performance of various prediction functions that have been optimized under the influence of such preferences.

**3. Overview of Optimization Methods.** We now turn our attention to the main focus of our study, namely, numerical algorithms for solving optimization problems that arise in large-scale machine learning. We begin by formalizing our problems of interest, which can be seen as generic statements of problems of the type described in section 2 for minimizing expected and empirical risks. We then provide an overview of two main classes of optimization methods—*stochastic* and *batch*—that can be applied to solve such problems, emphasizing some of the fundamental reasons why stochastic methods have inherent advantages. We close this section with a preview of some of the advanced optimization techniques that are discussed in detail in later sections, which borrow ideas from both stochastic and batch methods.

**3.1. Formal Optimization Problem Statements.** As seen in section 2, optimization problems in machine learning arise through the definition of prediction and loss functions that appear in measures of expected and empirical risk that one aims to minimize. Our discussions revolve around the following definitions.

**Prediction and Loss Functions.** Rather than consider a variational optimization problem over a generic family of prediction functions, we assume that the prediction function  $h$  has a fixed form and is parameterized by a real vector  $w \in \mathbb{R}^d$  over which the optimization is to be performed. Formally, for some given  $h(\cdot; \cdot) : \mathbb{R}^{d_x} \times \mathbb{R}^d \rightarrow \mathbb{R}^{d_y}$ , we consider the family of prediction functions

$$\mathcal{H} := \{h(\cdot; w) : w \in \mathbb{R}^d\}.$$

We aim to find the prediction function in this family that minimizes the losses incurred from inaccurate predictions. For this purpose, we assume a given loss function  $\ell : \mathbb{R}^{d_y} \times \mathbb{R}^{d_y} \rightarrow \mathbb{R}$  to be one that, given an input-output pair  $(x, y)$ , yields the loss  $\ell(h(x; w), y)$  when  $h(x; w)$  and  $y$  are the predicted and true outputs, respectively.

**Expected Risk.** Ideally, the parameter vector  $w$  is chosen to minimize the expected loss that would be incurred from *any* input-output pair. To state this idea formally, we assume that losses are measured with respect to a probability distribution  $P(x, y)$  representing the true relationship between inputs and outputs. That is, we assume that the input-output space  $\mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$  is endowed with  $P : \mathbb{R}^{d_x} \times \mathbb{R}^{d_y} \rightarrow [0, 1]$  and the objective function we wish to minimize is

$$(3.1) \quad R(w) = \int_{\mathbb{R}^{d_x} \times \mathbb{R}^{d_y}} \ell(h(x; w), y) \, dP(x, y) = \mathbb{E}[\ell(h(x; w), y)].$$

We say that  $R : \mathbb{R}^d \rightarrow \mathbb{R}$  yields the *expected risk* (i.e., expected loss) given a parameter vector  $w$  with respect to the probability distribution  $P$ .

**Empirical Risk.** While it may be desirable to minimize (3.1), such a goal is untenable when one does not have complete information about  $P$ . Thus, in practice, one seeks the solution of a problem that involves an estimate of the expected risk  $R$ . In supervised learning, one has access (either all at once or incrementally) to a set of  $n \in \mathbb{N}$  independently drawn input-output samples  $\{(x_i, y_i)\}_{i=1}^n \subseteq \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$ , with which one may define the *empirical risk* function  $R_n : \mathbb{R}^d \rightarrow \mathbb{R}$  by

$$(3.2) \quad R_n(w) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i; w), y_i).$$

Generally speaking, minimization of  $R_n$  may be considered the practical optimization problem of interest. For now, we consider the unregularized measure (3.2), remarking

that the optimization methods that we discuss in the subsequent sections can be applied readily when a smooth regularization term is included. (We leave a discussion of nonsmooth regularizers until section 8.)

Note that, in section 2, the functions  $R$  and  $R_n$  represented *misclassification error*; see (2.1) and (2.6). However, these new definitions of  $R$  and  $R_n$  measure the loss as determined by the function  $\ell$ . We use these latter definitions for the rest of the paper.

**Simplified Notation.** The expressions (3.1) and (3.2) show explicitly how the expected and empirical risks depend on the loss function, sample space, or sample set, etc. However, when discussing optimization methods, we will often employ a simplified notation that also offers some avenues for generalizing certain algorithmic ideas. In particular, let us represent a sample (or set of samples) by a random seed  $\xi$ ; e.g., one may imagine a realization of  $\xi$  as a single sample  $(x, y)$  from  $\mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$ , or a realization of  $\xi$  might be a set of samples  $\{(x_i, y_i)\}_{i \in \mathcal{S}}$ . In addition, let us refer to the loss incurred for a given  $(w, \xi)$  as  $f(w; \xi)$ , i.e.,

(3.3)  $f$  is the composition of the loss function  $\ell$  and the prediction function  $h$ .

In this manner, the expected risk for a given  $w$  is the expected value of this composite function taken with respect to the distribution of  $\xi$ :

$$(3.4) \quad \text{(Expected Risk)} \quad R(w) = \mathbb{E}[f(w; \xi)].$$

In a similar manner, when given a set of realizations  $\{\xi_{[i]}\}_{i=1}^n$  of  $\xi$  corresponding to a sample set  $\{(x_i, y_i)\}_{i=1}^n$ , let us define the loss incurred by the parameter vector  $w$  with respect to the  $i$ th sample as

$$(3.5) \quad f_i(w) := f(w; \xi_{[i]}),$$

and then write the empirical risk as the average of the sample losses:

$$(3.6) \quad \text{(Empirical Risk)} \quad R_n(w) = \frac{1}{n} \sum_{i=1}^n f_i(w).$$

For future reference, we use  $\xi_{[i]}$  to denote the  $i$ th element of a fixed set of realizations of a random variable  $\xi$ , whereas, starting in section 4, we will use  $\xi_k$  to denote the  $k$ th element of a sequence of random variables.

**3.2. Stochastic vs. Batch Optimization Methods.** Let us now introduce some fundamental optimization algorithms for minimizing risk. For the moment, since it is the typical setting in practice, we introduce two algorithm classes in the context of minimizing the empirical risk measure  $R_n$  in (3.6). Note, however, that much of our later discussion will focus on the performance of algorithms when considering the true measure of interest, namely, the expected risk  $R$  in (3.4).

Optimization methods for machine learning fall into two broad categories. We refer to them as *stochastic* and *batch*. The prototypical stochastic optimization method is the stochastic gradient method (SG) [130], which, in the context of minimizing  $R_n$  and with  $w_1 \in \mathbb{R}^d$  given, is defined by

$$(3.7) \quad w_{k+1} \leftarrow w_k - \alpha_k \nabla f_{i_k}(w_k).$$

Here, for all  $k \in \mathbb{N} := \{1, 2, \dots\}$ , the index  $i_k$  (corresponding to the seed  $\xi_{[i_k]}$ , i.e., the sample pair  $(x_{i_k}, y_{i_k})$ ) is chosen *randomly* from  $\{1, \dots, n\}$  and  $\alpha_k$  is a positive

stepsize. Each iteration of this method is thus very cheap, involving only the computation of the gradient  $\nabla f_{i_k}(w_k)$  corresponding to one sample. The method is notable in that the iterate sequence is not determined uniquely by the function  $R_n$ , the starting point  $w_1$ , and the sequence of stepsizes  $\{\alpha_k\}$ , as it would be in a deterministic optimization algorithm. Rather,  $\{w_k\}$  is a stochastic process whose behavior is determined by the random sequence  $\{i_k\}$ . Still, as we shall see in our analysis in section 4, while each direction  $-\nabla f_{i_k}(w_k)$  might not be one of descent from  $w_k$  (in the sense of yielding a negative directional derivative for  $R_n$  from  $w_k$ ), if it is a descent direction *in expectation*, then the sequence  $\{w_k\}$  can be guided toward a minimizer of  $R_n$ .

For many in the optimization research community, a *batch* approach is a more natural and well-known idea. The simplest such method in this class is the steepest descent algorithm—also referred to as the gradient, batch gradient, or full gradient method—which is defined by the iteration

$$(3.8) \quad w_{k+1} \leftarrow w_k - \alpha_k \nabla R_n(w_k) = w_k - \frac{\alpha_k}{n} \sum_{i=1}^n \nabla f_i(w_k).$$

Computing the step  $-\alpha_k \nabla R_n(w_k)$  in such an approach is more expensive than computing the step  $-\alpha_k \nabla f_{i_k}(w_k)$  in SG, though one may expect that a better step is computed when all samples are considered in an iteration.

Stochastic and batch approaches offer different trade-offs in terms of per-iteration costs and expected per-iteration improvement in minimizing empirical risk. Why, then, has SG risen to such prominence in the context of large-scale machine learning? Understanding the reasoning behind this requires careful consideration of the computational trade-offs between stochastic and batch methods, as well as a deeper look into their abilities to guarantee improvement in the underlying expected risk  $R$ . We start to investigate these topics in the next subsection.

We remark in passing that the stochastic and batch approaches mentioned here have analogues in the simulation and stochastic optimization communities, where they are referred to as *stochastic approximation* (SA) and *sample average approximation* (SAA), respectively [63].

**3.3. Motivation for Stochastic Methods.** Before discussing the strengths of stochastic methods such as SG, one should not lose sight of the fact that batch approaches possess some intrinsic advantages. First, when one has reduced the stochastic problem of minimizing the expected risk  $R$  to focus exclusively on the deterministic problem of minimizing the empirical risk  $R_n$ , the use of full gradient information at each iterate opens the door for many deterministic gradient-based optimization methods. That is, in a batch approach, one has at one's disposal the wealth of nonlinear optimization techniques that have been developed over the past decades, including not only the full gradient method (3.8), but also accelerated gradient, conjugate gradient, quasi-Newton, and inexact Newton methods [114]. (See sections 6 and 7 for discussion of these techniques.) Second, due to the sum structure of  $R_n$ , a batch method can easily benefit from parallelization since the bulk of the computation lies in evaluations of  $R_n$  and  $\nabla R_n$ . Calculations of these quantities can even be done in a distributed manner.

Despite these advantages, there are intuitive, practical, and theoretical reasons for following a stochastic approach. Let us motivate them by contrasting the hallmark SG iteration (3.7) with the full batch gradient iteration (3.8).

The paper by Robbins and Monro [130] represents a landmark in the history of numerical optimization methods. Together with the invention of back propagation [134, 135], it also represents one of the most notable developments in the field of machine learning. The SG method was first proposed in [130], not as a gradient method, but as a Markov chain.

Viewed more broadly, the works by Robbins and Monro [130] and Kalman [83] mark the beginning of the field of *stochastic approximation*, which studies the behavior of iterative methods that use noisy signals. The initial focus on optimization led to the study of algorithms that track the solution of the ordinary differential equation  $\dot{w} = -\nabla F(w)$ . Stochastic approximation theory has had a major impact in signal processing and in areas closer to the subject of this paper, such as pattern recognition [4] and neural networks [20].

After receiving his Ph.D. Herbert Robbins became a lecturer at New York University, where he coauthored with Richard Courant the popular book *What is Mathematics?* [39], which is still in print after more than seven decades [40]. Robbins went on to become one of the most prominent mathematicians of the second half of the twentieth century, known for his contributions to probability, algebra, and graph theory.

---

**Inset 3.1** *Herbert Robbins and stochastic approximation.*

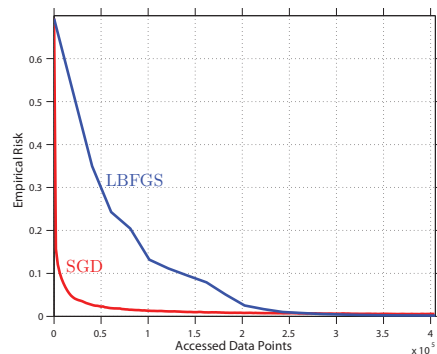
---

**Intuitive Motivation.** On an intuitive level, SG employs information more efficiently than a batch method. To see this, consider a situation in which a training set, call it  $\mathcal{S}$ , consists of ten copies of a set  $\mathcal{S}_{\text{sub}}$ . A minimizer of empirical risk for the larger set  $\mathcal{S}$  is clearly given by a minimizer for the smaller set  $\mathcal{S}_{\text{sub}}$ , but if one were to apply a batch approach to minimize  $R_n$  over  $\mathcal{S}$ , then each iteration would be *ten times* more expensive than if one only had one copy of  $\mathcal{S}_{\text{sub}}$ . On the other hand, SG performs the same computations in both scenarios, in the sense that the stochastic gradient computations involve choosing elements from  $\mathcal{S}_{\text{sub}}$  with the same probabilities. In reality, a training set typically does not consist of exact duplicates of sample data, but in many large-scale applications the data does involve a good deal of (approximate) redundancy. This suggests that using all of the sample data in every optimization iteration is inefficient.

A similar conclusion can be drawn by recalling the discussion in section 2 related to the use of training, validation, and testing sets. If one believes that working with only, say, half of the data in the training set is sufficient to make good predictions on unseen data, then one may argue against working with the entire training set in every optimization iteration. Repeating this argument, working with only a quarter of the training set may be useful at the start, or even with only an eighth of the data, and so on. In this manner, we arrive at the motivation for the idea that working with small samples, at least initially, can be quite appealing.

**Practical Motivation.** The intuitive benefits just described have been observed repeatedly in practice, where one often finds very real advantages of SG in many applications. As an example, Figure 3.1 compares the performance of a batch L-BFGS method [97, 113] (see section 6) and the SG method (3.7) with a constant stepsize (i.e.,  $\alpha_k = \alpha$  for all  $k \in \mathbb{N}$ ) on a binary classification problem using a logistic loss objective function and the data from the RCV1 dataset mentioned in section 2.1. The figure plots the empirical risk  $R_n$  as a function of the number of accesses of a sample from the training set, i.e., the number of evaluations of a sample gradient



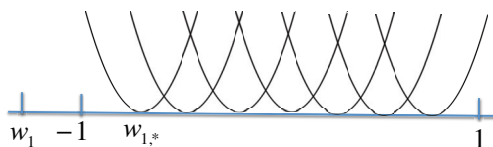


**Fig. 3.1** Empirical risk  $R_n$  as a function of the number of accessed data points (ADPs) for a batch L-BFGS method and the SG method (3.7) on a binary classification problem with a logistic loss objective and the RCV1 dataset. SG was run with a fixed stepsize of  $\alpha = 4$ .

$\nabla f_{i_k}(w_k)$ . Each set of  $n$  consecutive accesses is called an *epoch*. The batch method performs only one step per epoch, while SG performs  $n$  steps per epoch. The plot shows the behavior over the first ten epochs. The advantage of SG is striking and representative of typical behavior in practice. (One should note, however, that to obtain such efficient behavior, it was necessary to run SG repeatedly using different choices for the stepsize  $\alpha$  until a good choice was identified for this particular problem. We discuss theoretical and practical issues related to the choice of stepsize in our analysis in section 4.)

At this point, it is worthwhile to mention that the fast initial improvement achieved by SG, followed by a drastic slowdown after one or two epochs, is common in practice and fairly well understood. An intuitive way to explain this behavior is by considering the following example due to Bertsekas [15].

**EXAMPLE 3.1.** Suppose that each  $f_i$  in (3.6) is a convex quadratic with minimal value at zero and minimizers  $w_{i,*}$  evenly distributed in  $[-1, 1]$  such that the minimizer of  $R_n$  is  $w_* = 0$ ; see Figure 3.2. At  $w_1 \ll -1$ , SG will, with certainty, move to the right (toward  $w_*$ ). Indeed, even if a subsequent iterate lies slightly to the right of the minimizer  $w_{1,*}$  of the “leftmost” quadratic, it is likely (but not certain) that SG will continue moving to the right. However, as iterates near  $w_*$ , the algorithm enters a region of confusion in which there is a significant chance that a step will not move toward  $w_*$ . In this manner, progress will slow significantly. Only with more complete gradient information could the method know with certainty how to move toward  $w_*$ .



**Fig. 3.2** Simple illustration to motivate the fast initial behavior of the SG method for minimizing empirical risk (3.6), where each  $f_i$  is a convex quadratic. This example is adapted from [15].

Despite the issues illustrated by this example, we shall see in section 4 that one can nevertheless ensure convergence by employing a sequence of diminishing stepsizes to overcome any oscillatory behavior of the algorithm.

**Theoretical Motivation.** One can also cite theoretical arguments for a preference for SG over a batch approach. Let us now give a preview of these arguments, which are studied in more depth and further detail in section 4.

- It is well known that a batch approach can minimize  $R_n$  at a fast rate; e.g., if  $R_n$  is strongly convex (see Assumption 4.5) and one applies a batch gradient method, then there exists a constant  $\rho \in (0, 1)$  such that, for all  $k \in \mathbb{N}$ , the *training error* satisfies

$$(3.9) \quad R_n(w_k) - R_n^* \leq \mathcal{O}(\rho^k),$$

where  $R_n^*$  denotes the minimal value of  $R_n$ . The rate of convergence exhibited here is referred to as *R-linear convergence* in the optimization literature [117] and *geometric convergence* in the machine learning research community; we shall simply refer to it as *linear convergence*. From (3.9), one can conclude that, in the worst case, the total number of iterations in which the training error can be above a given  $\epsilon > 0$  is proportional to  $\log(1/\epsilon)$ . This means that, with a per-iteration cost proportional to  $n$  (due to the need to compute  $\nabla R_n(w_k)$  for all  $k \in \mathbb{N}$ ), the total work required to obtain  $\epsilon$ -optimality for a batch gradient method is proportional to  $n \log(1/\epsilon)$ .

- The rate of convergence of a basic stochastic method is slower than for a batch gradient method; e.g., if  $R_n$  is strictly convex and each  $i_k$  is drawn uniformly from  $\{1, \dots, n\}$ , then, for all  $k \in \mathbb{N}$ , the SG iterates defined by (3.7) satisfy the *sublinear convergence* property (see Theorem 4.7)

$$(3.10) \quad \mathbb{E}[R_n(w_k) - R_n^*] = \mathcal{O}(1/k).$$

However, it is crucial to note that *neither the per-iteration cost nor the right-hand side of (3.10) depends on the sample set size  $n$* . This means that the total work required to obtain  $\epsilon$ -optimality for SG is proportional to  $1/\epsilon$ . Admittedly, this can be larger than  $n \log(1/\epsilon)$  for moderate values of  $n$  and  $\epsilon$ , but, as discussed in detail in section 4.4, the comparison favors SG when one moves to the *big data* regime where  $n$  is large and one is merely limited by a computational time budget.

- Another important feature of SG is that, in a stochastic optimization setting, it yields the same convergence rate as in (3.10) for the error in expected risk,  $R - R^*$ , where  $R^*$  is the minimal value of  $R$ . Specifically, by applying the SG iteration (3.7), but with  $\nabla f_{i_k}(w_k)$  replaced by  $\nabla f(w_k; \xi_k)$  with each  $\xi_k$  drawn independently according to the distribution  $P$ , one finds that

$$(3.11) \quad \mathbb{E}[R(w_k) - R^*] = \mathcal{O}(1/k);$$

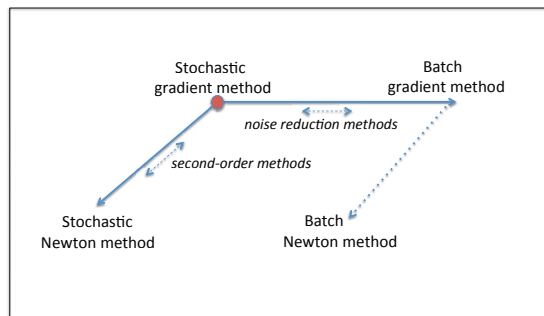
again a sublinear rate, but on the expected risk. Moreover, in this context, a batch approach is not even viable without the ability to compute  $\nabla R$ . Of course, this represents a different setting than one in which only a finite training set is available, but it reveals that if  $n$  is large with respect to  $k$ ,

then the behavior of SG in terms of minimizing the empirical risk  $R_n$  or the expected risk  $R$  is practically indistinguishable up to iteration  $k$ . This property cannot be claimed by a batch method.

In summary, there are intuitive, practical, and theoretical arguments in favor of stochastic over batch approaches in optimization methods for large-scale machine learning. For these reasons, and since SG is used so pervasively by practitioners, we frame our discussions about optimization methods in the context of their relationship with SG. We do not claim, however, that batch methods have no place in practice. For one thing, if Figure 3.1 were to consider a larger number of epochs, then one would see the batch approach eventually overtake the stochastic method and yield a lower training error. This motivates why many recently proposed methods try to combine the best properties of batch and stochastic algorithms. Moreover, the SG iteration is difficult to parallelize and requires excessive communication between nodes in a distributed computing setting, providing further impetus for the design of new and improved optimization algorithms.

**3.4. Beyond SG: Noise Reduction and Second-Order Methods.** Looking forward, one of the main questions being asked by researchers and practitioners alike is: what lies beyond SG that can serve as an efficient, reliable, and easy-to-use optimization method for the kinds of applications discussed in section 2?

To answer this question, we depict in Figure 3.3 methods that aim to improve upon SG as lying on a two-dimensional plane. At the origin of this organizational scheme is SG, representing the base from which all other methods may be compared.



**Fig. 3.3** Schematic of a two-dimensional spectrum of optimization methods for machine learning. The horizontal axis represents methods designed to control stochastic noise; the second axis, methods that deal with ill-conditioning.

From the origin along the horizontal axis, we place methods that are neither purely stochastic nor purely batch, but which attempt to combine the best properties of both approaches. For example, observing the iteration (3.7), one quickly realizes that there is no particular reason to employ information from only one sample point per iteration. Instead, one can employ a *minibatch* approach in which a small subset of samples, call it  $\mathcal{S}_k \subseteq \{1, \dots, n\}$ , is chosen randomly in each iteration, leading to

$$(3.12) \quad w_{k+1} \leftarrow w_k - \frac{\alpha_k}{|\mathcal{S}_k|} \sum_{i \in \mathcal{S}_k} \nabla f_i(w_k).$$

Such an approach falls under the framework set out by Robbins and Monro [130]

and allows some degree of parallelization to be exploited in the computation of minibatch gradients. In addition, one often finds that, due to the reduced variance of the stochastic gradient estimates, the method is easier to tune in terms of choosing the stepsizes  $\{\alpha_k\}$ . Such a minibatch SG method has been widely used in practice.

Along this horizontal axis, one finds other methods as well. In our investigation, we classify two main groups as *dynamic sample size* and *gradient aggregation* methods, both of which aim to improve the rate of convergence from sublinear to linear. These methods do not simply compute minibatches of fixed size, nor do they compute full gradients in every iteration. Instead, they dynamically replace or incorporate new gradient information in order to construct a more reliable step with smaller variance than an SG step. For this reason, we refer to the methods along the horizontal axis as *noise reduction methods*. We will discuss methods of this type in section 5.

Along the second axis in Figure 3.3 are algorithms that, in a broad sense, attempt to overcome the adverse effects of high nonlinearity and ill-conditioning. For such algorithms, we use the term *second-order methods*, which encompasses a variety of strategies; see section 6. We will discuss well-known inexact Newton and quasi-Newton methods, as well as (generalized) Gauss–Newton methods [14, 141], the natural gradient method [5], and scaled gradient iterations [152, 54].

We caution that the schematic representation of methods presented in Figure 3.3 should not be taken too literally since it is not possible to truly organize algorithms so simply, or to include all methods along only two such axes. For example, one could argue that iterate averaging methods do not belong neatly in the category of second-order methods, even though we place them there, and one could argue that gradient methods with momentum [123] or acceleration [107, 108] do belong in this category, even though we discuss them separately in section 7. Nevertheless, Figure 3.3 provides a useful road map as we describe and analyze a large collection of optimization methods of various forms and characteristics. Moreover, our two-dimensional roadmap is useful in that it suggests that optimization methods do not need to exist along the coordinate axes only; e.g., a batch Newton method is placed at the lower-right corner, and one may consider various combinations of second-order and noise reduction schemes.

**4. Analyses of Stochastic Gradient Methods.** In this section, we provide insights into the behavior of an SG method by establishing its convergence properties and worst-case iteration complexity bounds. A preview of such properties was given in (3.10)–(3.11), but now we prove these and other interesting results in detail, all within the context of a generalized SG algorithm. We start by analyzing our SG algorithm when it is invoked to minimize a strongly convex objective function, where it is possible to establish a global rate of convergence to the optimal objective value. This is followed by analyses when our SG algorithm is employed to minimize a generic nonconvex objective. To emphasize the generality of the results proved in this section, we remark that the objective function under consideration could be the expected risk (3.4) or the empirical risk (3.6); i.e., we refer to the objective function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$ , which represents either

$$(4.1) \quad F(w) = \begin{cases} R(w) = \mathbb{E}[f(w; \xi)] \\ \text{or} \\ R_n(w) = \frac{1}{n} \sum_{i=1}^n f_i(w). \end{cases}$$

Our analyses apply equally to both objectives; the only difference lies in the way that one picks the stochastic gradient estimates in the method.<sup>3</sup>

We define our generalized SG method as Algorithm 4.1. The algorithm merely presumes that three computational tools exist: (i) a mechanism for generating a realization of a random variable  $\xi_k$  (with  $\{\xi_k\}$  representing a sequence of jointly independent random variables); (ii) given an iterate  $w_k \in \mathbb{R}^d$  and the realization of  $\xi_k$ , a mechanism for computing a stochastic vector  $g(w_k, \xi_k) \in \mathbb{R}^d$ ; and (iii) given an iteration number  $k \in \mathbb{N}$ , a mechanism for computing a scalar stepsize  $\alpha_k > 0$ .

---

**Algorithm 4.1** Stochastic Gradient (SG) Method
 

---

- 1: Choose an initial iterate  $w_1$ .
  - 2: **for**  $k = 1, 2, \dots$  **do**
  - 3:   Generate a realization of the random variable  $\xi_k$ .
  - 4:   Compute a stochastic vector  $g(w_k, \xi_k)$ .
  - 5:   Choose a stepsize  $\alpha_k > 0$ .
  - 6:   Set the new iterate as  $w_{k+1} \leftarrow w_k - \alpha_k g(w_k, \xi_k)$ .
  - 7: **end for**
- 

The generality of Algorithm 4.1 can be seen in various ways. First, the value of the random variable  $\xi_k$  need only be viewed as a seed for generating a stochastic direction; as such, a realization of it may represent the choice of a single training sample as in the simple SG method stated as (3.7), or may represent a set of samples as in the minibatch SG method (3.12). Second,  $g(w_k, \xi_k)$  could represent a stochastic gradient—i.e., an unbiased estimator of  $\nabla F(w_k)$ , as in the classical method of Robbins and Monro [130]—or it could represent a stochastic Newton or quasi-Newton direction; see section 6. That is, our analyses cover the choices

$$(4.2) \quad g(w_k, \xi_k) = \begin{cases} \nabla f(w_k; \xi_k), \\ \frac{1}{n_k} \sum_{i=1}^{n_k} \nabla f(w_k; \xi_{k,i}), \\ H_k \frac{1}{n_k} \sum_{i=1}^{n_k} \nabla f(w_k; \xi_{k,i}), \end{cases}$$

where, for all  $k \in \mathbb{N}$ , one has flexibility in the choice of minibatch size  $n_k$  and symmetric positive definite scaling matrix  $H_k$ . No matter what choice is made, we shall come to see that all of our theoretical results hold as long as the expected angle between  $g(w_k, \xi_k)$  and  $\nabla F(w_k)$  is sufficiently positive. Third, Algorithm 4.1 allows various choices of the stepsize sequence  $\{\alpha_k\}$ . Our analyses focus on two choices, one involving a fixed stepsize and one involving diminishing stepsizes, as both are interesting in

---

<sup>3</sup>Picking samples uniformly from a finite training set, replacing them in the set for each iteration, corresponds to sampling from a discrete distribution giving equal weight to every sample. In this case, the SG algorithm in this section optimizes the empirical risk  $F = R_n$ . Alternatively, picking samples in each iteration according to the distribution  $P$ , the SG algorithm optimizes the expected risk  $F = R$ . One could also imagine picking samples *without replacement* until one exhausts a finite training set. In this case, the SG algorithm here can be viewed as optimizing either  $R_n$  or  $R$ , but only until the training set is exhausted. After that point, our analysis no longer apply. Generally speaking, the analysis of such *incremental algorithms* often requires specialized techniques [15, 72].

theory and in practice. Finally, we note that Algorithm 4.1 also covers active learning techniques in which the iterate  $w_k$  influences the sample selection.<sup>4</sup>

Notwithstanding all of this generality, we henceforth refer to Algorithm 4.1 as *SG*. The particular instance (3.7) will be referred to as *simple SG* or *basic SG*, whereas the instance (3.12) will be referred to as *minibatch SG*.

Beyond our convergence and complexity analyses, a complete appreciation for the properties of SG is not possible without highlighting its theoretical advantages over batch methods in terms of computational complexity. Thus, we include in section 4.4 a discussion of the trade-offs between rate of convergence and computational effort among prototypical stochastic and batch methods for large-scale learning.

**4.1. Two Fundamental Lemmas.** Our approach for establishing convergence guarantees for SG is built upon an assumption of smoothness of the objective function. (Alternative foundations are possible; see Appendix A.) This, and an assumption about the first and second moments of the stochastic vectors  $\{g(w_k, \xi_k)\}$ , lead to two fundamental lemmas from which all of our results will be derived.

Our first assumption is formally stated as follows. Recall that, as already mentioned in (4.1),  $F$  can represent either expected or empirical risk.

**ASSUMPTION 4.1** (Lipschitz-continuous objective gradients). *The objective function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  is continuously differentiable and the gradient function of  $F$ , namely,  $\nabla F : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , is Lipschitz continuous with Lipschitz constant  $L > 0$ , i.e.,*

$$\|\nabla F(w) - \nabla F(\bar{w})\|_2 \leq L\|w - \bar{w}\|_2 \text{ for all } \{w, \bar{w}\} \subset \mathbb{R}^d.$$

Intuitively, Assumption 4.1 ensures that the gradient of  $F$  does not change arbitrarily quickly with respect to the parameter vector. Such an assumption is essential for convergence analyses of most gradient-based methods; without it, the gradient would not provide a good indicator for how far to move to decrease  $F$ . An important consequence of Assumption 4.1 is that

$$(4.3) \quad F(w) \leq F(\bar{w}) + \nabla F(\bar{w})^T(w - \bar{w}) + \frac{1}{2}L\|w - \bar{w}\|_2^2 \text{ for all } \{w, \bar{w}\} \subset \mathbb{R}^d.$$

This inequality is proved in Appendix B, but note that it also follows immediately if  $F$  is twice continuously differentiable and the Hessian function  $\nabla^2 F : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$  satisfies  $\|\nabla^2 F(w)\|_2 \leq L$  for all  $w \in \mathbb{R}^d$ .

Under Assumption 4.1 alone, we obtain the following lemma. In the result, we use  $\mathbb{E}_{\xi_k}[\cdot]$  to denote an expected value taken with respect to the distribution of the random variable  $\xi_k$  given  $w_k$ . Therefore,  $\mathbb{E}_{\xi_k}[F(w_{k+1})]$  is a meaningful quantity since  $w_{k+1}$  depends on  $\xi_k$  through the update in step 6 of Algorithm 4.1.

**LEMMA 4.2.** *Under Assumption 4.1, the iterates of SG (Algorithm 4.1) satisfy the following inequality for all  $k \in \mathbb{N}$ :*

$$(4.4) \quad \mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k) \leq -\alpha_k \nabla F(w_k)^T \mathbb{E}_{\xi_k}[g(w_k, \xi_k)] + \frac{1}{2} \alpha_k^2 L \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2].$$

<sup>4</sup>We have assumed that the elements of the random variable sequence  $\{\xi_k\}$  are independent in order to avoid requiring certain machinery from the analysis of stochastic processes. Viewing  $\xi_k$  as a seed instead of a sample during iteration  $k$  makes this restriction minor. However, it is worthwhile to mention that all of the results in this section still hold if, instead,  $\{\xi_k\}$  forms an adapted (nonanticipating) stochastic process and expectations taken with respect to  $\xi_k$  are replaced by expectations taken with respect to the conditional distribution of  $\xi_k$  given  $\{\xi_1, \dots, \xi_{k-1}\}$ .

*Proof.* By Assumption 4.1, the iterates generated by SG satisfy

$$\begin{aligned} F(w_{k+1}) - F(w_k) &\leq \nabla F(w_k)^T (w_{k+1} - w_k) + \frac{1}{2} L \|w_{k+1} - w_k\|_2^2 \\ &\leq -\alpha_k \nabla F(w_k)^T g(w_k, \xi_k) + \frac{1}{2} \alpha_k^2 L \|g(w_k, \xi_k)\|_2^2. \end{aligned}$$

Taking expectations in these inequalities with respect to the distribution of  $\xi_k$ , and noting that  $w_{k+1}$ —but not  $w_k$ —depends on  $\xi_k$ , we obtain the desired bound.  $\square$

This lemma shows that, regardless of how SG arrived at  $w_k$ , the expected decrease in the objective function yielded by the  $k$ th step is bounded above by a quantity involving (i) the expected directional derivative of  $F$  at  $w_k$  along  $-g(w_k, \xi_k)$ , and (ii) the second moment of  $g(w_k, \xi_k)$ . For example, if  $g(w_k, \xi_k)$  is an unbiased estimate of  $\nabla F(w_k)$ , then it follows from Lemma 4.2 that

$$(4.5) \quad \mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k) \leq -\alpha_k \|\nabla F(w_k)\|_2^2 + \frac{1}{2} \alpha_k^2 L \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2].$$

We shall see that convergence of SG is guaranteed as long as the stochastic directions and stepsizes are chosen such that the right-hand side of (4.4) is bounded above by a *deterministic* quantity that asymptotically ensures sufficient descent in  $F$ . One can ensure this in part by stating additional requirements on the first and second moments of the stochastic directions  $\{g(w_k, \xi_k)\}$ . In particular, in order to limit the harmful effect of the last term in (4.5), we restrict the variance of  $g(w_k, \xi_k)$ , i.e.,

$$(4.6) \quad \mathbb{V}_{\xi_k}[g(w_k, \xi_k)] := \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2] - \|\mathbb{E}_{\xi_k}[g(w_k, \xi_k)]\|_2^2.$$

ASSUMPTION 4.3 (first and second moment limits). *The objective function and SG (Algorithm 4.1) satisfy the following conditions:*

- (a) *The sequence of iterates  $\{w_k\}$  is contained in an open set over which  $F$  is bounded below by a scalar  $F_{\inf}$ .*
- (b) *There exist scalars  $\mu_G \geq \mu > 0$  such that, for all  $k \in \mathbb{N}$ ,*

$$(4.7a) \quad \nabla F(w_k)^T \mathbb{E}_{\xi_k}[g(w_k, \xi_k)] \geq \mu \|\nabla F(w_k)\|_2^2 \quad \text{and}$$

$$(4.7b) \quad \|\mathbb{E}_{\xi_k}[g(w_k, \xi_k)]\|_2 \leq \mu_G \|\nabla F(w_k)\|_2.$$

- (c) *There exist scalars  $M \geq 0$  and  $M_V \geq 0$  such that, for all  $k \in \mathbb{N}$ ,*

$$(4.8) \quad \mathbb{V}_{\xi_k}[g(w_k, \xi_k)] \leq M + M_V \|\nabla F(w_k)\|_2^2.$$

The first condition, Assumption 4.3(a), merely requires the objective function to be bounded below over the region explored by the algorithm. The second requirement, Assumption 4.3(b), states that, in expectation, the vector  $-g(w_k, \xi_k)$  is a direction of sufficient descent for  $F$  from  $w_k$  with a norm comparable to the norm of the gradient. The properties in this requirement hold immediately with  $\mu_G = \mu = 1$  if  $g(w_k, \xi_k)$  is an unbiased estimate of  $\nabla F(w_k)$ , and are maintained if such an unbiased estimate is multiplied by a positive definite matrix  $H_k$  that is conditionally uncorrelated with  $g(w_k, \xi_k)$  given  $w_k$  and whose eigenvalues lie in a fixed positive interval for all  $k \in \mathbb{N}$ . The third requirement, Assumption 4.3(c), states that the variance of  $g(w_k, \xi_k)$  is restricted, but in a relatively minor manner. For example, if  $F$  is a convex quadratic function, then the variance is allowed to be nonzero at any stationary point for  $F$  and is allowed to grow quadratically in any direction.

Taken together, Assumption 4.3, combined with the definition (4.6), requires that the second moment of  $g(w_k, \xi_k)$  satisfies

$$(4.9) \quad \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2] \leq M + M_G \|\nabla F(w_k)\|_2^2 \quad \text{with} \quad M_G := M_V + \mu_G^2 \geq \mu^2 > 0.$$

In fact, all of our analyses in this section would hold if this bound on the second moment were to be assumed directly. (We have stated Assumption 4.3 in the form above merely to facilitate our discussion in section 5.)

The following lemma builds on Lemma 4.2 under the additional conditions now set forth in Assumption 4.3.

LEMMA 4.4. *Under Assumptions 4.1 and 4.3, the iterates of SG (Algorithm 4.1) satisfy the following inequalities for all  $k \in \mathbb{N}$ :*

$$(4.10a) \quad \mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k) \leq -\mu\alpha_k \|\nabla F(w_k)\|_2^2 + \frac{1}{2}\alpha_k^2 L \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2]$$

$$(4.10b) \quad \leq -(\mu - \frac{1}{2}\alpha_k L M_G) \alpha_k \|\nabla F(w_k)\|_2^2 + \frac{1}{2}\alpha_k^2 L M.$$

*Proof.* By Lemma 4.2 and (4.7a), it follows that

$$\begin{aligned} \mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k) &\leq -\alpha_k \nabla F(w_k)^T \mathbb{E}_{\xi_k}[g(w_k, \xi_k)] + \frac{1}{2}\alpha_k^2 L \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2] \\ &\leq -\mu\alpha_k \|\nabla F(w_k)\|_2^2 + \frac{1}{2}\alpha_k^2 L \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2], \end{aligned}$$

which is (4.10a). Assumption 4.3, giving (4.9), then yields (4.10b).  $\square$

As mentioned, this lemma reveals that regardless of how the method arrived at the iterate  $w_k$ , the optimization process continues in a *Markovian* manner in the sense that  $w_{k+1}$  is a random variable that depends only on the iterate  $w_k$ , the seed  $\xi_k$ , and the stepsize  $\alpha_k$  and *not on any past iterates*. This can be seen in the fact that the difference  $\mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k)$  is bounded above by a deterministic quantity. Note also that the first term in (4.10b) is strictly negative for small  $\alpha_k$  and suggests a decrease in the objective function by a magnitude proportional to  $\|\nabla F(w_k)\|_2^2$ . However, the second term in (4.10b) could be large enough to allow the objective value to increase. Balancing these terms is critical in the design of SG methods.

**4.2. SG for Strongly Convex Objectives.** The most benign setting for analyzing the SG method is in the context of minimizing a strongly convex objective function. For the reasons described in Inset 4.1, when not considering a generic nonconvex objective  $F$ , we focus on the strongly convex case and only briefly mention the (not strongly) convex case on certain occasions.

We formalize a strong convexity assumption as follows.

ASSUMPTION 4.5 (strong convexity). *The objective function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  is strongly convex in that there exists a constant  $c > 0$  such that*

$$(4.11) \quad F(\bar{w}) \geq F(w) + \nabla F(w)^T (\bar{w} - w) + \frac{1}{2}c \|\bar{w} - w\|_2^2 \quad \text{for all } (\bar{w}, w) \in \mathbb{R}^d \times \mathbb{R}^d.$$

Hence,  $F$  has a unique minimizer, denoted as  $w_* \in \mathbb{R}^d$  with  $F_* := F(w_*)$ .

A useful fact from convex analysis (proved in Appendix B) is that, under Assumption 4.5, one can bound the optimality gap at a given point in terms of the squared  $\ell_2$ -norm of the gradient of the objective at that point:

$$(4.12) \quad 2c(F(w) - F_*) \leq \|\nabla F(w)\|_2^2 \quad \text{for all } w \in \mathbb{R}^d.$$

We use this inequality in several proofs. We also observe that, from (4.3) and (4.11), the constants in Assumptions 4.1 and 4.5 must satisfy  $c \leq L$ .

We now state our first convergence theorem for SG, describing its behavior when minimizing a strongly convex objective function when employing a fixed stepsize. In this case, it will not be possible to prove convergence to the solution, but only to



All of the convergence rate and complexity results presented in this paper relate to the minimization of *strongly convex* functions. This is in contrast with a large portion of the literature on optimization methods for machine learning, in which much effort is made to strengthen convergence guarantees for methods applied to functions that are convex, but not strongly convex. We have made this choice for a few reasons. First, it leads to a focus on results that are relevant to actual machine learning practice, since in many situations when a convex model is employed—such as in logistic regression—it is often regularized by a strongly convex function to facilitate the solution process. Second, there exists a variety of situations in which the objective function is not globally (strongly) convex, but is so in the neighborhood of local minimizers, meaning that our results can represent the behavior of the algorithm in such regions of the search space. Third, one can argue that related results when minimizing non-strongly-convex models can be derived as extensions of the results presented here [3], making our analyses a starting point for deriving a more general theory.

We have also taken a pragmatic approach in the types of convergence guarantees that we provide. In particular, in our analyses, we focus on results that reveal the properties of SG iterates *in expectation*. The stochastic approximation literature, on the other hand, often relies on martingale techniques to establish *almost sure convergence* [66, 131] under the same assumptions [21]. For our purposes, we omit these complications since, in our view, they do not provide significant additional insights into the forces driving convergence of the method.

**Inset 4.1** *Perspectives on SG analyses.*

a neighborhood of the optimal value. (Intuitively, this limitation should be clear from (4.10b) since the first term on the right-hand side decreases in magnitude as the solution is approached—i.e., as  $\nabla F(w_k)$  tends to zero—but the last term remains constant. Thus, after some point, a reduction in the objective cannot be expected.) We use  $\mathbb{E}[\cdot]$  to denote an expected value taken with respect to the joint distribution of all random variables. For example, since  $w_k$  is completely determined by the realizations of the independent random variables  $\{\xi_1, \xi_2, \dots, \xi_{k-1}\}$ , the *total expectation* of  $F(w_k)$  for any  $k \in \mathbb{N}$  can be taken as

$$\mathbb{E}[F(w_k)] = \mathbb{E}_{\xi_1} \mathbb{E}_{\xi_2} \dots \mathbb{E}_{\xi_{k-1}} [F(w_k)].$$

The theorem shows that if the stepsize is not too large, then, in expectation, the sequence of function values  $\{F(w_k)\}$  converges near the optimal value.

**THEOREM 4.6** (strongly convex objective, fixed stepsize). *Under Assumptions 4.1, 4.3, and 4.5 (with  $F_{\inf} = F_*$ ), suppose that the SG method (Algorithm 4.1) is run with a fixed stepsize,  $\alpha_k = \bar{\alpha}$  for all  $k \in \mathbb{N}$ , satisfying*

$$(4.13) \quad 0 < \bar{\alpha} \leq \frac{\mu}{LM_G}.$$

*Then the expected optimality gap satisfies the following inequality for all  $k \in \mathbb{N}$ :*

$$(4.14) \quad \mathbb{E}[F(w_k) - F_*] \leq \frac{\bar{\alpha}LM}{2c\mu} + (1 - \bar{\alpha}c\mu)^{k-1} \left( F(w_1) - F_* - \frac{\bar{\alpha}LM}{2c\mu} \right) \\ \xrightarrow{k \rightarrow \infty} \frac{\bar{\alpha}LM}{2c\mu}.$$

*Proof.* Using Lemma 4.4 with (4.13) and (4.12), we have for all  $k \in \mathbb{N}$  that

$$\begin{aligned}\mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k) &\leq -(\mu - \tfrac{1}{2}\bar{\alpha}LM_G)\bar{\alpha}\|\nabla F(w_k)\|_2^2 + \tfrac{1}{2}\bar{\alpha}^2LM \\ &\leq -\tfrac{1}{2}\bar{\alpha}\mu\|\nabla F(w_k)\|_2^2 + \tfrac{1}{2}\bar{\alpha}^2LM \\ &\leq -\bar{\alpha}c\mu(F(w_k) - F_*) + \tfrac{1}{2}\bar{\alpha}^2LM.\end{aligned}$$

Subtracting  $F_*$  from both sides, taking total expectations, and rearranging, this yields

$$\mathbb{E}[F(w_{k+1}) - F_*] \leq (1 - \bar{\alpha}c\mu)\mathbb{E}[F(w_k) - F_*] + \tfrac{1}{2}\bar{\alpha}^2LM.$$

Subtracting the constant  $\bar{\alpha}LM/(2c\mu)$  from both sides, one obtains

$$\begin{aligned}\mathbb{E}[F(w_{k+1}) - F_*] - \frac{\bar{\alpha}LM}{2c\mu} &\leq (1 - \bar{\alpha}c\mu)\mathbb{E}[F(w_k) - F_*] + \frac{\bar{\alpha}^2LM}{2} - \frac{\bar{\alpha}LM}{2c\mu} \\ (4.15) \qquad \qquad \qquad &= (1 - \bar{\alpha}c\mu)\left(\mathbb{E}[F(w_k) - F_*] - \frac{\bar{\alpha}LM}{2c\mu}\right).\end{aligned}$$

Observe that (4.15) is a contraction inequality since, by (4.13) and (4.9),

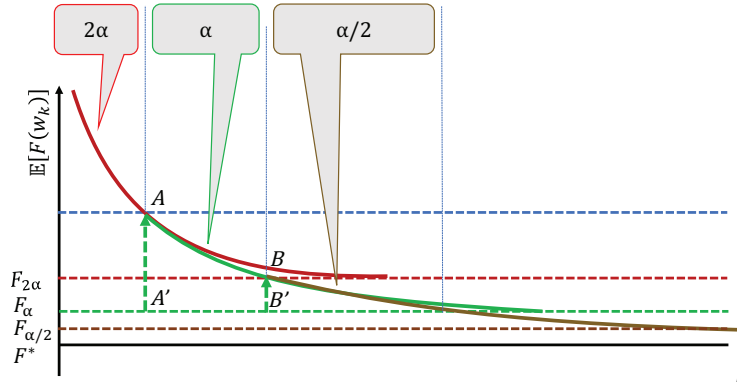
$$(4.16) \qquad \qquad \qquad 0 < \bar{\alpha}c\mu \leq \frac{c\mu^2}{LM_G} \leq \frac{c\mu^2}{L\mu^2} = \frac{c}{L} \leq 1.$$

The result thus follows by applying (4.15) repeatedly through iteration  $k \in \mathbb{N}$ .  $\square$

If  $g(w_k, \xi_k)$  is an unbiased estimate of  $\nabla F(w_k)$ , then  $\mu = 1$ , and if there is no noise in  $g(w_k, \xi_k)$ , then we may presume that  $M_G = 1$  (due to (4.9)). In this case, (4.13) reduces to  $\bar{\alpha} \in (0, 1/L]$ , a classical stepsize requirement of interest for a steepest descent method.

Theorem 4.6 illustrates the interplay between the stepsizes and bound on the variance of the stochastic directions. If there were no noise in the gradient computation or if noise were to decay with  $\|\nabla F(w_k)\|_2^2$  (i.e., if  $M = 0$  in (4.8) and (4.9)), then one could obtain linear convergence to the optimal value. This is a standard result for the full gradient method with a sufficiently small positive stepsize. On the other hand, when the gradient computation is noisy, one clearly loses this property. One can still use a fixed stepsize and be sure that the expected objective values will converge linearly to a neighborhood of the optimal value, but, after some point, the noise in the gradient estimates prevent further progress; recall Example 3.1. It is apparent from (4.14) that selecting a smaller stepsize worsens the contraction constant in the convergence rate, but allows one to arrive closer to the optimal value.

These observations provide a foundation for a strategy often employed in practice in which SG is run with a fixed stepsize, and, if progress appears to stall, a smaller stepsize is selected and the process is repeated. A straightforward instance of such an approach can be motivated with the following sketch. Suppose that  $\alpha_1 \in (0, \frac{\mu}{LM_G}]$  is chosen as in (4.13) and the SG method is run with this stepsize from iteration  $k_1 = 1$  until iteration  $k_2$ , where  $w_{k_2}$  is the first iterate at which the expected suboptimality gap is smaller than twice the asymptotic value in (4.14), i.e.,  $\mathbb{E}[F(w_{k_2}) - F_*] \leq 2F_{\alpha_1}$ , where  $F_{\alpha} := \frac{\alpha LM}{2c\mu}$ . Suppose further that, at this point, the stepsize is halved and the process is repeated; see Figure 4.1. This leads to the stepsize schedule  $\{\alpha_{r+1}\} = \{\alpha_1 2^{-r}\}$ , index sequence  $\{k_r\}$ , and bound sequence  $\{F_{\alpha_r}\} = \{\frac{\alpha_r LM}{2c\mu}\} \searrow 0$  such that,



**Fig. 4.1** Depiction of the strategy of halving the stepsize  $\alpha$  when the expected suboptimality gap is smaller than twice the asymptotic limit  $F_\alpha$ . In the figure, the segment  $B-B'$  has one third of the length of  $A-A'$ . This is the amount of decrease that must be made in the exponential term in (4.14) by raising the contraction factor to the power of the number of steps during which one maintains a given constant stepsize; see (4.18). Since the contraction factor is  $(1 - \alpha c \mu)$ , the number of steps must be proportional to  $\alpha$ . Therefore, whenever the stepsize is halved, one must maintain it twice as long. Overall, doubling the number of iterations halves the suboptimality gap each time, yielding an effective rate of  $\mathcal{O}(1/k)$ .

for all  $r \in \{2, 3, \dots\}$ ,

$$(4.17) \quad \mathbb{E}[F(w_{k_{r+1}}) - F_*] \leq 2F_{\alpha_r}, \quad \text{where} \quad \mathbb{E}[F(w_{k_r}) - F_*] \approx 2F_{\alpha_{r-1}} = 4F_{\alpha_r}.$$

In this manner, the expected suboptimality gap converges to zero.

However, this does *not* occur by halving the stepsize in every iteration, but only once the gap itself has been cut in half from a previous threshold. To see what might be the appropriate *effective rate* of stepsize decrease, we may invoke Theorem 4.6, from which it follows that to achieve the first bound in (4.17) one needs

$$(4.18) \quad \begin{aligned} (1 - \alpha_r c \mu)^{(k_{r+1} - k_r)} (4F_{\alpha_r} - F_{\alpha_r}) &\leq F_{\alpha_r} \\ \implies k_{r+1} - k_r &\geq \frac{\log(1/3)}{\log(1 - \alpha_r c \mu)} \approx \frac{\log(3)}{\alpha_r c \mu} = \mathcal{O}(2^r). \end{aligned}$$

In other words, each time the stepsize is cut in half, double the number of iterations are required. This is a *sublinear* rate of stepsize decrease—e.g., if  $\{k_r\} = \{2^{r-1}\}$ , then  $\alpha_k = \alpha_1/k$  for all  $k \in \{2^r\}$ —which, from  $\{F_{\alpha_r}\} = \{\frac{\alpha_r LM}{2c\mu}\}$  and (4.17), means that a sublinear convergence rate of the suboptimality gap is achieved.

In fact, these conclusions can be obtained in a more rigorous manner that also allows more flexibility in the choice of stepsize sequence. The following result harks back to the seminal work of Robbins and Monro [130], where the stepsize requirement takes the form

$$(4.19) \quad \sum_{k=1}^{\infty} \alpha_k = \infty \quad \text{and} \quad \sum_{k=1}^{\infty} \alpha_k^2 < \infty.$$

**THEOREM 4.7** (strongly convex objective, diminishing stepsizes). *Under Assumptions 4.1, 4.3, and 4.5 (with  $F_{\inf} = F_*$ ), suppose that the SG method (Algorithm 4.1)*

is run with a stepsize sequence such that, for all  $k \in \mathbb{N}$ ,

$$(4.20) \quad \alpha_k = \frac{\beta}{\gamma + k} \text{ for some } \beta > \frac{1}{c\mu} \text{ and } \gamma > 0 \text{ such that } \alpha_1 \leq \frac{\mu}{LM_G}.$$

Then, for all  $k \in \mathbb{N}$ , the expected optimality gap satisfies

$$(4.21) \quad \mathbb{E}[F(w_k) - F_*] \leq \frac{\nu}{\gamma + k},$$

where

$$(4.22) \quad \nu := \max \left\{ \frac{\beta^2 LM}{2(\beta c\mu - 1)}, (\gamma + 1)(F(w_1) - F_*) \right\}.$$

*Proof.* By (4.20), the inequality  $\alpha_k LM_G \leq \alpha_1 LM_G \leq \mu$  holds for all  $k \in \mathbb{N}$ . Hence, along with Lemma 4.4 and (4.12), one has for all  $k \in \mathbb{N}$  that

$$\begin{aligned} \mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k) &\leq -(\mu - \tfrac{1}{2}\alpha_k LM_G)\alpha_k \|\nabla F(w_k)\|_2^2 + \tfrac{1}{2}\alpha_k^2 LM \\ &\leq -\tfrac{1}{2}\alpha_k \mu \|\nabla F(w_k)\|_2^2 + \tfrac{1}{2}\alpha_k^2 LM \\ &\leq -\alpha_k c\mu (F(w_k) - F(w_*)) + \tfrac{1}{2}\alpha_k^2 LM. \end{aligned}$$

Subtracting  $F_*$  from both sides, taking total expectations, and rearranging, this yields

$$(4.23) \quad \mathbb{E}[F(w_{k+1}) - F_*] \leq (1 - \alpha_k c\mu) \mathbb{E}[F(w_k) - F_*] + \tfrac{1}{2}\alpha_k^2 LM.$$

We now prove (4.21) by induction. First, the definition of  $\nu$  ensures that it holds for  $k = 1$ . Then, assuming (4.21) holds for some  $k \geq 1$ , it follows from (4.23) that

$$\begin{aligned} \mathbb{E}[F(w_{k+1}) - F_*] &\leq \left(1 - \frac{\beta c\mu}{\hat{k}}\right) \frac{\nu}{\hat{k}} + \frac{\beta^2 LM}{2\hat{k}^2} \quad (\text{with } \hat{k} := \gamma + k) \\ &= \left(\frac{\hat{k} - \beta c\mu}{\hat{k}^2}\right) \nu + \frac{\beta^2 LM}{2\hat{k}^2} \\ &= \left(\frac{\hat{k} - 1}{\hat{k}^2}\right) \nu - \underbrace{\left(\frac{\beta c\mu - 1}{\hat{k}^2}\right) \nu + \frac{\beta^2 LM}{2\hat{k}^2}}_{\text{nonpositive by the definition of } \nu} \leq \frac{\nu}{\hat{k} + 1}, \end{aligned}$$

where the last inequality follows because  $\hat{k}^2 \geq (\hat{k} + 1)(\hat{k} - 1)$ .  $\square$

Let us now remark on what can be learned from Theorems 4.6 and 4.7.

**Role of Strong Convexity.** Observe the crucial role played by the strong convexity parameter  $c > 0$ , the positivity of which is needed to argue that (4.15) and (4.23) contract the expected optimality gap. However, the strong convexity constant impacts the stepsizes in different ways in Theorems 4.6 and 4.7. In the case of constant stepsizes, the possible values of  $\bar{\alpha}$  are constrained by the upper bound (4.13) that does not depend on  $c$ . In the case of diminishing stepsizes, the initial stepsize  $\alpha_1$  is subject to the same upper bound (4.20), but the stepsize parameter  $\beta$  must be larger than  $1/(c\mu)$ . This additional requirement is critical to ensure the  $\mathcal{O}(1/k)$  convergence rate. How critical? Consider, e.g., [105], in which the authors provide a simple example (with unbiased gradient estimates and  $\mu = 1$ ) involving the minimization of a deterministic quadratic function *with only one optimization variable* in which  $c$  is overestimated, which results in  $\beta$  being underestimated. In the example, even after  $10^9$  iterations, the distance to the solution remains greater than  $10^{-2}$ .

**Role of the Initial Point.** Also observe the role played by the initial point, which determines the initial optimality gap, namely,  $F(w_1) - F_*$ . When using a fixed stepsize, the initial gap appears with an exponentially decreasing factor; see (4.14). In the case of diminishing stepsizes, the gap appears prominently in the second term defining  $\nu$  in (4.22). However, with an appropriate initialization phase, one can easily diminish the role played by this term.<sup>5</sup> For example, suppose that one begins by running SG with a fixed stepsize  $\bar{\alpha}$  until one (approximately) obtains a point, call it  $w_1$ , with  $F(w_1) - F_* \leq \bar{\alpha}LM/(2c\mu)$ . A guarantee for this bound can be argued from (4.14). Starting here with  $\alpha_1 = \bar{\alpha}$ , the choices for  $\beta$  and  $\gamma$  in Theorem 4.7 yield

$$(\gamma + 1)\mathbb{E}[F(w_1) - F_*] \leq \beta\alpha_1^{-1} \frac{\alpha_1 LM}{2c\mu} = \frac{\beta LM}{2c\mu} < \frac{\beta^2 LM}{2(\beta c\mu - 1)},$$

meaning that the value for  $\nu$  is dominated by the first term in (4.22).

On a related note, we claim that for practical purposes the initial stepsize should be chosen as large as allowed, i.e.,  $\alpha_1 = \mu/(LM_G)$ . Given this choice of  $\alpha_1$ , the best asymptotic regime with decreasing stepsizes (4.21) is achieved by making  $\nu$  as small as possible. Since we have argued that only the first term matters in the definition of  $\nu$ , this leads to choosing  $\beta = 2/(c\mu)$ . Under these conditions, one has

$$(4.24) \quad \nu = \frac{\beta^2 LM}{2(\beta c\mu - 1)} = \frac{2}{\mu^2} \left(\frac{L}{c}\right) \left(\frac{M}{c}\right).$$

We shall see the (potentially large) ratios  $L/c$  and  $M/c$  arise again later.

**Trade-Offs of (Mini-)Batching.** As a final observation about what can be learned from Theorems 4.6 and 4.7, let us take a moment to compare the theoretical performance of two fundamental algorithms—the simple SG iteration (3.7) and the mini-batch SG iteration (3.12)—when these results are applied for minimizing empirical risk, i.e., when  $F = R_n$ . This provides a glimpse into how such results can be used to compare algorithms in terms of their computational trade-offs.

The most elementary instance of our SG algorithm is simple SG, which, as we have seen, consists of picking a random sample index  $i_k$  at each iteration and computing

$$(4.25) \quad g(w_k, \xi_k) = \nabla f_{i_k}(w_k).$$

By contrast, instead of picking a single sample, minibatch SG consists of randomly selecting a subset  $\mathcal{S}_k$  of the sample indices and computing

$$(4.26) \quad g(w_k, \xi_k) = \frac{1}{|\mathcal{S}_k|} \sum_{i \in \mathcal{S}_k} \nabla f_i(w_k).$$

To compare these methods, let us assume for simplicity that we employ the same number of samples in each iteration so that the minibatches are of constant size, i.e.,  $|\mathcal{S}_k| = n_{\text{mb}}$ . There are then two distinct regimes to consider, namely, when  $n_{\text{mb}} \ll n$  and when  $n_{\text{mb}} \approx n$ . Our goal here is to use the results of Theorems 4.6 and 4.7 to show that, in the former scenario, the theoretical benefit of minibatching can appear to be somewhat ambiguous, meaning that one must leverage certain computational

<sup>5</sup>In fact, the bound (4.21) slightly overstates the asymptotic influence of the initial optimality gap. Applying Chung's lemma [36] to the contraction equation (4.23) shows that the first term in the definition of  $\nu$  effectively determines the asymptotic convergence rate of  $\mathbb{E}[F(w_k) - F_*]$ .

tools to benefit from minibatching in practice. As for the scenario when  $n_{\text{mb}} \approx n$ , the comparison is more complex due to a trade-off between per-iteration costs and overall convergence rate of the method (recall section 3.3). We leave a more formal treatment of this scenario, specifically with the goals of large-scale machine learning in mind, for section 4.4.

Suppose then that the minibatch size is  $n_{\text{mb}} \ll n$ . The computation of the stochastic direction  $g(w_k, \xi_k)$  in (4.26) is clearly  $n_{\text{mb}}$  times more expensive than in (4.25). In return, the variance of the direction is reduced by a factor of  $1/n_{\text{mb}}$ . (See section 5.2 for further discussion of this fact.) That is, with respect to our analysis, the constants  $M$  and  $M_V$  that appear in Assumption 4.3 (see (4.8)) are reduced by the same factor, becoming  $M/n_{\text{mb}}$  and  $M_V/n_{\text{mb}}$  for minibatch SG. It is natural to ask whether this reduction in the variance pays for the higher per-iteration cost.

Consider, for instance, the case of employing a sufficiently small constant stepsize  $\bar{\alpha} > 0$ . For minibatch SG, Theorem 4.6 leads to

$$\mathbb{E}[F(w_k) - F_*] \leq \frac{\bar{\alpha}LM}{2c\mu n_{\text{mb}}} + [1 - \bar{\alpha}c\mu]^{k-1} \left( F(w_1) - F_* - \frac{\bar{\alpha}LM}{2c\mu n_{\text{mb}}} \right).$$

Using the simple SG method with stepsize  $\bar{\alpha}/n_{\text{mb}}$  leads to a similar asymptotic gap:

$$\mathbb{E}[F(w_k) - F_*] \leq \frac{\bar{\alpha}LM}{2c\mu n_{\text{mb}}} + \left[ 1 - \frac{\bar{\alpha}c\mu}{n_{\text{mb}}} \right]^{k-1} \left( F(w_1) - F_* - \frac{\bar{\alpha}LM}{2c\mu n_{\text{mb}}} \right).$$

The worse contraction constant (indicated using square brackets) means that one needs to run  $n_{\text{mb}}$  times more iterations of the simple SG algorithm to obtain an equivalent optimality gap. That said, since the computation in a simple SG iteration is  $n_{\text{mb}}$  times cheaper, this amounts to effectively the same total computation as for the minibatch SG method. A similar analysis employing the result of Theorem 4.7 can be performed when decreasing stepsizes are used.

These observations suggest that the methods can be comparable. However, an important consideration remains. In particular, the convergence theorems require that the initial stepsize be smaller than  $\mu/(LM_G)$ . Since (4.9) shows that  $M_G \geq \mu^2$ , the largest this stepsize can be is  $1/(\mu L)$ . Therefore, one cannot simply assume that the minibatch SG method is allowed to employ a stepsize that is  $n_{\text{mb}}$  times larger than the one used by SG. In other words, one cannot always compensate for the higher per-iteration cost of a minibatch SG method by selecting a larger stepsize.

One can, however, realize the benefits of minibatching in practice since it offers the important opportunities for software optimization and parallelization; e.g., using sizable minibatches is often the only way to fully leverage a GPU processor. Dynamic minibatch sizes can also be used as a substitute for decreasing stepsizes; see section 5.2.

**4.3. SG for General Objectives.** As mentioned in our case study of DNNs in section 2.2, many important machine learning models lead to nonconvex optimization problems, which are currently having a profound impact in practice. Analyzing the SG method when minimizing nonconvex objectives is more challenging than in the convex case since such functions may possess multiple local minima and other stationary points. Still, we show in this subsection that one can provide meaningful guarantees for the SG method in nonconvex settings.

Paralleling section 4.2, we present two results—one for employing a fixed positive stepsize and one for diminishing stepsizes. We maintain the same assumptions about

the stochastic directions  $g(w_k, \xi_k)$ , but do not assume convexity of  $F$ . As before, the results in this section still apply to a wide class of methods since  $g(w_k, \xi_k)$  could be defined as a (minibatch) stochastic gradient or a Newton-like direction; recall (4.2).

Our first result describes the behavior of the sequence of gradients of  $F$  when fixed stepsizes are employed. Recall from Assumption 4.3 that the sequence of function values  $\{F(w_k)\}$  is assumed to be bounded below by a scalar  $F_{\inf}$ .

**THEOREM 4.8** (nonconvex objective, fixed stepsize). *Under Assumptions 4.1 and 4.3, suppose that the SG method (Algorithm 4.1) is run with a fixed stepsize,  $\alpha_k = \bar{\alpha}$  for all  $k \in \mathbb{N}$ , satisfying*

$$(4.27) \quad 0 < \bar{\alpha} \leq \frac{\mu}{LM_G}.$$

*Then, the expected sum-of-squares and average-squared gradients of  $F$  corresponding to the SG iterates satisfy the following inequalities for all  $K \in \mathbb{N}$ :*

$$(4.28a) \quad \mathbb{E} \left[ \sum_{k=1}^K \|\nabla F(w_k)\|_2^2 \right] \leq \frac{K\bar{\alpha}LM}{\mu} + \frac{2(F(w_1) - F_{\inf})}{\mu\bar{\alpha}}$$

$$(4.28b) \quad \text{and therefore } \mathbb{E} \left[ \frac{1}{K} \sum_{k=1}^K \|\nabla F(w_k)\|_2^2 \right] \leq \frac{\bar{\alpha}LM}{\mu} + \frac{2(F(w_1) - F_{\inf})}{K\mu\bar{\alpha}} \\ \xrightarrow{K \rightarrow \infty} \frac{\bar{\alpha}LM}{\mu}.$$

*Proof.* Taking the total expectation of (4.10b) and using (4.27),

$$\begin{aligned} \mathbb{E}[F(w_{k+1})] - \mathbb{E}[F(w_k)] &\leq -(\mu - \tfrac{1}{2}\bar{\alpha}LM_G)\bar{\alpha}\mathbb{E}[\|\nabla F(w_k)\|_2^2] + \tfrac{1}{2}\bar{\alpha}^2LM \\ &\leq -\tfrac{1}{2}\mu\bar{\alpha}\mathbb{E}[\|\nabla F(w_k)\|_2^2] + \tfrac{1}{2}\bar{\alpha}^2LM. \end{aligned}$$

Summing both sides of this inequality for  $k \in \{1, \dots, K\}$  and recalling Assumption 4.3(a) gives

$$F_{\inf} - F(w_1) \leq \mathbb{E}[F(w_{K+1})] - F(w_1) \leq -\tfrac{1}{2}\mu\bar{\alpha} \sum_{k=1}^K \mathbb{E}[\|\nabla F(w_k)\|_2^2] + \tfrac{1}{2}K\bar{\alpha}^2LM.$$

Rearranging yields (4.28a), and dividing further by  $K$  yields (4.28b).  $\square$

If  $M = 0$ , meaning that there is no noise or that noise reduces proportionally to  $\|\nabla F(w_k)\|_2^2$  (see (4.8) and (4.9)), then (4.28a) captures a classical result for the full gradient method applied to nonconvex functions, namely, that the sum of squared gradients remains finite, implying that  $\{\|\nabla F(w_k)\|_2\} \rightarrow 0$ . In the presence of noise (i.e.,  $M > 0$ ), Theorem 4.8 illustrates the interplay between the stepsize  $\bar{\alpha}$  and the variance of the stochastic directions. While one cannot bound the expected optimality gap as in the convex case, inequality (4.28b) bounds the average norm of the gradient of the objective function observed on  $\{w_k\}$  visited during the first  $K$  iterations. This quantity gets smaller when  $K$  increases, indicating that the SG method spends increasingly more time in regions where the objective function has a (relatively) small gradient. Moreover, the asymptotic result that one obtains from (4.28b) illustrates that noise in the gradients inhibits further progress, as in (4.14) for the convex case. The average norm of the gradients can be made arbitrarily small by selecting a small

stepsize, but doing so reduces the speed at which the norm of the gradient approaches its limiting distribution.

We now turn to the case when the SG method is applied to a nonconvex objective with a decreasing sequence of stepsizes satisfying the classical conditions (4.19). While not the strongest result that one can prove in this context—and, in fact, we prove a stronger result below—the following theorem is perhaps the easiest to interpret and remember. Hence, we state it first.

**THEOREM 4.9** (nonconvex objective, diminishing stepsizes). *Under Assumptions 4.1 and 4.3, suppose that the SG method (Algorithm 4.1) is run with a stepsize sequence satisfying (4.19). Then*

$$(4.29) \quad \liminf_{k \rightarrow \infty} \mathbb{E}[\|\nabla F(w_k)\|_2^2] = 0.$$

The proof of this theorem follows based on the results given in Theorem 4.10 below. A “liminf” result of this type should be familiar to those knowledgeable about the nonlinear optimization literature. After all, such a result is all that can be shown for certain important methods, such as the nonlinear conjugate gradient method [114]. The intuition that one should gain from the statement of Theorem 4.9 is that for the SG method with diminishing stepsizes, the expected gradient norms cannot stay bounded away from zero.

The following result characterizes more precisely the convergence property of SG.

**THEOREM 4.10** (nonconvex objective, diminishing stepsizes). *Under Assumptions 4.1 and 4.3, suppose that the SG method (Algorithm 4.1) is run with a stepsize sequence satisfying (4.19). Then, with  $A_K := \sum_{k=1}^K \alpha_k$ ,*

$$(4.30a) \quad \lim_{K \rightarrow \infty} \mathbb{E} \left[ \sum_{k=1}^K \alpha_k \|\nabla F(w_k)\|_2^2 \right] < \infty$$

$$(4.30b) \quad \text{and therefore } \mathbb{E} \left[ \frac{1}{A_K} \sum_{k=1}^K \alpha_k \|\nabla F(w_k)\|_2^2 \right] \xrightarrow{K \rightarrow \infty} 0.$$

*Proof.* The second condition in (4.19) ensures that  $\{\alpha_k\} \rightarrow 0$ , meaning that, without loss of generality, we may assume that  $\alpha_k LM_G \leq \mu$  for all  $k \in \mathbb{N}$ . Then, taking the total expectation of (4.10b),

$$\begin{aligned} \mathbb{E}[F(w_{k+1})] - \mathbb{E}[F(w_k)] &\leq -(\mu - \tfrac{1}{2}\alpha_k LM_G)\alpha_k \mathbb{E}[\|\nabla F(w_k)\|_2^2] + \tfrac{1}{2}\alpha_k^2 LM \\ &\leq -\tfrac{1}{2}\mu\alpha_k \mathbb{E}[\|\nabla F(w_k)\|_2^2] + \tfrac{1}{2}\alpha_k^2 LM. \end{aligned}$$

Summing both sides of this inequality for  $k \in \{1, \dots, K\}$  gives

$$F_{\inf} - \mathbb{E}[F(w_1)] \leq \mathbb{E}[F(w_{K+1})] - \mathbb{E}[F(w_1)] \leq -\tfrac{1}{2}\mu \sum_{k=1}^K \alpha_k \mathbb{E}[\|\nabla F(w_k)\|_2^2] + \tfrac{1}{2}LM \sum_{k=1}^K \alpha_k^2.$$

Dividing by  $\mu/2$  and rearranging the terms, we obtain

$$\sum_{k=1}^K \alpha_k \mathbb{E}[\|\nabla F(w_k)\|_2^2] \leq \frac{2(\mathbb{E}[F(w_1)] - F_{\inf})}{\mu} + \frac{LM}{\mu} \sum_{k=1}^K \alpha_k^2.$$

The second condition in (4.19) implies that the right-hand side of this inequality converges to a finite limit when  $K$  increases, proving (4.30a). Then, (4.30b) follows since the first condition in (4.19) ensures that  $A_K \rightarrow \infty$  as  $K \rightarrow \infty$ .  $\square$



Theorem 4.10 establishes results about a weighted sum-of-squares and a weighted average of squared gradients of  $F$  similar to those in Theorem 4.8. However, unlike (4.28b), the conclusion (4.30b) states that the weighted average norm of the squared gradients converges to zero even if the gradients are noisy, i.e., if  $M > 0$ . The fact that (4.30b) only specifies a property of a weighted average (with weights dictated by  $\{\alpha_k\}$ ) is only of minor importance since one can still conclude that the expected gradient norms cannot asymptotically stay far from zero.

We can now see that Theorem 4.9 is a direct consequence of Theorem 4.10, for if (4.29) did not hold, it would contradict Theorem 4.10.

The next result gives a stronger conclusion than Theorem 4.9, at the expense of only showing a property of the gradient of  $F$  at a randomly selected iterate.

**COROLLARY 4.11.** *Suppose the conditions of Theorem 4.10 hold. For any  $K \in \mathbb{N}$ , let  $k(K) \in \{1, \dots, K\}$  represent a random index chosen with probabilities proportional to  $\{\alpha_k\}_{k=1}^K$ . Then,  $\|\nabla F(w_{k(K)})\|_2 \xrightarrow{K \rightarrow \infty} 0$  in probability.*

*Proof.* Using Markov's inequality and (4.30a), for any  $\varepsilon > 0$ , we can write

$$\mathbb{P}\{\|\nabla F(w_k)\|_2 \geq \varepsilon\} = \mathbb{P}\{\|\nabla F(w_k)\|_2^2 \geq \varepsilon^2\} \leq \varepsilon^{-2} \mathbb{E}[\mathbb{E}_k[\|\nabla F(w_k)\|_2^2]] \xrightarrow{K \rightarrow \infty} 0,$$

which is the definition of convergence in probability.  $\square$

Finally, we present the following result (with proof in Appendix B) to illustrate that stronger convergence results also follow under additional regularity conditions.

**COROLLARY 4.12.** *Under the conditions of Theorem 4.10, if we further assume that the objective function  $F$  is twice differentiable, and that the mapping  $w \mapsto \|\nabla F(w)\|_2^2$  has Lipschitz-continuous derivatives, then*

$$\lim_{k \rightarrow \infty} \mathbb{E}[\|\nabla F(w_k)\|_2^2] = 0.$$

**4.4. Work Complexity for Large-Scale Learning.** Our discussion thus far has focused on the convergence properties of SG when minimizing a given objective function representing either expected or empirical risk. However, our investigation would be incomplete without considering how these properties impact the computational workload associated with solving an underlying machine learning problem. As previewed in section 3, there are arguments that a more slowly convergent algorithm such as SG, with its sublinear rate of convergence, is more efficient for large-scale learning than (full, batch) gradient-based methods that have a linear rate of convergence. The purpose of this section is to present these arguments in further detail.

As a first attempt for setting up a framework in which one might compare optimization algorithms for large-scale learning, one might be tempted to consider the situation in which one has a given training set size  $n$  and asks what type of algorithm—e.g., a simple SG or batch gradient method—would provide the best guarantees in terms of achieving a low expected risk. However, such a comparison is difficult to make when one cannot determine the precise trade-off between per-iteration costs and the overall progress of the optimization method that one can guarantee.

An easier way to approach the issue is to consider a *big data* scenario with an infinite supply of training examples, but a limited computational time budget. One can then ask whether running a simple optimization algorithm such as SG works better than running a more sophisticated batch optimization algorithm. We use such an approach next, following the work in [22, 23].

Suppose that both the expected risk  $R$  and the empirical risk  $R_n$  attain their minima with parameter vectors  $w_* \in \arg \min R(w)$  and  $w_n \in \arg \min R_n(w)$ , respectively. In addition, let  $\tilde{w}_n$  be the approximate empirical risk minimizer returned by a given optimization algorithm when the time budget  $\mathcal{T}_{\max}$  is exhausted. The trade-offs associated with this scenario can be formalized as choosing the family of prediction functions  $\mathcal{H}$ , the number of examples  $n$ , and the optimization accuracy  $\epsilon := \mathbb{E}[R_n(\tilde{w}_n) - R_n(w_n)]$  in order to minimize, within time  $\mathcal{T}_{\max}$ , the *total error*

$$(4.31) \quad \mathbb{E}[R(\tilde{w}_n)] = \underbrace{R(w_*)}_{\mathcal{E}_{app}(\mathcal{H})} + \underbrace{\mathbb{E}[R(w_n) - R(w_*)]}_{\mathcal{E}_{est}(\mathcal{H}, n)} + \underbrace{\mathbb{E}[R(\tilde{w}_n) - R(w_n)]}_{\mathcal{E}_{opt}(\mathcal{H}, n, \epsilon)}.$$

To minimize this error, one needs to balance the contributions from each of the three terms on the right-hand side. For instance, if one decides to make the optimization more accurate—i.e., reducing  $\epsilon$  in the hope of also reducing the *optimization error*  $\mathcal{E}_{opt}(\mathcal{H}, n, \epsilon)$  (evaluated with respect to  $R$  rather than  $R_n$ )—one might need to make up for the additional computing time by (i) reducing the sample size  $n$ , potentially increasing the *estimation error*  $\mathcal{E}_{est}(\mathcal{H}, n)$ , or (ii) simplifying the function family  $\mathcal{H}$ , potentially increasing the *approximation error*  $\mathcal{E}_{app}(\mathcal{H})$ .

Useful guidelines for achieving an optimal balance between these errors can be obtained by setting aside the choice of  $\mathcal{H}$  and carrying out a worst-case analysis of the influence of the sample size  $n$  and optimization tolerance  $\epsilon$ , which together only influence the estimation and optimization errors. This simplified setup can be formalized in terms of the macroscopic optimization problem

$$(4.32) \quad \min_{n, \epsilon} \mathcal{E}(n, \epsilon) = \mathbb{E}[R(\tilde{w}_n) - R(w_*)] \text{ s.t. } \mathcal{T}(n, \epsilon) \leq \mathcal{T}_{\max}.$$

The computing time  $\mathcal{T}(n, \epsilon)$  depends on the details of the optimization algorithm in interesting ways. For example, the computing time of a batch algorithm increases linearly (at least) with the number of examples  $n$ , whereas, crucially, the computing time of a stochastic algorithm is independent of  $n$ . With a batch optimization algorithm, one could consider increasing  $\epsilon$  in order to make time to use more training examples. However, with a stochastic algorithm, one should always use as many examples as possible because the per-iteration computing time is independent of  $n$ .

To be specific, let us compare the solutions of (4.32) for prototypical stochastic and batch methods—namely, simple SG and a batch gradient method—using simplified forms for the worst cases of the error function  $\mathcal{E}$  and the time function  $\mathcal{T}$ . For the error function, a direct application of the uniform laws of large numbers [155] yields

$$\begin{aligned} \mathcal{E}(n, \epsilon) = \mathbb{E}[R(\tilde{w}_n) - R(w_*)] &= \underbrace{\mathbb{E}[R(\tilde{w}_n) - R_n(\tilde{w}_n)]}_{= \mathcal{O}\left(\sqrt{\log(n)/n}\right)} + \underbrace{\mathbb{E}[R_n(\tilde{w}_n) - R_n(w_n)]}_{= \epsilon} \\ &\quad + \underbrace{\mathbb{E}[R_n(w_n) - R_n(w_*)]}_{\leq 0} + \underbrace{\mathbb{E}[R_n(w_*) - R(w_*)]}_{= \mathcal{O}\left(\sqrt{\log(n)/n}\right)}, \end{aligned}$$

which leads to the upper bound

$$(4.33) \quad \mathcal{E}(n, \epsilon) = \mathcal{O}\left(\sqrt{\frac{\log(n)}{n}} + \epsilon\right).$$

The inverse-square-root dependence on the number of examples  $n$  that appears here is typical of statistical problems. However, even faster convergence rates for reducing these terms with respect to  $n$  can be established under specific conditions. For instance, when the loss function is strongly convex [91] or when the data distribution satisfies certain assumptions [154], it is possible to show that

$$\mathcal{E}(n, \epsilon) = \mathcal{O} \left( \frac{\log(n)}{n} + \epsilon \right).$$

To simplify further, let us work with the asymptotic (i.e., for large  $n$ ) equivalence

$$(4.34) \quad \mathcal{E}(n, \epsilon) \sim \frac{1}{n} + \epsilon,$$

which is the fastest rate that remains compatible with elementary statistical results.<sup>6</sup> Under this assumption, noting that the time constraint in (4.32) will always be active (since one can always lower  $\epsilon$ , and hence  $\mathcal{E}(n, \epsilon)$ , by giving more time to the optimization algorithm), and recalling the worst-case computing time bounds introduced in section 3.3, one arrives at the following conclusions.

- A simple SG method can achieve  $\epsilon$ -optimality with a computing time of  $\mathcal{T}_{\text{stoch}}(n, \epsilon) \sim 1/\epsilon$ . Hence, within the time budget  $\mathcal{T}_{\text{max}}$ , the accuracy achieved is proportional to  $1/\mathcal{T}_{\text{max}}$ , regardless of the  $n$ . This means that, to minimize the error  $\mathcal{E}(n, \epsilon)$ , one should simply choose  $n$  as large as possible. Since the maximum number of examples that can be processed by SG during the time budget is proportional to  $\mathcal{T}_{\text{max}}$ , it follows that the optimal error is proportional to  $1/\mathcal{T}_{\text{max}}$ .
- A batch gradient method can achieve  $\epsilon$ -optimality with a computing time of  $\mathcal{T}_{\text{batch}}(n, \epsilon) \sim n \log(1/\epsilon)$ . This means that, within the time budget  $\mathcal{T}_{\text{max}}$ , it can achieve  $\epsilon$ -optimality by processing  $n \sim \mathcal{T}_{\text{max}}/\log(1/\epsilon)$  examples. One now finds that the optimal error is not necessarily achieved by choosing  $n$  as large as possible, but rather by choosing  $\epsilon$  (which dictates  $n$ ) to minimize (4.34). Differentiating  $\mathcal{E}(n, \epsilon) \sim \log(1/\epsilon)/\mathcal{T}_{\text{max}} + \epsilon$  with respect to  $\epsilon$  and setting the result equal to zero, one finds that optimality is achieved with  $\epsilon \sim 1/\mathcal{T}_{\text{max}}$ , from which it follows that the optimal error is proportional to  $\log(\mathcal{T}_{\text{max}})/\mathcal{T}_{\text{max}} + 1/\mathcal{T}_{\text{max}}$ .

These results are summarized in Table 4.1. Even though a batch approach possesses a better dependency on  $\epsilon$ , this advantage does not make up for its dependence on  $n$ . This is true even though we have employed (4.34), the most favorable rate that one may reasonably assume. In conclusion, we have found that a stochastic optimization algorithm performs better in terms of expected error, and, hence, makes a better learning algorithm in the sense considered here. These observations are supported by practical experience (recall Figure 3.1 in section 3.3).

**A Lower Bound.** The results reported in Table 4.1 are also notable because the SG algorithm matches a lower complexity bound that has been established for any optimization method employing a *noisy oracle*. To be specific, in the widely employed model for studying optimization algorithms proposed by Nemirovsky and

<sup>6</sup>For example, suppose that one is estimating the mean of a distribution  $P$  defined on the real line by minimizing the risk  $R(\mu) = \int (x - \mu)^2 dP(x)$ . The convergence rate (4.34) amounts to estimating the distribution mean with a variance proportional to  $1/n$ . A faster convergence rate would violate the Cramer–Rao bound.

**Table 4.1** The first row displays the computing times of idealized batch and stochastic optimization algorithms. The second row gives the corresponding solutions of (4.32), assuming (4.34).

	Batch	Stochastic
$\mathcal{T}(n, \epsilon)$	$\sim n \log\left(\frac{1}{\epsilon}\right)$	$\frac{1}{\epsilon}$
$\mathcal{E}^*$	$\sim \frac{\log(\mathcal{T}_{\max})}{\mathcal{T}_{\max}} + \frac{1}{\mathcal{T}_{\max}}$	$\frac{1}{\mathcal{T}_{\max}}$

Yudin [106], one assumes that information regarding the objective function is acquired by querying an *oracle*, ignoring the computational demands of doing so. Using such a model, it has been established, e.g., that the full gradient method applied to minimize a strongly convex objective function is not optimal in terms of the accuracy that can be achieved within a given number of calls to the oracle, but that one can achieve an optimal method through acceleration techniques; see section 7.2.

The case when only gradient estimates are available through a noisy oracle has been studied, e.g., in [1, 128]. Roughly speaking, these investigations show that, again when minimizing a strongly convex function, no algorithm that performs  $k$  calls to the oracle can guarantee accuracy better than  $\Omega(1/k)$ . As we have seen, SG achieves this lower bound up to constant factors. This analysis applies for the optimization of both expected risk and empirical risk.

**4.5. Commentary.** Although the analysis presented in section 4.4 can be quite compelling, it would be premature to conclude that SG is a perfect solution for large-scale machine learning problems. There is, in fact, a large gap between asymptotical behavior and practical realities. Next, we discuss issues related to this gap.

**Fragility of the Asymptotic Performance of SG.** The convergence speed given, e.g., by Theorem 4.7, holds when the stepsize constant  $\beta$  exceeds a quantity inversely proportional to the strong convexity parameter  $c$  (see (4.20)). In some settings, determining such a value is relatively easy, such as when the objective function includes a squared  $\ell_2$ -norm regularizer (e.g., as in (2.3)), in which case the regularization parameter provides a lower bound for  $c$ . However, despite the fact that this can work well in practice, it is not completely satisfactory because one should reduce the regularization parameter when the number of samples increases. It is therefore desirable to design algorithms that adapt to local convexity properties of the objective, so as to avoid having to place cumbersome restrictions on the stepsizes.

**SG and Ill-Conditioning.** The analysis of section 4.4 is compelling since, as long as the optimization problem is reasonably well-conditioned, the constant factors favor the SG algorithm. In particular, the minimalism of the SG algorithm allows for very efficient implementations that either fully leverage the sparsity of training examples (as in the case study on text classification in section 2.1) or harness the computational power of GPU processors (as in the case study on DNNs in section 2.2). In contrast, state-of-the-art batch optimization algorithms often carry more overhead. However, this advantage is eroded when the conditioning of the objective function worsens. Again, consider Theorem 4.7. This result involves constant factors that grow with both the condition number  $L/c$  and the ratio  $M/c$ . Both of these ratios can be improved greatly by adaptively rescaling the stochastic directions based on matrices that capture local curvature information of the objective function; see section 6.

**Opportunities for Distributed Computing.** Distributing the SG step computation can potentially reduce the computing time by a constant factor equal to the number of machines. However, such an improvement is difficult to realize. The SG algorithm is notoriously difficult to distribute efficiently because it accesses the shared parameter vector  $w$  with relatively high frequency. Consequently, even though it is very robust to additional noise and can be run with very relaxed synchronization [112, 45], distributed SG algorithms suffer from large communication overhead. Since this overhead is potentially much larger than the additional work associated with minibatch and other methods with higher per-iteration costs, distributed computing offers new opportunities for the success of such methods for machine learning.

**Alternatives with Faster Convergence.** As mentioned above, [1, 128] establish lower complexity bounds for optimization algorithms that only access information about the objective function through noisy estimates of  $F(w_k)$  and  $\nabla F(w_k)$  acquired in each iteration. The bounds apply, e.g., when SG is employed to minimize the expected risk  $R$  using gradient estimates evaluated on samples drawn from the distribution  $P$ . However, an algorithm that optimizes the empirical risk  $R_n$  has access to an additional piece of information: it knows when a gradient estimate is evaluated on a training example that has already been visited during previous iterations. Recent *gradient aggregation* methods (see section 5.3) make use of this information and improve upon the lower bound in [1] for the optimization of the empirical risk (though not for the expected risk). These algorithms enjoy linear convergence with low computing times in practice. Another avenue for improving the convergence rate is to employ a *dynamic sampling* approach (see section 5.2), which, as we shall see, can match the optimal asymptotic efficiency of SG in big data settings.

**5. Noise Reduction Methods.** The theoretical arguments in the previous section, together with extensive computational experience, have led many in the machine learning community to view SG as the ideal optimization approach for large-scale applications. We argue, however, that this is far from settled. SG suffers from, among other things, the adverse effect of noisy gradient estimates. This prevents it from converging to the solution when fixed stepsizes are used and leads to a slow, sublinear rate of convergence when a diminishing stepsize sequence  $\{\alpha_k\}$  is employed.

To address this limitation, methods endowed with *noise reduction* capabilities have been developed. These methods, which reduce the errors in the gradient estimates and/or iterate sequence, have proved to be effective in practice and enjoy attractive theoretical properties. Recalling the schematic of optimization methods in Figure 3.3, we depict these methods on the horizontal axis given in Figure 5.1.

The first two classes of methods that we consider achieve noise reduction in a manner that allows them to possess a linear rate of convergence to the optimal value using a fixed stepsize. The first type, *dynamic sampling* methods, achieve noise reduction by gradually increasing the minibatch size used in the gradient computation, thus employing increasingly more accurate gradient estimates as the optimization process proceeds. *Gradient aggregation* methods, on the other hand, improve the quality of the search directions by storing gradient estimates corresponding to samples employed in previous iterations, updating one (or some) of these estimates in each iteration, and defining the search direction as a weighted average of these estimates.

The third class of methods that we consider, *iterate averaging* methods, accomplish noise reduction not by averaging gradient estimates, but by maintaining an average of iterates computed during the optimization process. Employing a more aggressive stepsize sequence—of order  $\mathcal{O}(1/\sqrt{k})$  rather than  $\mathcal{O}(1/k)$ , which is appealing

Convergence results for SG and its variants are occasionally established using *regret bounds* as an intermediate step [164, 144, 54]. Regret bounds can be traced to Novikoff's analysis of the Perceptron [115] and to Cover's universal portfolios [41]. To illustrate, suppose that one aims to minimize a convex expected risk measure  $R(w) = \mathbb{E}[f(w; \xi)]$  over  $w \in \mathbb{R}^d$  with minimizer  $w_* \in \mathbb{R}^d$ . At a given iterate  $w_k$ , one obtains by convexity of  $f(w; \xi_k)$  (recall (A.1)) that

$$\|w_{k+1} - w_*\|^2 - \|w_k - w_*\|^2 \leq -2\alpha_k(f(w_k; \xi_k) - f(w_*; \xi_k)) + \alpha_k^2 \|\nabla f(w_k; \xi_k)\|_2^2.$$

Following [164], assuming that  $\|\nabla f(w_k; \xi_k)\|_2^2 \leq M$  and  $\|w_k - w_*\|_2^2 < B$  for some constants  $M > 0$  and  $B > 0$  for all  $k \in \mathbb{N}$ , one finds that

$$\begin{aligned} & \alpha_{k+1}^{-1} \|w_{k+1} - w_*\|_2^2 - \alpha_k^{-1} \|w_k - w_*\|_2^2 \\ & \leq -2(f(w_k; \xi_k) - f(w_*; \xi_k)) + \alpha_k M + (\alpha_{k+1}^{-1} - \alpha_k^{-1}) \|w_k - w_*\|_2^2 \\ & \leq -2(f(w_k; \xi_k) - f(w_*; \xi_k)) + \alpha_k M + (\alpha_{k+1}^{-1} - \alpha_k^{-1}) B. \end{aligned}$$

Summing for  $k = \{1, \dots, K\}$  with stepsizes  $\alpha_k = 1/\sqrt{k}$  leads to the regret bound

$$(4.35) \quad \left( \sum_{k=1}^K f(w_k; \xi_k) \right) \leq \left( \sum_{k=1}^K f(w_*; \xi_k) \right) + M\sqrt{K} + o(\sqrt{K}),$$

which bounds the losses incurred from  $\{w_k\}$  compared to those yielded by the fixed vector  $w_*$ . Such a bound is remarkable because its derivation holds for any sequence of noise variables  $\{\xi_k\}$ . This means that the average loss observed during the execution of the SG algorithm is *never* much worse than the best average loss one would have obtained if the optimal parameter  $w_*$  were known in advance. Further assuming that the noise variables are independent and using a martingale argument [34] leads to more traditional results of the form

$$\mathbb{E} \left[ \frac{1}{K} \sum_{k=1}^K F(w_k) \right] \leq F_* + \mathcal{O} \left( \frac{1}{\sqrt{K}} \right).$$

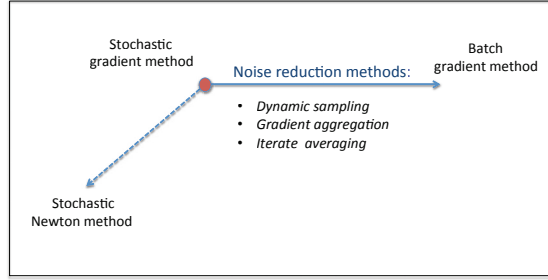
As long as one makes the same independent noise assumption, results obtained with this technique cannot be fundamentally different from the results that we have established. However, one should appreciate that the regret bound (4.35) itself remains meaningful when the noise variables are dependent or even adversarial. Such results reveal interesting connections between probability theory, machine learning, and game theory [143, 34].

---

**Inset 4.2** *Regret bounds.*

in itself—it is this sequence of averaged iterates that converges to the solution. These methods are closer in spirit to SG and their rate of convergence remains sublinear, but it can be shown that the variance of the sequence of average iterates is smaller than the variance of the SG iterates.

To formally motivate a concept of noise reduction, we begin this section by discussing a fundamental result that stipulates a rate of decrease in noise that allows an SG-type method to converge at a linear rate. We then show that a dynamic sampling method that enforces such noise reduction enjoys optimal complexity bounds, as defined in section 4.4. Next, we discuss three gradient aggregation methods—SVRG,



**Fig. 5.1** View of the schematic from Figure 3.3 with a focus on noise reduction methods.

SAGA, and SAG—the first of which can be viewed as a bridge between methods that control errors in the gradient with methods like SAGA and SAG in which noise reduction is accomplished in a more subtle manner. We conclude with a discussion of the practical and theoretical properties of iterate averaging methods.

**5.1. Reducing Noise at a Geometric Rate.** Let us recall the fundamental inequality (4.4), which we restate here for convenience:

$$\mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k) \leq -\alpha_k \nabla F(w_k)^T \mathbb{E}_{\xi_k}[g(w_k, \xi_k)] + \frac{1}{2} \alpha_k^2 L \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2].$$

(Recall that, as stated in (4.1), the objective  $F$  could stand for the expected risk  $R$  or the empirical risk  $R_n$ .) It is intuitively clear that if  $-g(w_k, \xi_k)$  is a descent direction in expectation (making the first term on the right-hand side negative), and if we are able to decrease  $\mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2]$  fast enough (along with  $\|\nabla F(w_k)\|_2^2$ ), then the effect of having noisy directions will not impede a fast rate of convergence. From another point of view, we can expect such behavior if, in Assumption 4.3, we suppose instead that the variance of  $g(w_k, \xi_k)$  vanishes sufficiently quickly.

We formalize this intuitive argument by considering the SG method with a fixed stepsize and showing that the sequence of expected optimality gaps vanishes at a linear rate as long as the variance of the stochastic vectors, denoted by  $\mathbb{V}_{\xi_k}[g(w_k, \xi_k)]$  (recall (4.6)), decreases *geometrically*.

**THEOREM 5.1** (strongly convex objective, noise reduction). *Suppose that Assumptions 4.1, 4.3, and 4.5 (with  $F_{\inf} = F_*$ ) hold, but with (4.8) refined to the existence of constants  $M \geq 0$  and  $\zeta \in (0, 1)$  such that, for all  $k \in \mathbb{N}$ ,*

$$(5.1) \quad \mathbb{V}_{\xi_k}[g(w_k, \xi_k)] \leq M \zeta^{k-1}.$$

*In addition, suppose that the SG method (Algorithm 4.1) is run with a fixed stepsize,  $\alpha_k = \bar{\alpha}$  for all  $k \in \mathbb{N}$ , satisfying*

$$(5.2) \quad 0 < \bar{\alpha} \leq \min \left\{ \frac{\mu}{L\mu_G^2}, \frac{1}{c\mu} \right\}.$$

*Then, for all  $k \in \mathbb{N}$ , the expected optimality gap satisfies*

$$(5.3) \quad \mathbb{E}[F(w_k) - F_*] \leq \omega \rho^{k-1},$$

where

$$(5.4a) \quad \omega := \max\left\{\frac{\bar{\alpha}LM}{c\mu}, F(w_1) - F_*\right\}$$

$$(5.4b) \quad \text{and } \rho := \max\left\{1 - \frac{\bar{\alpha}c\mu}{2}, \zeta\right\} < 1.$$

*Proof.* By Lemma 4.4 (specifically, (4.10a)), we have

$$\mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k) \leq -\mu\bar{\alpha}\|\nabla F(w_k)\|_2^2 + \frac{1}{2}\bar{\alpha}^2L\mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2].$$

Hence, from (4.6), (4.7b), (5.1), (5.2), and (4.12), we have

$$\begin{aligned} \mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k) &\leq -\mu\bar{\alpha}\|\nabla F(w_k)\|_2^2 + \frac{1}{2}\bar{\alpha}^2L(\mu_G^2\|\nabla F(w_k)\|_2^2 + M\zeta^{k-1}) \\ &\leq -(\mu - \frac{1}{2}\bar{\alpha}L\mu_G^2)\bar{\alpha}\|\nabla F(w_k)\|_2^2 + \frac{1}{2}\bar{\alpha}^2LM\zeta^{k-1} \\ &\leq -\frac{1}{2}\mu\bar{\alpha}\|\nabla F(w_k)\|_2^2 + \frac{1}{2}\bar{\alpha}^2LM\zeta^{k-1} \\ &\leq -\bar{\alpha}c\mu(F(w_k) - F_*) + \frac{1}{2}\bar{\alpha}^2LM\zeta^{k-1}. \end{aligned}$$

Adding and subtracting  $F_*$  and taking total expectations, this yields

$$(5.5) \quad \mathbb{E}[F(w_{k+1}) - F_*] \leq (1 - \bar{\alpha}c\mu)\mathbb{E}[F(w_k) - F_*] + \frac{1}{2}\bar{\alpha}^2LM\zeta^{k-1}.$$

We now use induction to prove (5.3). By the definition of  $\omega$ , it holds for  $k = 1$ . Then, assuming it holds for  $k \geq 1$ , we have from (5.5), (5.4a), and (5.4b) that

$$\begin{aligned} \mathbb{E}[F(w_{k+1}) - F_*] &\leq (1 - \bar{\alpha}c\mu)\omega\rho^{k-1} + \frac{1}{2}\bar{\alpha}^2LM\zeta^{k-1} \\ &= \omega\rho^{k-1}\left(1 - \bar{\alpha}c\mu + \frac{\bar{\alpha}^2LM}{2\omega}\left(\frac{\zeta}{\rho}\right)^{k-1}\right) \\ &\leq \omega\rho^{k-1}\left(1 - \bar{\alpha}c\mu + \frac{\bar{\alpha}^2LM}{2\omega}\right) \\ &\leq \omega\rho^{k-1}\left(1 - \bar{\alpha}c\mu + \frac{\bar{\alpha}c\mu}{2}\right) \\ &= \omega\rho^{k-1}\left(1 - \frac{\bar{\alpha}c\mu}{2}\right) \\ &\leq \omega\rho^k, \end{aligned}$$

as desired.  $\square$

Consideration of the typical magnitudes of the constants  $\mu$ ,  $L$ ,  $\mu_G$ , and  $c$  in (5.2) reveals that the admissible range of values of  $\bar{\alpha}$  is large, i.e., the restriction on the stepsize  $\bar{\alpha}$  is not unrealistic in practical situations.

Now, a natural question to ask is whether one can design efficient optimization methods for attaining the critical bound (5.1) on the variance of the stochastic directions. We show next that a dynamic sampling method is one such approach.

**5.2. Dynamic Sample Size Methods.** Consider the iteration

$$(5.6) \quad w_{k+1} \leftarrow w_k - \bar{\alpha}g(w_k, \xi_k),$$

where the stochastic directions are computed for some  $\tau > 1$  as

$$(5.7) \quad g(w_k, \xi_k) := \frac{1}{n_k} \sum_{i \in \mathcal{S}_k} \nabla f(w_k; \xi_{k,i}) \quad \text{with } n_k := |\mathcal{S}_k| = \lceil \tau^{k-1} \rceil.$$



That is, consider a minibatch SG iteration with a fixed stepsize in which the minibatch size used to compute unbiased stochastic gradient estimates increases geometrically as a function of the iteration counter  $k$ . To show that such an approach can fall under the class of methods considered in Theorem 5.1, suppose that the samples represented by the random variables  $\{\xi_{k,i}\}_{i \in \mathcal{S}_k}$  are drawn independently according to  $P$  for all  $k \in \mathbb{N}$ . If we assume that each stochastic gradient  $\nabla f(w_k; \xi_{k,i})$  has an expectation equal to the true gradient  $\nabla F(w_k)$ , then (4.7a) holds with  $\mu_G = \mu = 1$ . If, in addition, the variance of each such stochastic gradient is equal and is bounded by  $M \geq 0$ , then for arbitrary  $i \in \mathcal{S}_k$  we have (see [61, p. 183])

$$(5.8) \quad \mathbb{V}_{\xi_k}[g(w_k, \xi_k)] \leq \frac{\mathbb{V}_{\xi_k}[\nabla f(w_k; \xi_{k,i})]}{n_k} \leq \frac{M}{n_k}.$$

This bound, when combined with the rate of increase in  $n_k$  given in (5.7), yields (5.1). We have thus shown that if one employs a minibatch SG method with (unbiased) gradient estimates computed as in (5.7), then by Theorem 5.1, one obtains linear convergence to the optimal value of a strongly convex function. We state this formally as the following corollary to Theorem 5.1.

**COROLLARY 5.2.** *Let  $\{w_k\}$  be the iterates generated by (5.6)–(5.7) with unbiased gradient estimates, i.e.,  $\mathbb{E}_{\xi_{k,i}}[\nabla f(w_k; \xi_{k,i})] = \nabla F(w_k)$  for all  $k \in \mathbb{N}$  and  $i \in \mathcal{S}_k$ . Then the variance condition (5.1) is satisfied, and if all other assumptions of Theorem 5.1 hold, then the expected optimality gap vanishes linearly in the sense of (5.3).*

The reader may question whether it is meaningful to describe a method as linearly convergent if the per-iteration cost increases without bound. In other words, it is not immediately apparent that such an algorithm is competitive with a classical SG approach even though the desired reduction in the gradient variance is achieved. To address this question, let us estimate the total work complexity of the dynamic sampling algorithm, defined as the number of evaluations of the individual gradients  $\nabla f(w_k; \xi_{k,i})$  required to compute an  $\epsilon$ -optimal solution, i.e., to achieve

$$(5.9) \quad \mathbb{E}[F(w_k) - F_*] \leq \epsilon.$$

We have seen that the simple SG method (3.7) requires one such evaluation per iteration, and that its rate of convergence for diminishing stepsizes (i.e., the only setup in which convergence to the solution can be guaranteed) is given by (4.21). Therefore, as previously mentioned, the number of stochastic gradient evaluations required by the SG method to guarantee (5.9) is  $\mathcal{O}(\epsilon^{-1})$ . We now show that the method (5.6)–(5.7) can attain the same complexity.

**THEOREM 5.3.** *Suppose that the dynamic sampling SG method (5.6)–(5.7) is run with a stepsize  $\bar{\alpha}$  satisfying (5.2) and some  $\tau \in (1, (1 - \frac{\bar{\alpha}c\mu}{2})^{-1}]$ . In addition, suppose that Assumptions 4.1, 4.3, and Assumption 4.5 (with  $F_{\inf} = F_*$ ) hold. Then the total number of evaluations of a stochastic gradient of the form  $\nabla f(w_k; \xi_{k,i})$  required to obtain (5.9) is  $\mathcal{O}(\epsilon^{-1})$ .*

*Proof.* We have that the conditions of Theorem 5.1 hold with  $\zeta = 1/\tau$ . Hence, we have from (5.3) that there exists  $\bar{k} \in \mathbb{N}$  such that (5.9) holds for all  $k \geq \bar{k}$ . We can then use  $\omega \rho^{\bar{k}-1} \leq \epsilon$  to write  $(\bar{k} - 1) \log \rho \leq \log(\epsilon/\omega)$ , which along with  $\rho \in (0, 1)$  (recall (5.4b)) implies that

$$(5.10) \quad \bar{k} - 1 \geq \left\lceil \frac{\log(\epsilon/\omega)}{\log \rho} \right\rceil = \left\lceil \frac{\log(\omega/\epsilon)}{-\log \rho} \right\rceil.$$

Let us now estimate the total number of sample gradient evaluations required up to iteration  $\bar{k}$ . We claim that, without loss of generality, we may assume that  $\log(\omega/\epsilon)/(-\log \rho)$  is integer-valued and that (5.10) holds at equality. Then, by (5.7), the number of sample gradients required in iteration  $\bar{k}$  is  $\lceil \tau^{\bar{k}-1} \rceil$ , where

$$\begin{aligned}\tau^{\bar{k}-1} &= \tau^{\frac{\log(\omega/\epsilon)}{-\log \rho}} \\ &= \exp\left(\log\left(\tau^{\frac{\log(\omega/\epsilon)}{-\log \rho}}\right)\right) \\ &= \exp\left(\left(\frac{\log(\omega/\epsilon)}{-\log \rho}\right) \log \tau\right) \\ &= (\exp(\log(\omega/\epsilon)))^{\frac{\log \tau}{-\log \rho}} \\ &= \left(\frac{\omega}{\epsilon}\right)^\theta \quad \text{with } \theta := \frac{\log \tau}{-\log \rho}.\end{aligned}$$

Therefore, the total number of sample gradient evaluations for the first  $\bar{k}$  iterations is

$$\sum_{j=1}^{\bar{k}} \lceil \tau^{j-1} \rceil \leq 2 \sum_{j=1}^{\bar{k}} \tau^{j-1} = 2 \left( \frac{\tau^{\bar{k}} - \tau}{\tau - 1} \right) = 2 \left( \frac{\tau(\omega/\epsilon)^\theta - \tau}{\tau - 1} \right) \leq 2 \left( \frac{\omega}{\epsilon} \right)^\theta \left( \frac{1}{1 - 1/\tau} \right).$$

In fact, since  $\tau \leq (1 - \frac{\bar{\alpha}c\mu}{2})^{-1}$ , it follows from (5.4b) that  $\rho = \zeta = \tau^{-1}$ , which implies that  $\theta = 1$ . Specifically, with  $\tau = (1 - \sigma(\frac{\bar{\alpha}c\mu}{2}))^{-1}$  for some  $\sigma \in [0, 1]$ , then  $\theta = 1$  and

$$\sum_{j=1}^{\bar{k}} \lceil \tau^{j-1} \rceil \leq \frac{4\omega}{\sigma\epsilon\bar{\alpha}c\mu},$$

as desired.  $\square$

The discussion so far has focused on dynamic sampling strategies for a gradient method, but these ideas also apply for second-order methods that incorporate the matrix  $H$  as in (4.2).

This leads to the following question: given the rate of convergence of a batch optimization algorithm on strongly convex functions (i.e., linear, superlinear, etc.), what should be the sampling rate so that the overall algorithm is *efficient* in the sense that it results in the lowest computational complexity? To answer this question, certain results have been established [120]: (i) if the optimization method has a sublinear rate of convergence, then there is no sampling rate that makes the algorithm “efficient”; (ii) if the optimization algorithm is linearly convergent, then the sampling rate must be geometric (with restrictions on the constant in the rate) for the algorithm to be “efficient”; (iii) for superlinearly convergent methods, increasing the sample size at a rate that is slightly faster than geometric will yield an “efficient” method.

**5.2.1. Practical Implementation.** We now address the question of how to design practical algorithms that enjoy the theoretical benefits of noise reduction through the use of increasing sample sizes.

One approach is to follow the theoretical results described above and tie the rate of growth in the sample size  $n_k = |\mathcal{S}_k|$  to the rate of convergence of the optimization iteration [62, 77]. Such an approach, which *presets* the sampling rate before running the optimization algorithm, requires some experimentation. For example, for the iteration (5.6), one needs to find a value of the parameter  $\tau > 1$  in (5.7) that yields

good performance for the application at hand. In addition, one may want to delay the application of dynamic sampling to prevent the full sample set from being employed too soon (or at all). Such heuristic adaptations could be difficult in practice.

An alternative is to consider mechanisms for choosing the sample sizes not according to a prescribed sequence, but *adaptively* according to information gathered during the optimization process. One avenue that has been explored along these lines has been to design techniques that produce descent directions sufficiently often [30, 73]. Such an idea is based on two observations. First, one can show that *any* direction  $g(w_k, \xi_k)$  is a descent direction for  $F$  at  $w_k$  if, for some  $\chi \in [0, 1)$ , one has

$$(5.11) \quad \delta(w_k, \xi_k) := \|g(w_k, \xi_k) - \nabla F(w_k)\|_2 \leq \chi \|g(w_k, \xi_k)\|_2.$$

To see this, note that  $\|\nabla F(w_k)\|_2 \geq (1 - \chi)\|g(w_k, \xi_k)\|_2$ , which after squaring (5.11) implies

$$\begin{aligned} \nabla F(w_k)^T g(w_k, \xi_k) &\geq \frac{1}{2}(1 - \chi^2)\|g(w_k, \xi_k)\|_2^2 + \|\nabla F(w_k)\|_2^2 \\ &\geq \frac{1}{2}(1 - \chi^2 + (1 - \chi)^2)\|g(w_k, \xi_k)\|_2^2 \\ &\geq (1 - \chi)\|g(w_k, \xi_k)\|_2^2. \end{aligned}$$

The second observation is that while one cannot cheaply verify the inequality in (5.11) because it involves the evaluation of  $\nabla F(w_k)$ , one can estimate the left-hand side  $\delta(w_k, \xi_k)$  and then choose  $n_k$  so that (5.11) holds sufficiently often.

Specifically, if we assume that  $g(w_k, \xi_k)$  is an unbiased gradient estimate, then

$$\mathbb{E}[\delta(w_k, \xi_k)^2] = \mathbb{E}[\|g(w_k, \xi_k) - \nabla F(w_k)\|_2^2] = \mathbb{V}_{\xi_k}[g(w_k, \xi_k)].$$

This variance is expensive to compute, but one can approximate it with a sample variance. For example, if the samples are drawn without replacement from a set of (very large) size  $n_k$ , then one has the approximation

$$\mathbb{V}_{\xi_k}[g(w_k, \xi_k)] \approx \frac{\text{trace}(\text{Cov}(\{\nabla f(w_k; \xi_{k,i})\}_{i \in \mathcal{S}_k}))}{n_k} =: \varphi_k.$$

An *adaptive sampling* algorithm thus tests the following condition in place of (5.11):

$$(5.12) \quad \varphi_k \leq \chi^2 \|g(w_k, \xi_k)\|_2^2.$$

If this condition is not satisfied, then one may consider increasing the sample size—either immediately in iteration  $k$  or in a subsequent iteration—to a size that one might predict would satisfy such a condition. This technique is algorithmically attractive, but does not guarantee that the sample size  $n_k$  increases at a geometric rate. One can, however, employ a backup [30, 73]: if (5.12) increases the sampling rate more slowly than a preset geometric sequence, then a growth in the sample size is imposed.

Dynamic sampling techniques are not yet widely used in machine learning, and we suspect that the practical technique presented here might serve merely as a starting point for further investigations. Ultimately, an algorithm that performs like SG at the start and transitions to a regime of reduced variance in an efficient manner could prove to be a very powerful method for large-scale machine learning.

**5.3. Gradient Aggregation.** The dynamic sample size methods described in the previous subsection require a careful balance in order to achieve the desired linear rate of convergence without jeopardizing per-iteration costs. Alternatively, one can

attempt to achieve an improved rate by asking a different question: rather than compute increasingly more *new* stochastic gradient information in each iteration, is it possible to achieve a lower variance by *reusing* and/or *revising* previously computed information? After all, if the current iterate has not been displaced too far from previous iterates, then stochastic gradient information from previous iterates may still be useful. In addition, if one maintains indexed gradient estimates in storage, then one can revise specific estimates as new information is collected. Ideas such as these lead to concepts of *gradient aggregation*. In this subsection, we present ideas for gradient aggregation in the context of minimizing a finite sum such as an empirical risk measure  $R_n$ , for which they were originally designed.

Gradient aggregation algorithms for minimizing finite sums that possess cheap per-iteration costs have a long history. For example, Bertsekas [13] and coauthors have proposed *incremental gradient* methods, the randomized versions of which can be viewed as instances of a basic SG method for minimizing a finite sum. However, the basic variants of these methods only achieve a sublinear rate of convergence. By contrast, the methods on which we focus in this section are able to achieve a linear rate of convergence on strongly convex problems. This improved rate is achieved primarily by either an increase in computation or an increase in storage.

**5.3.1. SVRG.** The first method we consider operates in cycles. At the beginning of each cycle, an iterate  $w_k$  is available at which the algorithm computes a batch gradient  $\nabla R_n(w_k) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(w_k)$ . Then, after initializing  $\tilde{w}_1 \leftarrow w_k$ , a set of  $m$  inner iterations indexed by  $j$  with an update  $\tilde{w}_{j+1} \leftarrow \tilde{w}_j - \alpha \tilde{g}_j$  are performed, where

$$(5.13) \quad \tilde{g}_j \leftarrow \nabla f_{i_j}(\tilde{w}_j) - (\nabla f_{i_j}(w_k) - \nabla R_n(w_k))$$

and  $i_j \in \{1, \dots, n\}$  is chosen at random. This formula has a simple interpretation. Since the expected value of  $\nabla f_{i_j}(w_k)$  over all possible  $i_j \in \{1, \dots, n\}$  is equal to  $\nabla R_n(w_k)$ , one can view  $\nabla f_{i_j}(w_k) - \nabla R_n(w_k)$  as the bias in the gradient estimate  $\nabla f_{i_j}(w_k)$ . Thus, in every iteration, the algorithm randomly draws a stochastic gradient  $\nabla f_{i_j}(\tilde{w}_j)$  evaluated at the current inner iterate  $\tilde{w}_j$  and corrects it based on a perceived bias. Overall,  $\tilde{g}_j$  represents an unbiased estimator of  $\nabla R_n(\tilde{w}_j)$ , but with a variance that one can expect to be smaller than if one were simply to chose  $\tilde{g}_j = \nabla f_{i_j}(\tilde{w}_j)$  (as in simple SG). This is the reason that the method is referred to as the *stochastic variance reduced gradient* (SVRG) method.

A formal description of a few variants of SVRG is presented as Algorithm 5.1. For both options (b) and (c), it has been shown that when applied to minimize a strongly convex  $R_n$ , Algorithm 5.1 can achieve a linear rate of convergence [82]. More precisely, if the stepsize  $\alpha$  and the length of the inner cycle  $m$  are chosen so that

$$\rho := \frac{1}{1 - 2\alpha L} \left( \frac{1}{m\alpha} + 2L\alpha \right) < 1,$$

then, given that the algorithm has reached  $w_k$ , one obtains

$$\mathbb{E}[R_n(w_{k+1}) - R_n(w_*)] \leq \rho \mathbb{E}[R_n(w_k) - R_n(w_*)]$$

(where expectation is taken with respect to the random variables in the inner cycle). It should be emphasized that the resulting linear convergence rate applies to the outer iterates  $\{w_k\}$ , where each step from  $w_k$  to  $w_{k+1}$  requires  $2m + n$  evaluations of

component gradients: step 7 requires two stochastic gradient evaluations, while step 3 requires  $n$  (a full gradient). Therefore, one iteration of SVRG is much more expensive than one of SG, and in fact is comparable to a full gradient iteration.

---

**Algorithm 5.1** SVRG Methods for Minimizing an Empirical Risk  $R_n$ 


---

```

1: Choose an initial iterate  $w_1 \in \mathbb{R}^d$ , stepsize  $\alpha > 0$ , and positive integer  $m$ .
2: for  $k = 1, 2, \dots$  do
3:   Compute the batch gradient  $\nabla R_n(w_k)$ .
4:   Initialize  $\tilde{w}_1 \leftarrow w_k$ .
5:   for  $j = 1, \dots, m$  do
6:     Choose  $i_j$  uniformly from  $\{1, \dots, n\}$ .
7:     Set  $\tilde{g}_j \leftarrow \nabla f_{i_j}(\tilde{w}_j) - (\nabla f_{i_j}(w_k) - \nabla R_n(w_k))$ .
8:     Set  $\tilde{w}_{j+1} \leftarrow \tilde{w}_j - \alpha \tilde{g}_j$ .
9:   end for
10:  Option (a): Set  $w_{k+1} = \tilde{w}_{m+1}$ .
11:  Option (b): Set  $w_{k+1} = \frac{1}{m} \sum_{j=1}^m \tilde{w}_{j+1}$ .
12:  Option (c): Choose  $j$  uniformly from  $\{1, \dots, m\}$  and set  $w_{k+1} = \tilde{w}_{j+1}$ .
13: end for

```

---

In practice, SVRG appears to be quite effective in certain applications compared with SG if one requires high training accuracy. For the first epochs, SG is more efficient, but once the iterates approach the solution the benefits of the fast convergence rate of SVRG can be observed. Without explicit knowledge of  $L$  and  $c$ , the lengths of the inner cycle  $m$  and the stepsize  $\alpha$  are typically both chosen by experimentation.

**5.3.2. SAGA.** The second method we consider employs an iteration that is closer in form to SG in that it does not operate in cycles, nor does it compute batch gradients (except possibly at the initial point). Instead, in each iteration, it computes a stochastic vector  $g_k$  as the average of stochastic gradients evaluated at previous iterates. Specifically, in iteration  $k$ , the method will have stored  $\nabla f_i(w_{[i]})$  for all  $i \in \{1, \dots, n\}$ , where  $w_{[i]}$  represents the latest iterate at which  $\nabla f_i$  was evaluated. An integer  $j \in \{1, \dots, n\}$  is then chosen at random and the stochastic vector is set by

$$(5.14) \quad g_k \leftarrow \nabla f_j(w_k) - \nabla f_j(w_{[j]}) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(w_{[i]}).$$

Taking the expectation of  $g_k$  with respect to all choices of  $j \in \{1, \dots, n\}$ , one again has that  $\mathbb{E}[g_k] = \nabla R_n(w_k)$ . Thus, the method employs unbiased gradient estimates, but with variances that are expected to be less than the stochastic gradients that would be employed in a basic SG routine. A precise algorithm employing (5.14), referred to as SAGA [46], is given in Algorithm 5.2.

Beyond its initialization phase, the per-iteration cost of Algorithm 5.2 is the same as in a basic SG method. However, it has been shown that the method can achieve a linear rate of convergence when minimizing a strongly convex  $R_n$ . Specifically, with  $\alpha = 1/(2(cn + L))$ , one can show that

$$\mathbb{E}[\|w_k - w_*\|_2^2] \leq \left(1 - \frac{c}{2(cn + L)}\right)^k \left(\|w_1 - w_*\|_2^2 + \frac{nD}{cn + L}\right),$$

where  $D := R_n(w_1) - \nabla R_n(w_*)^T(w_1 - w_*) - R_n(w_*)$ .

Of course, attaining such a result requires knowledge of the strong convexity constant  $c$

**Algorithm 5.2** SAGA Method for Minimizing an Empirical Risk  $R_n$ 


---

```

1: Choose an initial iterate  $w_1 \in \mathbb{R}^d$  and stepsize  $\alpha > 0$ .
2: for  $i = 1, \dots, n$  do
3:   Compute  $\nabla f_i(w_1)$ .
4:   Store  $\nabla f_i(w_{[i]}) \leftarrow \nabla f_i(w_1)$ .
5: end for
6: for  $k = 1, 2, \dots$  do
7:   Choose  $j$  uniformly in  $\{1, \dots, n\}$ .
8:   Compute  $\nabla f_j(w_k)$ .
9:   Set  $g_k \leftarrow \nabla f_j(w_k) - \nabla f_j(w_{[j]}) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(w_{[i]})$ .
10:  Store  $\nabla f_j(w_{[j]}) \leftarrow \nabla f_j(w_k)$ .
11:  Set  $w_{k+1} \leftarrow w_k - \alpha g_k$ .
12: end for

```

---

and Lipschitz constant  $L$ . If  $c$  is not known, then the stepsize can instead be chosen to be  $\alpha = 1/(3L)$  and a similar convergence result can be established; see [46].

Alternative initialization techniques could be used in practice, which may be more effective than evaluating all the gradients  $\{\nabla f_i\}_{i=1}^n$  at the initial point. For example, one could perform one epoch of simple SG, or one could assimilate iterates one-by-one and compute  $g_k$  using only the gradients available up to that point.

One important drawback of Algorithm 5.2 is the need to store  $n$  stochastic gradient vectors, which would be prohibitive in many large-scale applications. Note, however, that if the component functions are of the form  $f_i(w_k) = \hat{f}(x_i^T w_k)$ , then

$$\nabla f_i(w_k) = \hat{f}'(x_i^T w_k) x_i.$$

That is, when the feature vectors  $\{x_i\}$  are already available in storage, one need only store the scalar  $\hat{f}'(x_i^T w_k)$  in order to construct  $\nabla f_i(w_k)$  at a later iteration. Such a functional form of  $f_i$  occurs in logistic and least squares regression.

Algorithm 5.2 has its origins in the *stochastic average gradient* (SAG) algorithm [139, 90], where the stochastic direction is defined as

$$(5.15) \quad g_k \leftarrow \frac{1}{n} \left( \nabla f_j(w_k) - \nabla f_j(w_{[j]}) + \sum_{i=1}^n \nabla f_i(w_{[i]}) \right).$$

Although this  $g_k$  is not an unbiased estimator of  $\nabla R_n(w_k)$ , the method enjoys a linear rate of convergence. One finds that the SAG algorithm is a randomized version of the incremental aggregated gradient (IAG) method proposed in [17] and analyzed in [72], where the index  $j$  of the component function updated at every iteration is chosen cyclically. Interestingly, randomizing this choice yields good practical benefits.

**5.3.3. Commentary.** Although the gradient aggregation methods discussed in this section enjoy a faster rate of convergence than SG (i.e., linear vs. sublinear), they should not be regarded as clearly superior to SG. After all, following similar analysis as in section 4.4, the computing time for SG can be shown to be  $T(n, \epsilon) \sim \kappa^2/\epsilon$  with  $\kappa := L/c$ . (In fact, a computing time of  $\kappa/\epsilon$  is often observed in practice.) On the other hand, the computing times for SVRG, SAGA, and SAG are

$$\mathcal{T}(n, \epsilon) \sim (n + \kappa) \log(1/\epsilon),$$

which grows with the number of examples  $n$ . Thus, following similar analysis as in section 4.4, one finds that, for very large  $n$ , gradient aggregation methods are comparable to batch algorithms and therefore cannot beat SG in this regime. For example, if  $\kappa$  is close to 1, then SG is clearly more efficient, since within a single epoch it reaches the optimal testing error [25]. On the other hand, there exists a regime with  $\kappa \gg n$  in which gradient aggregation methods may be superior, and perhaps even easier to tune. At present, it is not known how useful gradient aggregation methods will prove to be in the future of large-scale machine learning. That being said, they certainly represent a class of optimization methods of interest due to their clever use of past information.

**5.4. Iterate Averaging Methods.** Since its inception, it has been observed that SG generates noisy iterate sequences that tend to oscillate around minimizers during the optimization process. Hence, a natural idea is to compute a corresponding sequence of *iterate averages* that would automatically possess less noisy behavior. Specifically, for minimizing a continuously differentiable  $F$  with unbiased gradient estimates, the idea is to employ the iteration

$$(5.16) \quad \begin{aligned} w_{k+1} &\leftarrow w_k - \alpha_k g(w_k, \xi_k) \\ \text{and } \tilde{w}_{k+1} &\leftarrow \frac{1}{k+1} \sum_{j=1}^{k+1} w_j, \end{aligned}$$

where the averaged sequence  $\{\tilde{w}_k\}$  has no effect on the computation of the SG iterate sequence  $\{w_k\}$ . Early hopes were that this auxiliary averaged sequence might possess better convergence properties than the SG iterates themselves. However, such improved behavior was found to be elusive when using classical stepsize sequences that diminished with a rate of  $\mathcal{O}(1/k)$  [124].

A fundamental advancement in the use of iterate averaging came with the work of Polyak [125], which was subsequently advanced with Juditsky [126]; see also the work of Ruppert [136] and Nemirovski and Yudin [106]. Here, the idea remains to employ the iteration (5.16), but with stepsizes diminishing at a slower rate of  $\mathcal{O}(1/(k^a))$  for some  $a \in (\frac{1}{2}, 1)$ . When minimizing strongly convex objectives, it follows from this choice that  $\mathbb{E}[\|w_k - w_*\|_2^2] = \mathcal{O}(1/(k^a))$ , while  $\mathbb{E}[\|\tilde{w}_k - w_*\|_2^2] = \mathcal{O}(1/k)$ . What is interesting, however, is that in certain cases this combination of *long steps* and *averaging* yields an optimal constant in the latter rate (for the iterate averages) in the sense that no rescaling of the steps—through multiplication with a positive definite matrix (see (6.2) in the next section)—can improve the asymptotic rate or constant. This shows that, due to averaging, the adverse effects caused by ill-conditioning disappear. (In this respect, the effect of averaging has been characterized as being similar to that of using a second-order method in which the Hessian approximations approach the Hessian of the objective at the minimizer [125, 22]; see section 6.2.1.) This asymptotic behavior is difficult to achieve in practice, but is possible in some circumstances with careful selection of the stepsize sequence [161].

Iterate averaging has since been incorporated into various schemes in order to allow longer steps while maintaining desired rates of convergence. Examples include the *robust SA* and *mirror descent SA* methods presented in [105], as well as Nesterov's *primal-dual averaging* method proposed in [109, section 6]. This latter method is notable for this section as it employs gradient aggregation and yields an  $\mathcal{O}(1/k)$  rate of convergence for the averaged iterate sequence.

**6. Second-Order Methods.** In section 5, we looked beyond classical SG to methods that are less affected by noise in the stochastic directions. Another manner in which one can move beyond classical SG is to address the adverse effects of high non-linearity and ill-conditioning of the objective function through the use of second-order information. As we shall see, these methods improve convergence rates of batch methods or the constants involved in the sublinear convergence rate of stochastic methods.

A common way to motivate second-order algorithms is to observe that first-order methods, such as SG and the full gradient method, are not *scale invariant*. Consider, for example, the full gradient iteration for minimizing a continuously differentiable function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$ , namely,

$$(6.1) \quad w_{k+1} \leftarrow w_k - \alpha_k \nabla F(w_k).$$

An alternative iteration is obtained by applying a full gradient approach after a linear transformation of variables, i.e., by considering  $\min_{\tilde{w}} F(B\tilde{w})$  for some symmetric positive definite matrix  $B$ . The full gradient iteration for this problem has the form

$$\tilde{w}_{k+1} \leftarrow \tilde{w}_k - \alpha_k B \nabla F(B\tilde{w}_k),$$

which, after scaling by  $B$  and defining  $\{w_k\} := \{B\tilde{w}_k\}$ , corresponds to the iteration

$$(6.2) \quad w_{k+1} \leftarrow w_k - \alpha_k B^2 \nabla F(w_k).$$

Comparing (6.2) with (6.1), it is clear that the behavior of the algorithm will be different under this change of variables. For instance, when  $F$  is a strongly convex quadratic with unique minimizer  $w_*$ , the full gradient method (6.1) generally requires many iterations to approach the minimizer, but from any initial point  $w_1$  the iteration (6.2) with  $B = (\nabla^2 F(w_1))^{-1/2}$  and  $\alpha_1 = 1$  will yield  $w_2 = w_*$ . These latter choices correspond to a single iteration of *Newton's method* [52]. In general, it is natural to seek transformations that perform well in theory and in practice.

Another motivation for second-order algorithms comes from the observation that each iteration of the form (6.1) or (6.2) chooses the subsequent iterate by first computing the minimizer of a second-order Taylor series approximation  $q_k : \mathbb{R}^d \rightarrow \mathbb{R}$  of  $F$  at  $w_k$ , which has the form

$$(6.3) \quad q_k(w) = F(w_k) + \nabla F(w_k)^T (w - w_k) + \frac{1}{2} (w - w_k)^T B^{-2} (w - w_k).$$

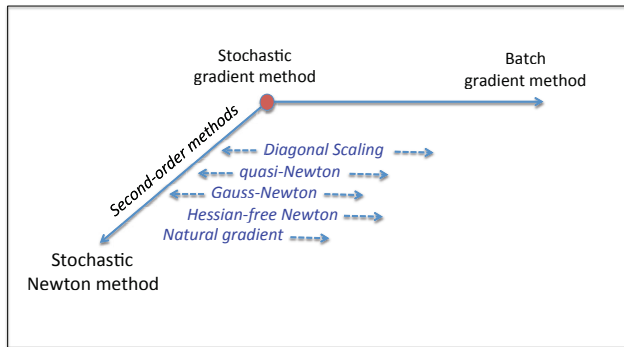
The full gradient iteration corresponds to  $B^{-2} = I$ , while Newton's method corresponds to  $B^{-2} = \nabla^2 F(w_k)$  (assuming this Hessian is positive definite). Thus, in general, a full gradient iteration works with a model that is only first-order accurate, while Newton's method applies successive local rescalings based on minimizing an exact second-order Taylor model of  $F$  at each iterate.

Deterministic (i.e., batch) methods are known to benefit from the use of second-order information; e.g., Newton's method achieves a quadratic rate of convergence if  $w_1$  is sufficiently close to a strong minimizer [52]. On the other hand, stochastic methods like the SG method in section 4 cannot achieve a convergence rate that is faster than sublinear, regardless of the choice of  $B$ ; see [1, 104]. (More on this in section 6.2.1.) Therefore, it is natural to ask, can there be a benefit to incorporating second-order information in stochastic methods? We address this question throughout this section by showing that the careful use of successive rescalings based on (approximate) second-order derivatives can be beneficial between the stochastic and batch regimes.



We begin this section by considering a *Hessian-free Newton* method that employs exact second-order information, but in a judicious manner that exploits the stochastic nature of the objective function. We then describe methods that attempt to mimic the behavior of a Newton algorithm through first-order information computed over sequences of iterates; these include *quasi-Newton*, *Gauss-Newton*, and related algorithms that employ only diagonal rescalings. We also discuss the *natural gradient* method, which defines a search direction in the space of realizable distributions, rather than in the space of the real parameter vector  $w$ . While Newton's method is invariant to linear transformations of the variables, the natural gradient method is invariant with respect to more general invertible transformations.

We depict the methods of interest in this section on the downward axis illustrated in Figure 6.1. We use double-sided arrows for the methods that can be effective throughout the spectrum between the stochastic and batch regimes. Single-sided arrows are used for those methods that one might consider to be effective only with at least a moderate batch size in the stochastic gradient estimates. We explain these distinctions as we describe the methods.



**Fig. 6.1** View of the schematic from Figure 3.3 with a focus on second-order methods. The dotted arrows indicate the effective regime of each method: the first three methods can employ minibatches of any size, whereas the last two methods are efficient only for moderate-to-large minibatch sizes.

**6.1. Hessian-Free Inexact Newton Methods.** Due to its scale invariance properties and its ability to achieve a quadratic rate of convergence in the neighborhood of a strong local minimizer, Newton's method represents an ideal in terms of optimization algorithms. It does not scale well with the dimension  $d$  of the optimization problem, but there are variants that can scale well while also being able to deal with nonconvexity.

When minimizing a twice-continuously differentiable  $F$ , a Newton iteration is

$$(6.4a) \quad w_{k+1} \leftarrow w_k + \alpha_k s_k,$$

$$(6.4b) \quad \text{where } s_k \text{ satisfies } \nabla^2 F(w_k) s_k = -\nabla F(w_k).$$

This iteration demands much in terms of computation and storage. However, rather than solve the Newton system (6.4b) exactly through matrix factorization techniques,

one can instead only solve it inexactly through an iterative approach such as the conjugate gradient (CG) method [69]. By ensuring that the linear solves are accurate enough, such an *inexact Newton-CG* method can enjoy a superlinear rate of convergence [47].

In fact, the computational benefits of inexact Newton-CG go beyond its ability to maintain classical convergence rate guarantees. Like many iterative linear system techniques, CG applied to (6.4b) does not require access to the Hessian itself, only Hessian-vector products [121]. It is in this sense that such a method may be called *Hessian-free*. This is ideal when such products can be coded directly without having to form an explicit Hessian, as Example 6.1 below demonstrates. Each product is at least as expensive as a gradient evaluation, but as long as the number of products—one per CG iteration—is not too large, the improved rate of convergence can compensate for the extra per-iteration work required over a simple full gradient method.

**EXAMPLE 6.1.** Consider the function of the parameter vector  $w = (w_1, w_2)$  given by  $F(w) = \exp(w_1 w_2)$ . Let us define, for any  $d \in \mathbb{R}^2$ , the function

$$\phi(w; d) = \nabla F(w)^T d = w_2 \exp(w_1 w_2) d_1 + w_1 \exp(w_1 w_2) d_2.$$

Computing the gradient of  $\phi$  with respect to  $w$ , we have

$$\nabla_w \phi(w; d) = \nabla^2 F(w) d = \begin{bmatrix} w_2^2 \exp(w_1 w_2) d_1 + (\exp(w_1 w_2) + w_1 w_2 \exp(w_1 w_2)) d_2 \\ (\exp(w_1 w_2) + w_1 w_2 \exp(w_1 w_2)) d_1 + w_1^2 \exp(w_1 w_2) d_2 \end{bmatrix}.$$

We have thus obtained, for any  $d \in \mathbb{R}^2$ , a formula for computing  $\nabla^2 F(w) d$  that does not require  $\nabla^2 F(w)$  explicitly. Note that by storing the scalars  $w_1 w_2$  and  $\exp(w_1 w_2)$  from the evaluation of  $F$ , the additional costs of computing the gradient-vector and Hessian-vector products are small.

The idea illustrated in this example can be applied in general; e.g., see also Example 6.2 below. For a smooth objective function  $F$ , one can compute  $\nabla^2 F(w) d$  at a cost that is a small multiple of the cost of evaluating  $\nabla F$ , and without forming the Hessian, which would require  $\mathcal{O}(d^2)$  storage. The savings in computation come at the expense of the storage of some additional quantities, as explained in Example 6.1.

In machine learning applications, including those involving multinomial logistic regression and DNNs, Hessian-vector products can be computed in this manner, and an inexact Newton-CG method can be applied. A concern is that, in certain cases, the cost of the CG iterations may render such a method uncompetitive with alternative approaches, such as an SG method or a limited memory BFGS method (see section 6.2), which have small computational overhead. Interestingly, however, the structure of the risk measures (3.1) and (3.2) can be exploited so that the resulting method has lighter computational overheads, as described next.

**6.1.1. Subsampled Hessian-Free Newton Methods.** The motivation for the method we now describe stems from the observation that, in inexact Newton methods, the Hessian matrix need not be as accurate as the gradient to yield an effective iteration. Translated into the context of large-scale machine learning applications, this means that the iteration is more tolerant to noise in the Hessian estimate than it is to noise in the gradient estimate.

Based on this idea, the technique we state here employs a smaller sample for defining the Hessian than for the stochastic gradient estimate. Following similar

notation as introduced in section 5.2, let the stochastic gradient estimate be

$$\nabla f_{\mathcal{S}_k}(w_k; \xi_k) = \frac{1}{|\mathcal{S}_k|} \sum_{i \in \mathcal{S}_k} \nabla f(w_k; \xi_{k,i})$$

and let the stochastic Hessian estimate be

$$(6.5) \quad \nabla^2 f_{\mathcal{S}_k^H}(w_k; \xi_k^H) = \frac{1}{|\mathcal{S}_k^H|} \sum_{i \in \mathcal{S}_k^H} \nabla^2 f(w_k; \xi_{k,i}),$$

where  $\mathcal{S}_k^H$  is conditionally uncorrelated with  $\mathcal{S}_k$  given  $w_k$ . If one chooses the *sub-sample* size  $|\mathcal{S}_k^H|$  small enough, then the cost of each product involving the Hessian approximation can be reduced significantly, thus reducing the cost of each CG iteration. On the other hand, one should choose  $|\mathcal{S}_k^H|$  large enough that the curvature information captured through the Hessian-vector products is productive. If done appropriately, *Hessian subsampling* is robust and effective [2, 28, 122, 132]. An inexact Newton method that incorporates this technique is outlined in Algorithm 6.1. The algorithm is stated with a backtracking (Armijo) line search [114], though other stepsize selection techniques could be considered as well.

---

**Algorithm 6.1** Subsampled Hessian-Free Inexact Newton Method

---

- 1: Choose an initial iterate  $w_1$ .
- 2: Choose constants  $\rho \in (0, 1)$ ,  $\gamma \in (0, 1)$ ,  $\eta \in (0, 1)$ , and  $\max_{cg} \in \mathbb{N}$ .
- 3: **for**  $k = 1, 2, \dots$  **do**
- 4:   Generate realizations of  $\xi_k$  and  $\xi_k^H$  corresponding to  $\mathcal{S}_k$  and  $\mathcal{S}_k^H$ .
- 5:   Compute  $s_k$  by applying Hessian-free CG to solve

$$(6.6) \quad \nabla^2 f_{\mathcal{S}_k^H}(w_k; \xi_k^H) s = -\nabla f_{\mathcal{S}_k}(w_k; \xi_k)$$

until  $\max_{cg}$  iterations have been performed or a trial solution yields

$$\|r_k\|_2 := \|\nabla^2 f_{\mathcal{S}_k^H}(w_k; \xi_k^H) s + \nabla f_{\mathcal{S}_k}(w_k; \xi_k)\|_2 \leq \rho \|\nabla f_{\mathcal{S}_k}(w_k; \xi_k)\|_2.$$

- 6:   Set  $w_{k+1} \leftarrow w_k + \alpha_k s_k$ , where  $\alpha_k \in \{\gamma^0, \gamma^1, \gamma^2, \dots\}$  is the largest element with

$$(6.7) \quad f_{\mathcal{S}_k}(w_{k+1}; \xi_k) \leq f_{\mathcal{S}_k}(w_k; \xi_k) + \eta \alpha_k \nabla f_{\mathcal{S}_k}(w_k; \xi_k)^T s_k.$$

- 7: **end for**
- 

As previously mentioned, the (subsampled) Hessian-vector products required in Algorithm 6.1 can be computed efficiently in the context of many machine learning applications. For instance, one such case in the following.

**EXAMPLE 6.2.** Consider a binary classification problem where the training function is given by the logistic loss with an  $\ell_2$ -norm regularization parameterized by  $\lambda > 0$ :

$$(6.8) \quad R_n(w) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \frac{\lambda}{2} \|w\|^2.$$

A (subsampled) Hessian-vector product can be computed efficiently by observing that

$$\nabla^2 f_{\mathcal{S}_k^H}(w_k; \xi_k^H) d = \frac{1}{|\mathcal{S}_k^H|} \sum_{i \in \mathcal{S}_k^H} \frac{\exp(-y_i w^T x_i)}{(1 + \exp(-y_i w^T x_i))^2} (x_i^T d) x_i + \lambda d.$$

To quantify the step computation cost in an inexact Newton-CG framework such as Algorithm 6.1, let  $g_{cost}$  be the cost of computing a gradient estimate  $\nabla f_{\mathcal{S}_k}(w_k; \xi_k)$  and let  $\text{factor} \times g_{cost}$  denote the cost of one Hessian-vector product. If the maximum number of CG iterations,  $\max_{cg}$ , is performed for every outer iteration, then the step computation cost in Algorithm 6.1 is

$$\max_{cg} \times \text{factor} \times g_{cost} + g_{cost}.$$

In a deterministic inexact Newton-CG method for minimizing the empirical risk  $R_n$ , i.e., when  $|\mathcal{S}_k^H| = |\mathcal{S}_k| = n$  for all  $k \in \mathbb{N}$ , the factor is at least 1 and  $\max_{cg}$  would typically be chosen as 5, 20, or more, leading to an iteration that is many times the cost of an SG iteration. However, in a stochastic framework using Hessian subsampling, the factor can be chosen to be sufficiently small such that  $\max_{cg} \times \text{factor} \approx 1$ , leading to a per-iteration cost proportional to that of SG.

Implicit in this discussion is the assumption that the gradient sample size  $|\mathcal{S}_k|$  is large enough that taking subsamples for the Hessian estimate is sensible. If, by contrast, the algorithm were to operate in the stochastic regime of SG, where  $|\mathcal{S}_k|$  is small and gradients are very noisy, then it may be necessary to choose  $|\mathcal{S}_k^H| > |\mathcal{S}_k|$  so that Hessian approximations do not corrupt the step. In such circumstances, the method would be far less attractive than SG. Therefore, the subsampled Hessian-free Newton method outlined here is only recommended when  $\mathcal{S}_k$  is large. This is why, in Figure 6.1, the Hessian-free Newton method is illustrated only with an arrow to the right, i.e., in the direction of larger sample sizes.

Convergence of Algorithm 6.1 is easy to establish when minimizing a strongly convex empirical risk measure  $F = R_n$  when  $\mathcal{S}_k \leftarrow \{1, \dots, n\}$  for all  $k \in \mathbb{N}$ , i.e., when full gradients are always used. In this case, a benefit of employing CG to solve (6.6) is that it immediately improves upon the direction employed in a steepest descent iteration. Specifically, when initialized at zero, it produces in its first iteration a scalar multiple of the steepest descent direction  $-\nabla F(w_k)$ , and further iterations monotonically improve upon this step (in terms of minimizing a quadratic model of the form in (6.3)) until the Newton step is obtained, which is achieved in at most  $d$  iterations of CG (in exact arithmetic). Therefore, by using any number of CG iterations, convergence can be established using standard techniques to choose the stepsize  $\alpha_k$  [114]. When exact Hessians are also used, the rate of convergence can be controlled through the accuracy with which the systems (6.4b) are solved. Defining  $r_k := \nabla^2 F(w_k)s_k + \nabla F(w_k)$  for all  $k \in \mathbb{N}$ , the iteration can enjoy a linear, superlinear, or quadratic rate of convergence by controlling  $\|r_k\|_2$ , where for the superlinear rates one must have  $\{\|r_k\|_2 / \|\nabla F(w_k)\|_2\} \rightarrow 0$  [47].

When the Hessians are subsampled (i.e.,  $\mathcal{S}_k^H \subset \mathcal{S}_k$  for all  $k \in \mathbb{N}$ ), it has not been shown that the rate of convergence is faster than linear; nevertheless, the reduction in the number of iterations required to produce a good approximate solution is often significantly lower than if no Hessian information is used in the algorithm.

**6.1.2. Dealing with Nonconvexity.** Hessian-free Newton methods are routinely applied for the solution of nonconvex problems. In such cases, it is common to employ a *trust region* [37] instead of a line search and to introduce an additional condition in step 5 of Algorithm 6.1: terminate CG if a candidate solution  $s_k$  is a direction of negative curvature, i.e.,  $s_k^T \nabla^2 f_{\mathcal{S}_k^H}(w_k; \xi_k^H) s_k < 0$  [149]. A number of more sophisticated strategies have been proposed throughout the years with some success, but none have proved to be totally satisfactory or universally accepted.

Instead of coping with indefiniteness, one can focus on strategies for ensuring positive (semi)definite Hessian approximations. One of the most attractive ways of doing this in the context of machine learning is to employ a (subsampled) Gauss–Newton approximation to the Hessian, which is a matrix of the form

$$(6.9) \quad G_{S_k^H}(w_k; \xi_k^H) = \frac{1}{|S_k^H|} \sum_{i \in S_k^H} J_h(w_k, \xi_{k,i})^T H_\ell(w_k, \xi_{k,i}) J_h(w_k, \xi_{k,i}).$$

Here, the matrix  $J_h$  captures the stochastic gradient information for the prediction function  $h(x; w)$ , whereas the matrix  $H_\ell$  only captures the second-order information for the (convex) loss function  $\ell(h, y)$ ; see section 6.3 for a detailed explanation. As before, one can directly code the product of this matrix times a vector without forming the matrix components explicitly. This approach has been applied successfully in the training of DNNs [12, 101].

We mention in passing that there has been much discussion about the role that negative curvature and saddle points play in the optimization of DNNs; see, e.g., [44, 70, 35]. Numerical tests designed to probe the geometry of the objective function in the neighborhood of a minimizer when training a DNN have shown the presence of negative curvature. It is believed that the inherent stochasticity of the SG method allows it to navigate efficiently through this complex landscape, but it is not known whether classical techniques to avoid approaching saddle points will prove to be successful for either batch or stochastic methods.

**6.2. Stochastic Quasi-Newton Methods.** One of the most important developments in the field of nonlinear optimization came with the advent of *quasi-Newton* methods. These methods construct approximations to the Hessian using only gradient information, and are applicable for convex and nonconvex problems. Versions that scale well with the number of variables, such as *limited memory* methods, have proved to be effective in a wide range of applications where the number of variables can be in the millions. It is therefore natural to ask whether quasi-Newton methods can be extended to the stochastic setting arising in machine learning. Before we embark on this discussion, let us review the basic principles underlying quasi-Newton methods, focusing on the most popular scheme, namely, BFGS [27, 60, 68, 146].

In the spirit of Newton’s method (6.4), the BFGS iteration for minimizing a twice continuously differentiable function  $F$  has the form

$$(6.10) \quad w_{k+1} \leftarrow w_k - \alpha_k H_k \nabla F(w_k),$$

where  $H_k$  is a symmetric positive definite approximation of  $(\nabla^2 F(w_k))^{-1}$ . This form of the iteration is consistent with (6.4), but the distinguishing feature of a quasi-Newton scheme is that the sequence  $\{H_k\}$  is updated dynamically by the algorithm rather than through a second-order derivative computation at each iterate. Specifically, in the BFGS method, the new inverse Hessian approximation is obtained by defining the iterate and gradient displacements

$$s_k := w_{k+1} - w_k \quad \text{and} \quad v_k := \nabla F(w_{k+1}) - \nabla F(w_k),$$

then setting

$$(6.11) \quad H_{k+1} \leftarrow \left( I - \frac{v_k s_k^T}{s_k^T v_k} \right)^T H_k \left( I - \frac{v_k s_k^T}{s_k^T v_k} \right) + \frac{s_k s_k^T}{s_k^T v_k}.$$

One important aspect of this update is that it ensures that the secant equation  $H_{k+1}^{-1}s_k = v_k$  holds, meaning that a second-order Taylor expansion is satisfied along the most recent displacement (though not necessarily along other directions).

A remarkable fact about the BFGS iteration (6.10)–(6.11) is that it enjoys a local superlinear rate of convergence [51], and this with only first-order information and without the need for any linear system solves (which are required by Newton's method for it to be quadratically convergent). However, a number of issues need to be addressed to have an effective method in practice. For one thing, the update (6.11) yields dense matrices, even when the exact Hessians are sparse, restricting its use to small and midsize problems. A common solution for this restriction is to employ a *limited memory* scheme, leading to a method such as L-BFGS [97, 113]. A key feature of this approach is that the matrices  $\{H_k\}$  need not be formed explicitly; instead, each product of the form  $H_k \nabla F(w_k)$  can be computed using a formula that only requires recent elements of the sequence of displacement pairs  $\{(s_k, v_k)\}$  that have been saved in storage. Such an approach incurs per-iteration costs of order  $\mathcal{O}(d)$ , and delivers practical performance that is significantly better than a full gradient iteration, though the rate of convergence is only provably linear.

**6.2.1. Deterministic to Stochastic.** Let us consider the extension of a quasi-Newton approach from the deterministic to the stochastic setting arising in machine learning. The iteration now takes the form

$$(6.12) \quad w_{k+1} \leftarrow w_k - \alpha_k H_k g(w_k, \xi_k).$$

Since we are interested in large-scale problems, we assume that (6.12) implements an L-BFGS scheme, which avoids the explicit construction of  $H_k$ . A number of questions arise when considering (6.12). We list them now with some proposed solutions.

**Theoretical Limitations.** The convergence rate of a stochastic iteration such as (6.12) cannot be faster than sublinear [1]. Given that SG also has a sublinear rate of convergence, what benefit, if any, could come from incorporating  $H_k$  into (6.12)? This is an important question. As it happens, one can see a benefit of  $H_k$  in terms of the *constant* that appears in the sublinear rate. Recall that for the SG method (Algorithm 4.1), the constant depends on  $L/c$ , which in turn depends on the conditioning of  $\{\nabla^2 F(w_k)\}$ . This is typical of first-order methods. In contrast, one can show [25] that if the sequence of Hessian approximations in (6.12) satisfies  $\{H_k\} \rightarrow \nabla^2 F(w_*)^{-1}$ , then the constant is independent of the conditioning of the Hessian. Although constructing Hessian approximations with this property might not be viable in practice, this fact suggests that stochastic quasi-Newton methods could be better equipped to cope with ill-conditioning than SG.

**Additional Per-Iteration Costs.** The SG iteration is very inexpensive, requiring only the evaluation of  $g(w_k, \xi_k)$ . The iteration (6.12), on the other hand, also requires the product  $H_k g(w_k, \xi_k)$ , which is known to require  $4md$  operations, where  $m$  is the *memory* in the L-BFGS updating scheme. Assuming for concreteness that the cost of evaluating  $g(w_k, \xi_k)$  is exactly  $d$  operations (using only one sample) and that the memory parameter is set to the typical value of  $m = 5$ , one finds that the stochastic quasi-Newton method is 20 times more expensive than SG. Can the iteration (6.12) yield fast enough progress as to offset this additional per-iteration cost? To address this question, one need only observe that the calculation just mentioned focuses on the gradient  $g(w_k, \xi_k)$  being based on a single sample. However, when employing minibatch gradient estimates, the additional cost of the iteration (6.12) is only

marginal. (Minibatches of size 256 are common in practice.) The use of minibatches may therefore be considered essential when one contemplates the use of a stochastic quasi-Newton method. This minibatch need not be large, as in the Hessian-free Newton method discussed in the previous section, but it should not be less than, say, 20 or 50, in light of the additional costs of computing the matrix-vector products.

**Conditioning of the Scaling Matrices.** The BFGS formula (6.11) for updating  $H_k$  involves differences in gradient estimates computed in consecutive iterations. In stochastic settings, the gradients  $\{g(w_k, \xi_k)\}$  are noisy estimates of  $\{\nabla F(w_k)\}$ . This can cause the updating process to yield poor curvature estimates, which may have a detrimental rather than beneficial effect on the quality of the computed steps. Since BFGS, like all quasi-Newton schemes, is an *overwriting* process, the effects of even a single bad update may linger for numerous iterations. How could such effects be avoided in the stochastic regime? There have been various proposals to avoid differencing noisy gradient estimates. One possibility is to employ the same sample when computing gradient differences [18, 142]. An alternative approach that allows greater freedom in the choice of the stochastic gradient is to *decouple* the step computation (6.12) and the Hessian update. In this manner, one can employ a larger sample, if necessary, when computing the gradient displacement vector. We discuss these ideas further in section 6.2.2.

It is worthwhile to note that if the gradient estimate  $g(w_k, \xi_k)$  does not have high variance, then standard BFGS updating can be applied without concern. Therefore, in the rest of this section, we focus on algorithms that employ noisy gradient estimates in the step computation (6.12). This means, e.g., that we are not considering the potential to tie the method to the noise reduction techniques described in section 5, though such an idea is natural and could be effective in practice.

**6.2.2. Algorithms.** A straightforward adaptation of L-BFGS only involves the replacement of deterministic gradients with stochastic gradients throughout the iterative process. The displacement pairs might then be defined as

$$(6.13) \quad s_k := w_{k+1} - w_k \quad \text{and} \quad v_k := \nabla f_{S_k}(w_{k+1}, \xi_k) - \nabla f_{S_k}(w_k, \xi_k).$$

Note the use of the same seed  $\xi_k$  in the two gradient estimates, in order to address the issues related to noise mentioned above. If each  $f_i$  is strongly convex, then  $s_k^T v_k > 0$ , and positive definiteness of the updates is also maintained. Such an approach is sometimes referred to as *online L-BFGS* [142, 103].

One disadvantage of this method is the need to compute two, as opposed to only one, gradient estimates per iteration: one to compute the gradient displacement (namely,  $g(w_{k+1}, \xi_k)$ ) and another (namely,  $g(w_{k+1}, \xi_{k+1})$ ) to compute the subsequent step. This is not too onerous, at least as long as the per-iteration improvement outweighs the extra per-iteration cost. A more worrisome feature is that updating the inverse Hessian approximations with *every* step may not be warranted, and may even be detrimental when the gradient displacement is based on a small sample, since it could easily represent a poor approximation of the action of the true Hessian of  $F$ .

An alternative strategy, which might better represent the action of the true Hessian even when  $g(w_k, \xi_k)$  has high variance, is to employ an alternative  $v_k$ . In particular, since  $\nabla F(w_{k+1}) - \nabla F(w_k) \approx \nabla^2 F(w_k)(w_{k+1} - w_k)$ , one could define

$$(6.14) \quad v_k := \nabla^2 f_{S_k^H}(w_k; \xi_k^H) s_k,$$

where  $\nabla^2 f_{S_k^H}(w_k; \xi_k^H)$  is a subsampled Hessian and  $|S_k^H|$  is large enough to provide

useful curvature information. As in the case of Hessian-free Newton from section 6.1.1, the product (6.14) can be performed without explicitly constructing  $\nabla^2 f_{\mathcal{S}_k^H}(w_k; \xi_k^H)$ .

Regardless of the definition of  $v_k$ , when  $|\mathcal{S}_k^H|$  is much larger than  $|\mathcal{S}_k|$ , the cost of quasi-Newton updating is excessive due to the cost of computing  $v_k$ . To address this issue, the computation of  $v_k$  can be performed only after a sequence of iterations, to amortize costs. This leads to the idea of decoupling the step computation from the quasi-Newton update. This approach, which we refer to for convenience as SQN, performs a sequence of iterations of (6.12) with  $H_k$  fixed, then computes a new displacement pair  $(s_k, v_k)$  with  $s_k$  defined as in (6.13) and  $v_k$  set using one of the strategies outlined above. This pair replaces one of the old pairs in storage, which in turn defines the limited memory BFGS step.

To formalize all of these alternatives, we state the general stochastic quasi-Newton method presented as Algorithm 6.2, with some notation borrowed from Algorithm 6.1. In the method, the step computation is based on a collection of  $m$  displacement pairs  $\mathcal{P} = \{s_j, v_j\}$  in storage and the current stochastic gradient  $\nabla f_{\mathcal{S}_k}(w_k; \xi_k)$ , where the matrix-vector product in (6.12) can be computed through a *two-loop recursion* [113, 114]. To demonstrate the generality of the method, we note that the *online L-BFGS* method sets  $\mathcal{S}_k^H \leftarrow \mathcal{S}_k$  and `update pairs = true` in every iteration. In SQN using (6.14), on the other hand,  $|\mathcal{S}_k^H|$  should be chosen larger than  $|\mathcal{S}_k|$  and one sets `update pairs = true` only every, say, 10 or 20 iterations.

---

**Algorithm 6.2** Stochastic Quasi-Newton Framework
 

---

```

1: Choose an initial iterate  $w_1$  and initialize  $\mathcal{P} \leftarrow \emptyset$ .
2: Choose a constant  $m \in \mathbb{N}$ .
3: Choose a stepsize sequence  $\{\alpha_k\} \subset \mathbb{R}_{++}$ .
4: for  $k = 1, 2, \dots$  do
5:   Generate realizations of  $\xi_k$  and  $\xi_k^H$  corresponding to  $\mathcal{S}_k$  and  $\mathcal{S}_k^H$ .
6:   Compute  $\hat{s}_k = H_k g(w_k, \xi_k)$  using the two-loop recursion based on the set  $\mathcal{P}$ .
7:   Set  $s_k \leftarrow -\alpha_k \hat{s}_k$ .
8:   Set  $w_{k+1} \leftarrow w_k + s_k$ .
9:   if update pairs then
10:    Compute  $s_k$  and  $v_k$  (based on the sample  $\mathcal{S}_k^H$ ).
11:    Add the new displacement pair  $(s_k, v_k)$  to  $\mathcal{P}$ .
12:    If  $|\mathcal{P}| > m$ , then remove eldest pair from  $\mathcal{P}$ .
13:   end if
14: end for

```

---

To guarantee that the BFGS update is well defined, each displacement pair  $(s_j, v_j)$  must satisfy  $s_j^T v_j > 0$ . In deterministic optimization, this issue is commonly addressed by either performing a line search (involving exact gradient computations) or modifying the displacement vectors (e.g., through *damping*) so that  $s_j^T v_j > 0$ , in which case one does ensure that (6.11) maintains positive definite approximations. However, these mechanisms have not been fully developed in the stochastic regime when exact gradient information is unavailable and the gradient displacement vectors are noisy. Simple ways to overcome these difficulties are to replace the Hessian matrix with a Gauss–Newton approximation or to introduce a combination of damping and regularization (say, through the addition of simple positive definite matrices).

There remains much to be explored in terms of stochastic quasi-Newton methods for machine learning applications. Experience has shown that some gains in perfor-



mance can be achieved, but the full potential of the quasi-Newton schemes discussed above (and potentially others) is not yet known.

**6.3. Gauss–Newton Methods.** The Gauss–Newton method is a classical approach for nonlinear least squares, i.e., minimization problems in which the objective function is a sum of squares. This method readily applies for optimization problems arising in machine learning involving a least squares loss function, but the idea generalizes for other popular loss functions as well. The primary advantage of Gauss–Newton is that it constructs an approximation to the Hessian using only first-order information, and this approximation is guaranteed to be positive semidefinite, even when the full Hessian itself may be indefinite. The price to pay for this convenient representation is that it ignores second-order interactions between elements of the parameter vector  $w$ , which might mean a loss of curvature information that could be useful for the optimization process.

**Classical Gauss–Newton.** Let us introduce the classical Gauss–Newton approach by considering a situation in which, for a given input-output pair  $(x, y)$ , the loss incurred by a parameter vector  $w$  is measured via a squared norm discrepancy between  $h(x; w) \in \mathbb{R}^d$  and  $y \in \mathbb{R}^d$ . Representing the input-output pair being chosen randomly via the subscript  $\xi$ , we may thus write

$$f(w; \xi) = \ell(h(x_\xi; w), y_\xi) = \frac{1}{2} \|h(x_\xi; w) - y_\xi\|_2^2.$$

Writing a second-order Taylor series model of this function in the vicinity of parameter vector  $w_k$  would involve its gradient and Hessian at  $w_k$ , and minimizing the resulting model (recall (6.3)) would lead to a Newton iteration. Alternatively, a Gauss–Newton approximation of the function is obtained by making an affine approximation of the prediction function inside the quadratic loss function. Letting  $J_h(\cdot; \xi)$  represent the Jacobian of  $h(x_\xi; \cdot)$  with respect to  $w$ , we have the approximation

$$h(x_\xi; w) \approx h(x_\xi; w_k) + J_h(w_k; \xi)(w - w_k),$$

which leads to

$$\begin{aligned} f(w; \xi) &\approx \frac{1}{2} \|h(x_\xi; w_k) + J_h(w_k; \xi)(w - w_k) - y_\xi\|_2^2 \\ &= \frac{1}{2} \|h(x_\xi; w_k) - y_\xi\|_2^2 + (h(x_\xi; w_k) - y_\xi)^T J_h(w_k; \xi)(w - w_k) \\ &\quad + \frac{1}{2} (w - w_k)^T J_h(w_k; \xi)^T J_h(w_k; \xi)(w - w_k). \end{aligned}$$

In fact, this approximation is similar to a second-order Taylor series model, except that the terms involving the second derivatives of the prediction function  $h$  with respect to the parameter vector have been dropped. The remaining second-order terms are those resulting from the positive curvature of the quadratic loss  $\ell$ . This leads to replacing the subsampled Hessian matrix (6.5) by the Gauss–Newton matrix

$$(6.15) \quad G_{\mathcal{S}_k^H}(w_k; \xi_k^H) = \frac{1}{|\mathcal{S}_k^H|} \sum_{i \in \mathcal{S}_k^H} J_h(w_k; \xi_{k,i})^T J_h(w_k; \xi_{k,i}).$$

Since the Gauss–Newton matrix only differs from the true Hessian by terms that involve the factors  $h(x_\xi; w_k) - y_\xi$ , these two matrices are the same when the loss is equal to zero, i.e., when  $h(x_\xi; w_k) = y_\xi$ .

A challenge in the application of a Gauss–Newton scheme is that the Gauss–Newton matrix is often singular or nearly singular. In practice, this is typically handled by regularizing it by adding to it a positive multiple of the identity matrix. For least squares loss functions, the inexact Hessian-free Newton methods of section 6.1 and the stochastic quasi-Newton methods of section 6.2 with gradient displacement vectors defined as in (6.14) can be applied with (regularized) Gauss–Newton approximations. This has the benefit that the scaling matrices are guaranteed to be positive definite.

The computational cost of the Gauss–Newton method depends on the dimensionality of the prediction function. When the prediction function is scalar-valued, the Jacobian matrix  $J_h$  is a single row whose elements are already being computed as an intermediate step in the computation of the stochastic gradient  $\nabla f(w; \xi)$ . However, this is no longer true when the dimensionality is larger than one since then computing the stochastic gradient vector  $\nabla f(w; \xi)$  does not usually require the explicit computation of all rows of the Jacobian matrix. This happens, for instance, in DNNs when one uses back propagation [134, 135].

**Generalized Gauss–Newton.** Gauss–Newton ideas can also be generalized for other standard loss functions [141]. To illustrate, let us consider a slightly more general situation in which loss is measured by a composition of an arbitrary convex loss function  $\ell(h, y)$  and a prediction function  $h(x; w)$ . Combining the affine approximation of the prediction function  $h(x_\xi; w)$  with a second-order Taylor expansion of the loss function  $\ell$  leads to the generalized Gauss–Newton matrix

$$(6.16) \quad G_{S_k^H}(w_k; \xi_k^H) = \frac{1}{|S_k^H|} \sum_{i \in S_k^H} J_h(w_k; \xi_{k,i})^T H_\ell(w_k; \xi_{k,i}) J_h(w_k; \xi_{k,i})$$

(recall (6.9)), where  $H_\ell(w_k; \xi) = \frac{\partial^2 \ell}{\partial h^2}(h(x_\xi; w_k), y_\xi)$  captures the curvature of the loss function  $\ell$ . This can be seen as a generalization of (6.15) in which  $H_\ell = I$ .

When training a DNN, one may exploit this generalized strategy by redefining  $\ell$  and  $h$  so that as much as possible of the network’s computation is formally performed by  $\ell$  rather than by  $h$ . If this can be done in such a way that convexity of  $\ell$  is maintained, then one can faithfully capture second-order terms for  $\ell$  using the generalized Gauss–Newton scheme. Interestingly, in many useful situations, this strategy gives simpler and more elegant expressions for  $H_\ell$ . For instance, probability estimation problems often reduce to using logarithmic losses of the form  $f(w; \xi) = -\log(h(x_\xi; w))$ . The generalized Gauss–Newton matrix then reduces to

$$(6.17) \quad \begin{aligned} G_{S_k^H}(w_k; \xi_k^H) &= \frac{1}{|S_k^H|} \sum_{i \in S_k^H} J_h(w_k; \xi_{k,i})^T \frac{1}{h(w; \xi_{k,i})^2} J_h(w_k; \xi_{k,i}) \\ &= \frac{1}{|S_k^H|} \sum_{i \in S_k^H} \nabla f(w; \xi_{k,i}) \nabla f(w; \xi_{k,i})^T, \end{aligned}$$

which does not require explicit computation of the Jacobian  $J_h$ .

**6.4. Natural Gradient Method.** We have seen that Newton’s method is invariant to *linear* transformations of the parameter vector  $w$ . By contrast, the natural gradient method [5, 6] aims to be invariant with respect to all differentiable and invertible transformations. The essential idea consists of formulating the gradient descent algorithm in the space of prediction functions rather than specific parameters.

Of course, the actual computation takes place with respect to the parameters, but accounts for the anisotropic relation between the parameters and the decision function. That is, in parameter space, the natural gradient algorithm will move the parameters more quickly along directions that have a small impact on the decision function, and more cautiously along directions that have a large impact on the decision function.

We remark at the outset that many authors [119, 99] propose *quasi-natural-gradient* methods that are strikingly similar to the *quasi-Newton* methods described in section 6.2. The natural gradient approach therefore offers a different justification for these algorithms, one that involves qualitatively different approximations. It should also be noted that research on the design of methods inspired by the natural gradient is ongoing and may lead to markedly different algorithms [33, 76, 100].

**Information Geometry.** In order to directly formulate the gradient descent in the space of prediction functions, we must elucidate the geometry of this space. Amari's work on information geometry [6] demonstrates this for parametric density estimation. The space  $\mathcal{H}$  of prediction functions for such a problem is a family of densities  $h_w(x)$  parametrized by  $w \in \mathcal{W}$  and satisfying the normalization condition

$$\int h_w(x) dx = 1 \quad \text{for all } w \in \mathcal{W}.$$

Assuming sufficient regularity, the derivatives of such densities satisfy the identity

$$(6.18) \quad \forall t > 0, \quad \int \frac{\partial^t h_w(x)}{\partial w^t} dx = \frac{\partial^t}{\partial w^t} \int h_w(x) dx = \frac{\partial^t 1}{\partial w^t} = 0.$$

To elucidate the geometry of the space  $\mathcal{H}$ , we seek to quantify how the density  $h_w$  changes when one adds a small quantity  $\delta w$  to its parameter. We achieve this in a statistically meaningful way by observing the Kullback–Leibler (KL) divergence

$$(6.19) \quad D_{KL}(h_w \| h_{w+\delta w}) = \mathbb{E}_{h_w} \left[ \log \left( \frac{h_w(x)}{h_{w+\delta w}(x)} \right) \right],$$

where  $\mathbb{E}_{h_w}$  denotes the expectation with respect to the distribution  $h_w$ . Note that (6.19) only depends on the values of the two density functions  $h_w$  and  $h_{w+\delta w}$  and therefore is invariant with respect to any invertible transformation of the parameter  $w$ . Approximating the divergence with a second-order Taylor expansion, one obtains

$$\begin{aligned} D_{KL}(h_w \| h_{w+\delta w}) &= \mathbb{E}_{h_w} [\log(h_w(x)) - \log(h_{w+\delta w}(x))] \\ &\approx -\delta w^T \mathbb{E}_{h_w} \left[ \frac{\partial \log(h_w(x))}{\partial w} \right] - \frac{1}{2} \delta w^T \mathbb{E}_{h_w} \left[ \frac{\partial^2 \log(h_w(x))}{\partial w^2} \right] \delta w, \end{aligned}$$

which, after observing that (6.18) implies that the first-order term is null, yields

$$(6.20) \quad D_{KL}(h_w \| h_{w+\delta w}) \approx \frac{1}{2} \delta w^T G(w) \delta w.$$

This is a quadratic form defined by the *Fisher information* matrix

$$(6.21) \quad G(w) := -\mathbb{E}_{h_w} \left[ \frac{\partial^2 \log(h_w(x))}{\partial w^2} \right] = -\mathbb{E}_{h_w} \left[ \left( \frac{\partial \log(h_w(x))}{\partial w} \right) \left( \frac{\partial \log(h_w(x))}{\partial w} \right)^T \right],$$

where the latter equality follows again from (6.18). The second form of  $G(w)$  is often preferred because it makes clear that the Fisher information matrix  $G(w)$  is symmetric and always positive semidefinite, though not necessarily positive definite.

The relation (6.20) means that the KL divergence behaves locally like a norm associated with  $G(w)$ . Therefore, every small region of  $\mathcal{H}$  looks like a small region of a Euclidean space. However, as we traverse larger regions of  $\mathcal{H}$ , we cannot ignore that the matrix  $G(w)$  changes. Such a construction defines a *Riemannian geometry*.<sup>7</sup>

Suppose, for instance, that we move along a smooth path connecting two densities, call them  $h_{w_0}$  and  $h_{w_1}$ . A parametric representation of the path can be given by a differentiable function, for which we define

$$\phi : t \in [0, 1] \mapsto \phi(t) \in \mathcal{W} \quad \text{with} \quad \phi(0) = w_0 \quad \text{and} \quad \phi(1) = w_1.$$

We can compute the length of the path by viewing it as a sequence of infinitesimal segments  $[\phi(t), \phi(t + dt)]$  whose length is given by (6.20), i.e., the total length is

$$D_\phi = \int_0^1 \sqrt{\left(\frac{d\phi}{dt}(t)\right)^T G(\phi(t)) \left(\frac{d\phi}{dt}(t)\right)} dt.$$

An important tool for the study of Riemannian geometries is the characterization of its *geodesics*, i.e., the shortest paths connecting two points. In a Euclidean space, the shortest path between two points is always the straight line segment connecting them. In a Riemannian space, on the other hand, the shortest path between two points can be curved and does not need to be unique. Such considerations are relevant to optimization since every iterative optimization algorithm can be viewed as attempting to follow a particular path connecting the initial point  $w_0$  to the optimum  $w_*$ . In particular, following the shortest path is attractive because it means that the algorithm reaches the optimum after making the fewest number of changes to the optimization variables, hopefully requiring the least amount of computation.

**Natural Gradient.** Let us now assume that the space  $\mathcal{H}$  of prediction functions  $\{h_w : w \in \mathcal{W}\}$  has a Riemannian geometry locally described by an identity of the form (6.20). We seek an algorithm that minimizes a functional  $F : h_w \in \mathcal{H} \mapsto F(h_w) = F(w) \in \mathbb{R}$  and is invariant with respect to differentiable invertible transformations of the parameters represented by the vector  $w$ .

Each iteration of a typical iterative optimization algorithm computes a new iterate  $h_{w_{k+1}}$  on the basis of information pertaining to the current iterate  $h_{w_k}$ . Since we can only expect this information to be valid in a small region surrounding  $h_{w_k}$ , we restrict our attention to algorithms that make a step from  $h_{w_k}$  to  $h_{w_{k+1}}$  of some small length  $\eta_k > 0$ . The number of iterations needed to reach the optimum then depends directly on the length of the path followed by the algorithm, which is desired to be as short as possible. Unfortunately, it is rarely possible to exactly follow a geodesic using only local information. We can, however, formulate the greedy strategy that

$$(6.22) \quad h_{w_{k+1}} = \arg \min_{h \in \mathcal{H}} F(h) \quad \text{s.t.} \quad D(h_{w_k} \| h) \leq \eta_k^2,$$

and use (6.20) to reformulate this problem in terms of the parameters:

$$(6.23) \quad w_{k+1} = \arg \min_{w \in \mathcal{W}} F(w) \quad \text{s.t.} \quad \frac{1}{2}(w - w_k)^T G(w_k) (w - w_k) \leq \eta_k^2.$$

<sup>7</sup>The objective of information geometry [6] is to exploit the Riemannian structure of parametric families of density functions to gain geometrical insights into the fundamental statistical phenomena. The natural gradient algorithm is only a particular aspect of this broader goal [5].

The customary derivation of the natural gradient algorithm handles the constraint in (6.23) using a Lagrangian formulation with Lagrange multiplier  $1/\alpha_k$ . In addition, since  $\eta_k$  is assumed small, it replaces  $F(w)$  in (6.23) by the first-order approximation  $F(w_k) + \nabla F(w_k)^T(w - w_k)$ . These two choices lead to the expression

$$w_{k+1} = \arg \min_{w \in \mathcal{W}} \nabla F(w_k)^T(w - w_k) + \frac{1}{2\alpha_k}(w - w_k)^T G(w_k)(w - w_k),$$

the optimization of the right-hand side of which leads to the natural gradient iteration

$$(6.24) \quad w_{k+1} = w_k - \alpha_k G^{-1}(w_k) \nabla F(w_k).$$

We can also replace  $F(w)$  in (6.23) by a noisy first-order approximation, leading to a stochastic natural gradient iteration where  $\nabla F(w_k)$  in (6.24) is replaced by a stochastic gradient estimate.

Both batch and stochastic versions of (6.24) resemble the quasi-Newton update rules discussed in section 6.2. Instead of multiplying the gradient by the inverse of an approximation of the Hessian (which is not necessarily positive definite), it employs the positive semidefinite matrix  $G(w_k)$  that expresses the local geometry of the space of prediction functions. In principle, this matrix does not even take into account the objective function  $F$ . However, as we shall now describe, one finds that these choices are all closely related in practice.

**Practical Natural Gradient.** Because the discovery of the natural gradient algorithm is closely associated with information geometry, nearly all its applications involve density estimation [5, 33] or conditional probability estimation [119, 76, 99] using objective functions that are closely connected to the KL divergence. Natural gradient in this context is closely related to Fisher's scoring algorithm [118]. For instance, in the case of density estimation, the objective is usually the negative log likelihood

$$F(w) = \frac{1}{n} \sum_{i=1}^n -\log(h_w(x_i)) \approx \text{constant} + D_{KL}(P \| h_w),$$

where  $\{x_1, \dots, x_n\}$  represent independent training samples from an unknown distribution  $P$ . Recalling the expression of the Fisher information matrix (6.21) then clarifies its connection with the Hessian, as one finds that

$$G(w) = -\mathbb{E}_{h_w} \left[ \frac{\partial^2 \log(h_w(x))}{\partial w^2} \right] \quad \text{and} \quad \nabla^2 F(w) = -\mathbb{E}_P \left[ \frac{\partial^2 \log(h_w(x))}{\partial w^2} \right].$$

These two expressions do not coincide in general because the expectations involve different distributions. However, when the natural gradient algorithm approaches the optimum, the parametric density  $h_{w_k}$  ideally approaches the true distribution  $P$ , in which case the Fisher information matrix  $G(w_k)$  approaches the Hessian matrix  $\nabla^2 F(w_k)$ . This means that the natural gradient algorithm and Newton's method perform very similarly as optimality is approached.

Although it is occasionally possible to determine a convenient analytic expression [33, 76], the numerical computation of the Fisher information matrix  $G(w_k)$  in large learning systems is generally very challenging. Moreover, estimating the expectation (6.21) with, say, a Monte Carlo approach is usually prohibitive due to the cost of sampling the current density estimate  $h_{w_k}$ .

Several authors [119, 99] suggest using instead a subset of training examples and computing a quantity of the form

$$\tilde{G}(w_k) = \frac{1}{|S_k|} \sum_{i \in S_k} \left( \frac{\partial \log(h_w(x_i))}{\partial w} \Big|_{w_k} \right) \left( \frac{\partial \log(h_w(x_i))}{\partial w} \Big|_{w_k} \right)^T.$$

Although such algorithms are essentially equivalent to the generalized Gauss–Newton schemes described in section 6.3, the natural gradient perspective comes with an interesting insight into the relationship between the generalized Gauss–Newton matrix (6.17) and the Hessian matrix (6.5). Similar to the equality (6.21), these two matrices would be equal if the expectation was taken with respect to the model distribution  $h_w$  instead of the empirical sample distribution.

**6.5. Methods that Employ Diagonal Scalings.** The methods that we have discussed so far in this section are forced to overcome the fact that when employing an iteration involving an  $\mathbb{R}^d \times \mathbb{R}^d$  scaling matrix, one needs to ensure that the improved per-iteration progress outweighs the added per-iteration cost. We have seen that these added costs can be as little as  $4md$  operations and therefore amount to a moderate multiplicative factor on the cost of each iteration.

A strategy to further reduce this multiplicative factor, while still incorporating second-order-type information, is to restrict attention to *diagonal* or *block-diagonal* scaling matrices. Rather than perform a more general linear transformation through a symmetric positive definite matrix (i.e., corresponding to a scaling and rotation of the direction), the incorporation of a diagonal scaling matrix only has the effect of scaling the individual search direction components. This can be efficiently achieved by multiplying each coefficient of the gradient vector by the corresponding diagonal term of the scaling matrix, or, when the prediction function is linear, by adaptively renormalizing the input pattern coefficients [133].

**Computing Diagonal Curvature.** A first family of algorithms that we consider directly computes the diagonal terms of the Hessian or Gauss–Newton matrix, then divides each coefficient of the stochastic gradient vector  $g(w_k, \xi_k)$  by the corresponding diagonal term. Since the computation overhead of this operation is very small, it becomes important to make sure that the estimation of the diagonal terms of the curvature matrix is very efficient.

For instance, in the context of DNNs, [12] describes a back-propagation algorithm to efficiently compute the diagonal terms of the squared Jacobian matrix  $J_h(w_k; \xi_k)^T J_h(w_k; \xi_k)$  that appears in the expression of the Gauss–Newton matrix (6.15). Each iteration of the proposed algorithm picks a training example, computes the stochastic gradient  $g(w_k, \xi_k)$ , updates a running estimate of the diagonal coefficients of the Gauss–Newton matrix by

$$[G_k]_i = (1 - \lambda)[G_{k-1}]_i + \lambda[J_h(w_k; \xi_k)^T J_h(w_k; \xi_k)]_{ii} \quad \text{for some } 0 < \lambda < 1,$$

then performs the scaled stochastic weight update

$$[w_{k+1}]_i = [w_k]_i - \left( \frac{\alpha}{[G_k]_i + \mu} \right) [g(w_k, \xi_k)]_i.$$

The small regularization constant  $\mu > 0$  is introduced to handle situations where the Gauss–Newton matrix is singular or nearly singular. Since the computation of

the diagonal of the squared Jacobian has a cost that is comparable to the cost of the computation of the stochastic gradient, each iteration of this algorithm is roughly twice as expensive as a first-order SG iteration. The experience described in [12] shows that improvement in per-iteration progress can be sufficient to modestly outperform a well-tuned SG algorithm.

After describing this algorithm in later work [89, section 9.1], the authors make two comments that illustrate well how this algorithm was used in practice. They first observe that the curvature only changes very slowly in the particular type of neural network considered. Due to this observation, a natural idea is to further reduce the computational overhead of the method by estimating the ratios  $\alpha/([G_{k+1}]_i + \mu)$  only once every few epochs, for instance, using a small subset of examples as in (6.15). The authors also mention that, as a rule of thumb, this diagonal scheme typically improves the convergence speed by a factor of only three relative to SG. Therefore, it might be more enlightening to view such an algorithm as a scheme to periodically retune a first-order SG approach rather than as a complete second-order method.

**Estimating Diagonal Curvature.** Instead of explicitly computing the diagonal terms of the curvature matrix, one can follow the template of section 6.2 and directly estimate the diagonal  $[H_k]_i$  of the inverse Hessian using displacement pairs  $\{(s_k, v_k)\}$  as defined in (6.13). For instance, [18] proposes to compute the scaling terms  $[H_k]_i$  with the running average

$$[H_{k+1}]_i = (1 - \lambda)[H_k]_i + \lambda \text{Proj} \left( \frac{[s_k]_i}{[v_k]_i} \right),$$

where  $\text{Proj}(\cdot)$  represents a projection onto a predefined positive interval. It was later found that a direct application of (6.13) after a parameter update introduces a correlated noise that ruins the curvature estimate [19, section 3]. Moreover, correcting this problem made the algorithm perform substantially worse because the chaotic behavior of the rescaling factors  $[H_k]_i$  makes the choice of the stepsize  $\alpha$  very difficult.

These problems can be addressed with a combination of two ideas [19, section 5]. The first idea consists of returning to estimating the diagonal of the Hessian instead of the diagonal of this inverse, which amounts to working with the ratio  $[v_k]_i/[s_k]_i$  instead of  $[s_k]_i/[v_k]_i$ . The second idea ensures that the effective stepsizes are monotonically decreasing by replacing the running average by the sum

$$[G_{k+1}]_i = [G_k]_i + \text{Proj} \left( \frac{[v_k]_i}{[s_k]_i} \right).$$

This effectively constructs a separate diminishing stepsize sequence  $\alpha/[G_k]_i$  for each coefficient of the parameter vector. Keeping the curvature estimates in a fixed positive interval ensures that the effective stepsizes decrease at the rate  $\mathcal{O}(1/k)$  as prescribed by Theorem 4.7, while taking the local curvature into account. This combination was shown to perform very well when the input pattern coefficients have very different variances [19], something that often happens, e.g., in text classification problems.

**Diagonal Rescaling without Curvature.** The algorithms described above often require some form of regularization to handle situations where the Hessian matrix is (nearly) singular. To illustrate why this is needed, consider, e.g., optimization of the

convex objective function

$$F(w_1, w_2) = \frac{1}{2}w_1^2 + \log(e^{w_2} + e^{-w_2}),$$

for which one finds

$$\nabla F(w_1, w_2) = \begin{bmatrix} w_1 \\ \tanh(w_2) \end{bmatrix} \quad \text{and} \quad \nabla^2 F(w_1, w_2) = \begin{bmatrix} 1 & 0 \\ 0 & 1/\cosh^2(w_2) \end{bmatrix}.$$

Performing a first-order gradient method update from a starting point of  $(3, 3)$  yields the negative gradient step  $-\nabla F \approx [-3, -1]$ , which unfortunately does not point toward the optimum, namely, the origin. Moreover, rescaling the step with the inverse Hessian actually gives a worse update direction  $-(\nabla^2 F)^{-1}\nabla F \approx [-3, -101]$  whose large second component requires a small stepsize to keep the step well contained. Batch second-order optimization algorithms can avoid having to guess a good stepsize by using, e.g., line search techniques. Stochastic second-order algorithms, on the other hand, cannot rely on such procedures as easily.

This problem is of great concern in situations where the objective function is nonconvex. For instance, optimization algorithms for DNNs must navigate around saddle points and handle near-singular curvature matrices. It is therefore tempting to consider diagonal rescaling techniques that simply ensure equal progress along each axis, rather than attempt to approximate curvature very accurately.

For instance, RMSPROP [152] estimates the average magnitude of each element of the stochastic gradient vector  $g(w_k, \xi_k)$  by maintaining the running averages

$$[R_k]_i = (1 - \lambda)[R_{k-1}]_i + \lambda[g(w_k, \xi_k)]_i^2.$$

The rescaling operation then consists in dividing each component of  $g(w_k, \xi_k)$  by the square root of the corresponding running average, ensuring that the expected second moment of each coefficient of the rescaled gradient is close to the unity:

$$[w_{k+1}]_i = [w_k]_i - \frac{\alpha}{\sqrt{[R_k]_i + \mu}} [g(w_k, \xi_k)]_i.$$

This surprisingly simple approach has been shown to be very effective for the optimization of DNNs. Various improvements have been proposed [162, 84] on an empirical basis. The theoretical explanation of this performance on nonconvex problems is still the object of active research [43].

The popular ADAGRAD algorithm [54] can be viewed as a member of this family that replaces the running average by a sum:

$$[R_k]_i = [R_{k-1}]_i + [g(w_k, \xi_k)]_i^2.$$

In this manner, the approach constructs a sequence of diminishing effective stepsizes  $\alpha/\sqrt{[R_k]_i + \mu}$  for each coefficient of the parameter vector. This algorithm was initially proposed and analyzed for the optimization of (not necessarily strongly) convex functions for which SG theory suggests diminishing stepsizes that scale with  $\mathcal{O}(1/\sqrt{k})$ . ADAGRAD is also known to perform well on deep learning networks, but one often finds that its stepsizes decrease too aggressively early in the optimization [162].

**Structural Methods.** The performance of DNN training can of course be improved by employing better optimization algorithms. However, it can also be improved by changing the structure of the network in a manner that facilitates the



optimization [80, 74]. We now describe one of these techniques, batch normalization [80], and discuss its relationship to diagonal second-order methods.

Consider a particular fully connected layer in a DNN of the form discussed in section 2.2. Using the notation of (2.4), the vector  $x_i^{(j)}$  represents the input values of layer  $j$  when the network is processing the  $i$ th training example. Omitting the layer index for simplicity, let  $\hat{x}_i = (x_i^{(j)}, 1)$  denote the input vector augmented with an additional unit coefficient and let  $\hat{w}_r = (W_{r1}, \dots, W_{rd_{j-1}}, b_r)$  be the  $r$ th row of the matrix  $W_j$  augmented with the  $r$ th coefficient of the bias vector  $b_j$ . The layer outputs are then obtained by applying the activation function to the quantities  $s_r = \hat{w}_r^T \hat{x}_i$  for  $r \in \{1, \dots, d_j\}$ . Assuming for simplicity that all other parameters of the network are kept fixed, we can write

$$F(\hat{w}_1, \dots, \hat{w}_{d_j}) = \frac{1}{n} \sum_{i=1}^n \ell(h(\hat{w}_1^T \hat{x}_i, \hat{w}_2^T \hat{x}_i, \dots, \hat{w}_{d_j}^T \hat{x}_i), y_i),$$

where  $h(s_1, \dots, s_{d_j})$  encapsulates all subsequent layers in the network. The diagonal block of the Gauss–Newton matrix (6.16) corresponding to the parameters  $\hat{w}_r$  then has the form

$$(6.25) \quad G_{[r]} = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left[ \left( \frac{dh}{ds_r} \right)^T \left( \frac{\partial^2 \ell}{\partial h^2} \right) \left( \frac{dh}{ds_r} \right) \right] \hat{x}_i \hat{x}_i^T,$$

which can be viewed as a weighted second moment matrix of the augmented input vectors  $\{\hat{x}_i\}_{i \in \mathcal{S}}$ . In particular, this matrix is perfectly conditioned if the weighted distribution of the layer inputs is white, i.e., they have zero mean and a unit covariance matrix. This could be achieved by first preprocessing the inputs by an affine transform that whitens their weighted distribution.

Two simplifications can drastically reduce the computational cost of this operation. First, we can ignore the bracketed coefficient in (6.25) and assume that we can use the same whitening transformation for all outputs  $r \in \{1, \dots, d_j\}$ . Second, we can ignore the input cross-correlations and simply ensure that each input variable has zero mean and unit variance by replacing the input vector coefficients  $\hat{x}_i[t]$  for each  $t \in \{1, \dots, d_{j-1}\}$  by the linearly transformed values  $\alpha_t \hat{x}_i[t] + \beta_t$ . Despite these simplifications, this normalization operation is very likely to improve the second-order properties of the objective function. An important detail here is the computation of the normalization constants  $\alpha_t$  and  $\beta_t$ . Estimating the mean and the standard deviation of each input with a simple running average works well if one expects these quantities to change very slowly. This is unfortunately not true in recent neural networks.<sup>8</sup>

Batch normalization [80] defines a special kind of neural network layer that performs this normalization using statistics collected with the current minibatch of examples. The back-propagation algorithm that computes the gradients must, of course, be adjusted to account for the on-the-fly computation of the normalization coefficients. Assuming that one uses sufficiently large minibatches, computing the statistics in this manner ensures that the normalization constants are very fresh. This comes at the

<sup>8</sup>This used to be true in the 1990s because neural networks were using bounded activation functions such as the sigmoid  $s(x) = 1/(1 + e^{-x})$ . However, many recent results were achieved using the ReLU activation function  $s(x) = \max\{0, x\}$ , which is unbounded and homogeneous. The statistics of the intermediate variables in such a network can change extremely quickly during the first phases of the optimization process [86].

price of making the output of the neural network on a particular training pattern dependent on the other patterns in the minibatch. Since these other examples are a priori random, this amounts to generating additional noise in the SG optimization. Although the variance of this noise is poorly controlled, inserting batch normalization layers in various points of a DNN is extremely effective and is now standard practice. Whether one can achieve the same improvement with more controlled techniques remains to be seen.

**7. Other Popular Methods.** Some optimization methods for machine learning are not well characterized as being within the two-dimensional schematic introduced in section 3.4 (see Figure 3.3), yet they represent fundamentally unique approaches that offer theoretical and/or practical advantages. The purpose of this section is to discuss a few such ideas, namely, gradient methods with momentum, accelerated gradient methods, and coordinate descent methods. For ease of exposition, we introduce these techniques under the assumption that one is minimizing a continuously differentiable (not necessarily convex) function  $F : \mathbb{R}^n \rightarrow \mathbb{R}$  and that full gradients can be computed in each iteration. Then, after each technique is introduced, we discuss how they may be applied in stochastic settings.

**7.1. Gradient Methods with Momentum.** Gradient methods with momentum are procedures in which each step is chosen as a combination of the steepest descent direction and the most recent iterate displacement. Specifically, with an initial point  $w_1$ , scalar sequences  $\{\alpha_k\}$  and  $\{\beta_k\}$  that are either predetermined or set dynamically, and  $w_0 := w_1$ , these methods are characterized by the iteration

$$(7.1) \quad w_{k+1} \leftarrow w_k - \alpha_k \nabla F(w_k) + \beta_k (w_k - w_{k-1}).$$

Here, the right-hand term is referred to as the momentum term, which, recursively, maintains the algorithm's movement along previous search directions.

The iteration (7.1) can be motivated in various ways; e.g., it is named after the fact that it represents a discretization of a certain second-order ordinary differential equation with friction. Of course, when  $\beta_k = 0$  for all  $k \in \mathbb{N}$ , it reduces to the steepest descent method. When  $\alpha_k = \alpha$  and  $\beta_k = \beta$  for some constants  $\alpha > 0$  and  $\beta > 0$  for all  $k \in \mathbb{N}$ , it is referred to as the heavy ball method [123], which is known to yield a superior rate of convergence as compared to steepest descent with a fixed stepsize for certain functions of interest. For example, when  $F$  is a strictly convex quadratic with minimum and maximum eigenvalues given by  $c > 0$  and  $L \geq c$ , respectively, steepest descent and the heavy ball method each yield a linear rate of convergence (in terms of the distance to the solution converging to zero) with contraction constants, respectively, given by

$$(7.2) \quad \frac{\kappa - 1}{\kappa + 1} \quad \text{and} \quad \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}, \quad \text{where} \quad \kappa := \frac{L}{c} \geq 1.$$

Choosing  $(\alpha, \beta)$  to achieve these rates requires knowledge of  $(c, L)$ , which might be unavailable. Still, even without this knowledge, the heavy ball method often outperforms steepest descent.

Additional connections with (7.1) can be made when  $F$  is a strictly convex quadratic. In particular, if  $(\alpha_k, \beta_k)$  is chosen optimally for all  $k \in \mathbb{N}$ , in the sense that the pair is chosen to solve

$$(7.3) \quad \min_{(\alpha, \beta)} F(w_k - \alpha \nabla F(w_k) + \beta (w_k - w_{k-1})),$$

then (7.1) is exactly the linear CG algorithm. While the heavy ball method is a stationary iteration (in the sense that the pair  $(\alpha, \beta)$  is fixed), the CG algorithm is nonstationary and its convergence behavior is relatively more complex; in particular, the step-by-step behavior of CG depends on the eigenvalue distribution of the Hessian of  $F$  [69]. That said, in contrast to the heavy ball method, CG has a finite convergence guarantee. This, along with the fact that problems with favorable eigenvalue distributions are quite prevalent, has led to the great popularity of CG in a variety of situations. More generally, nonlinear CG methods, which also follow the procedure in (7.1), can be viewed as techniques that approximate the optimal values defined by (7.3) when  $F$  is not quadratic.

An alternative view of the heavy ball method is obtained by expanding (7.1):

$$w_{k+1} \leftarrow w_k - \alpha \sum_{j=1}^k \beta^{k-j} \nabla F(w_j);$$

thus, each step can be viewed as an exponentially decaying average of past gradients. By writing the iteration this way, one can see that the steps tend to accumulate contributions in directions of persistent descent, while directions that oscillate tend to be cancelled, or at least remain small.

This latter interpretation provides some intuitive explanation as to why a stochastic heavy ball method, and SG methods with momentum in general, might be successful in various settings. In particular, their practical performance has made them popular in the community working on training DNNs [150]. Replacing the true gradient with a stochastic gradient in (7.1), one obtains an iteration that, over the long run, tends to continue moving in directions that the stochastic gradients suggest are ones of improvement, whereas movement is limited along directions along which contributions of many stochastic gradients cancel each other out. Theoretical guarantees about the inclusion of momentum in stochastic settings are elusive, and although practical gains have been reported [92, 150], more experimentation is needed.

**7.2. Accelerated Gradient Methods.** A method with an iteration similar to (7.1), but with its own unique properties, is the accelerated gradient method proposed by Nesterov [107]. Written as a two-step procedure, it involves the updates

$$(7.4) \quad \begin{aligned} \tilde{w}_k &\leftarrow w_k + \beta_k(w_k - w_{k-1}) \\ \text{and } w_{k+1} &\leftarrow \tilde{w}_k - \alpha_k \nabla F(\tilde{w}_k), \end{aligned}$$

which lead to the condensed form

$$(7.5) \quad w_{k+1} \leftarrow w_k - \alpha_k \nabla F(w_k + \beta_k(w_k - w_{k-1})) + \beta_k(w_k - w_{k-1}).$$

In this manner, it is easy to compare the approach with (7.1). In particular, one can describe their difference as being a reversal in the order of computation. In (7.1), one can imagine taking the steepest descent step and then applying the momentum term, whereas (7.5) results when one follows the momentum term first, then applies a steepest descent step (with the gradient evaluated at  $\tilde{w}_k$ , not at  $w_k$ ).

While this difference may appear to be minor, it is well known that (7.5) with appropriately chosen  $\alpha_k = \alpha > 0$  for all  $k \in \mathbb{N}$  and  $\{\beta_k\} \nearrow 1$  leads to an *optimal* iteration complexity when  $F$  is convex and continuously differentiable with a Lipschitz continuous gradient. Specifically, while in such cases a steepest descent method

converges with a distance to the optimal value decaying with a rate  $\mathcal{O}(\frac{1}{k})$ , the iteration (7.5) converges with a rate  $\mathcal{O}(\frac{1}{k^2})$ , which is provably the best rate that can be achieved by a gradient method. Unfortunately, no intuitive explanation as to how Nesterov's method achieves this optimal rate has been widely accepted. Still, one cannot deny the analysis and the practical gains that the technique has offered.

Acceleration ideas have been applied in a variety of other contexts as well, including for the minimization of nonsmooth convex functions; see [96]. For now, we merely mention that when applied in stochastic settings—with stochastic gradients employed in place of full gradients—one can only hope that acceleration might improve the constants in the convergence rate offered in Theorem 4.7; i.e., the rate itself cannot be improved [1].

**7.3. Coordinate Descent Methods.** Coordinate descent (CD) methods are among the oldest in the optimization literature. As their name suggests, they operate by taking steps along coordinate directions: one attempts to minimize the objective with respect to a single variable while all others are kept fixed, then other variables are updated similarly in an iterative manner. Such a simple idea is easy to implement, so it is not surprising that CD methods have a long history in many settings. Their limitations have been documented and well understood for many years (more on these below), but one can argue that their advantages were not fully recognized until recent work in machine learning and statistics demonstrated their ability to take advantage of certain structures commonly arising in practice.

The CD method for minimizing  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  is given by the iteration

$$(7.6) \quad w_{k+1} \leftarrow w_k - \alpha_k \nabla_{i_k} F(w_k) e_{i_k}, \quad \text{where} \quad \nabla_{i_k} F(w_k) := \frac{\partial F}{\partial w^{i_k}}(w_k),$$

$w^{i_k}$  represents the  $i_k$ th element of the parameter vector, and  $e_{i_k}$  represents the  $i_k$ th coordinate vector for some  $i_k \in \{1, \dots, d\}$ . In other words, the solution estimates  $w_{k+1}$  and  $w_k$  differ only in their  $i_k$ th element as a result of a move in the  $i_k$ th coordinate from  $w_k$ .

Specific versions of the CD method are defined by the manner in which the sequences  $\{\alpha_k\}$  and  $\{i_k\}$  are chosen. In some applications, it is possible to choose  $\alpha_k$  as the global minimizer of  $F$  from  $w_k$  along the  $i_k$ th coordinate direction. An important example of this, which has contributed to the recent revival of CD methods, occurs when the objective function has the form  $F(w) = q(w) + \|w\|_1$ , where  $q$  is a convex quadratic. Here, the exact minimization along each coordinate is not only possible, but desirable as it promotes the generation of sparse iterates; see also section 8. More often, an exact one-dimensional minimization of  $F$  is not practical, in which case one is typically satisfied with  $\alpha_k$  yielding a sufficient reduction in  $F$  from  $w_k$ . For example, so-called second-order CD methods compute  $\alpha_k$  as the minimizer of a quadratic model of  $F$  along the  $i_k$ th coordinate direction.

Concerning the choice of  $i_k$ , one could select it in each iteration in at least three different ways: by cycling through  $\{1, \dots, d\}$ ; by cycling through a random reordering of these indices (with the indices reordered after each set of  $d$  steps); or simply by choosing an index randomly with replacement in each iteration. Randomized CD algorithms (represented by the latter two strategies) have superior theoretical properties compared to the cyclic method (represented by the first strategy) as they are less likely to choose an unfortunate series of coordinates; more on this below. However, it remains an open question as to whether such randomized algorithms are more effective in typical applications.

We mention in passing that it is also natural to consider a *block-coordinate descent* method in which a handful of elements is chosen in each iteration. This is particularly effective when the objective function is (partially) block separable, which occurs in matrix factorization problems and least squares and logistic regression when each sample only depends on a few features. Clearly, in such settings, there are great advantages of a block-coordinate descent approach. However, since their basic properties are similar to the case of using a single index in each iteration, we focus on the iteration (7.6).

**Convergence Properties.** Contrary to what intuition might suggest, a CD method is not guaranteed to converge when applied to minimize any given continuously differentiable function. Powell [127] gives an example of a *nonconvex* continuously differentiable function of three variables for which a cyclic CD method, with  $\alpha_k$  chosen by exact one-dimensional minimization, cycles without converging to a solution, i.e., at any limit point the gradient of  $F$  is nonzero. Although one can argue that failures of this type are unlikely to occur in practice, particularly for a randomized CD method, they show the weakness of the myopic strategy in a CD method that considers only one variable at a time. This is in contrast with the full gradient method, which guarantees convergence to stationarity even when the objective is nonconvex.

On the other hand, if the objective  $F$  is strongly convex, the CD method will not fail and one can establish a linear rate of convergence. The analysis is very simple when using a constant stepsize and we present one such result to provide some insights into the trade-offs that arise with a CD approach. Let us assume that  $\nabla F$  is coordinate-wise Lipschitz continuous in the sense that, for all  $w \in \mathbb{R}^d$ ,  $i \in \{1, \dots, d\}$ , and  $\Delta w^i \in \mathbb{R}$ , there exists a constant  $L_i > 0$  such that

$$(7.7) \quad |\nabla_i F(w + \Delta w^i e_i) - \nabla_i F(w)| \leq L_i |\Delta w^i|.$$

We then define the maximum coordinate-wise Lipschitz constant as

$$\widehat{L} := \max_{i \in \{1, \dots, d\}} L_i.$$

Loosely speaking,  $\widehat{L}$  is a bound on the curvature of the function along all coordinates.

**THEOREM 7.1.** *Suppose that the objective function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  is continuously differentiable, strongly convex with constant  $c > 0$ , and has a gradient that is coordinate-wise Lipschitz continuous with constants  $\{L_1, \dots, L_d\}$ . In addition, suppose that  $\alpha_k = 1/\widehat{L}$  and that  $i_k$  is chosen independently and uniformly from  $\{1, \dots, d\}$  for all  $k \in \mathbb{N}$ . Then, for all  $k \in \mathbb{N}$ , the iteration (7.6) yields*

$$(7.8) \quad \mathbb{E}[F(w_{k+1})] - F_* \leq \left(1 - \frac{c}{\widehat{L}}\right)^k (F(w_1) - F_*).$$

*Proof.* As Assumption 4.1 leads to (4.3), coordinate-wise Lipschitz continuity of  $\nabla F$  yields

$$F(w_{k+1}) \leq F(w_k) + \nabla_{i_k} F(w_k)(w_{k+1}^{i_k} - w_k^{i_k}) + \frac{1}{2} \widehat{L} (w_{k+1}^{i_k} - w_k^{i_k})^2.$$

Thus, with the stepsize chosen as  $\alpha_k = 1/\widehat{L}$ , it follows that

$$F(w_{k+1}) - F(w_k) \leq -\frac{1}{L} \nabla_{i_k} F(w_k)^2 + \frac{1}{2L} \nabla_{i_k} F(w_k)^2 = -\frac{1}{2L} \nabla_{i_k} F(w_k)^2.$$

Taking expectations with respect to the distribution of  $i_k$ , one obtains

$$\begin{aligned}\mathbb{E}_{i_k}[F(w_{k+1})] - F(w_k) &\leq -\frac{1}{2L}\mathbb{E}_{i_k}[\nabla_{i_k}F(w_k)^2] \\ &= -\frac{1}{2L}\left(\frac{1}{d}\sum_{i=1}^d \nabla_i F(w_k)^2\right) \\ &= -\frac{1}{2Ld}\|\nabla F(w_k)\|_2^2.\end{aligned}$$

Subtracting  $F_*$ , taking total expectations, recalling (4.12), and applying the above inequality repeatedly over the first  $k \in \mathbb{N}$  iterations yields (7.8), as desired.  $\square$

A brief overview of the convergence properties of other CD methods under the assumptions of Theorem 7.1 and in other settings is also worthwhile. First, it is interesting to compare the result of Theorem 7.1 with a result obtained using the *deterministic* Gauss–Southwell rule, in which  $i_k$  is chosen in each iteration according to the largest (in absolute value) component of the gradient. Using this approach, one obtains a similar result in which  $c$  is replaced by  $\hat{c}$ , the strong convexity parameter as measured by the  $\ell_1$ -norm [116]. Since  $\frac{c}{n} \leq \hat{c} \leq c$ , the Gauss–Southwell rule can be up to  $n$  times faster than a randomized strategy, but in the worst case it is no better (and yet is more expensive due to the need to compute the full gradient vector in each iteration). Alternative methods have also been proposed in which, in each iteration, the index  $i_k$  is chosen randomly with probabilities proportional to  $L_i$  or according to the largest ratio  $|\nabla_i F(w_k)|/\sqrt{L_i}$  [93, 110]. These strategies also lead to linear convergence rates with constants that are better in some cases.

**Favorable Problem Structures.** Theorem 7.1 shows that a simple randomized CD method is linearly convergent with constant dependent on the parameter dimension  $d$ . At first glance, this appears to imply that such a method is less efficient than a standard full gradient method. However, in situations in which  $d$  coordinate updates can be performed at a cost similar to the evaluation of one full gradient, the method is competitive with a full gradient method both theoretically and in practice. Classes of problems in this category include those in which the objective function is

$$(7.9) \quad F(w) = \frac{1}{n} \sum_{j=1}^n \tilde{F}_j(x_j^T w) + \sum_{i=1}^d \hat{F}_i(w^i),$$

where, for all  $j \in \{1, \dots, n\}$ , the function  $\tilde{F}_j$  is continuously differentiable and dependent on the *sparse* data vector  $x_j$ , and, for all  $i \in \{1, \dots, d\}$ , the function  $\hat{F}_i$  is a (nonsmooth) regularization function. Such a form arises in least squares and logistic regression; see also section 8.

For example, consider an objective function of the form

$$f(w) = \frac{1}{2}\|Xw - y\|_2^2 + \sum_{i=1}^d \hat{F}_i(w^i) \quad \text{with} \quad X = \begin{bmatrix} x_1 & \cdots & x_n \end{bmatrix},$$

which might be the original function of interest or might represent a model of (7.9) in which the first term is approximated by a convex quadratic model. In this setting,

$$\nabla_{i_k} f(w_{k+1}) = x_{i_k}^T r_{k+1} + \hat{F}'_{i_k}(w_{k+1}^{i_k}) \quad \text{with} \quad r_{k+1} := Aw_{k+1} - b,$$

where, with  $w_{k+1} = w_k + \beta_k e_{i_k}$ , one may observe that  $r_{k+1} = r_k + \beta_k x_{i_k}$ . That is, since the residuals  $\{r_k\}$  can be updated with cost proportional to the number of

nonzeros in  $x_{i_k}$ , call it  $\text{nnz}(x_{i_k})$ , the overall cost of computing the search direction in iteration  $k + 1$  is also  $\mathcal{O}(\text{nnz}(x_{i_k}))$ . On the other hand, an evaluation of the entire gradient incurs a cost of  $\mathcal{O}(\sum_{i=1}^n \text{nnz}(x_i))$ .

Overall, there exist various types of objective functions for which minimization by a CD method (with exact gradient computations) can be effective. These include objectives that are (partially) block separable (which arise in dictionary learning and nonnegative matrix factorization problems), have structures that allow for the efficient computation of individual gradient components, or are diagonally dominant in the sense that each step along a coordinate direction yields a reduction in the objective proportional to that which would have been obtained by a step along the steepest descent direction. Additional settings in which CD methods are effective are online problems where gradient information with respect to a group of variables becomes available in time, in which case it is natural to update these variables as soon as information is received.

**Stochastic Dual Coordinate Ascent.** What about *stochastic* CD methods? As an initial thought, one might consider the replacement of  $\nabla_{i_k} F(w_k)$  in (7.6) with a stochastic approximation, but this is not typical since one can usually as easily compute a  $d$ -dimensional stochastic gradient to apply an SG method. However, an interesting setting for the application of stochastic CD methods arises when one considers approaches to minimize a convex objective function of the form (7.9) by maximizing its *dual*. In particular, defining the convex conjugate of  $\tilde{F}_j$  as  $\tilde{F}_j^*(u) := \max_w (w^T u - \tilde{F}_j(w))$ , the Fenchel–Rockafellar dual of (7.9) when  $\hat{F}_i(\cdot) = \frac{\lambda}{2}(\cdot)^2$  for all  $i \in \{1, \dots, d\}$  is given by

$$F_{\text{dual}}(v) = \frac{1}{n} \sum_{j=1}^n -\tilde{F}_j^*(-v_j) - \frac{\lambda}{2} \left\| \frac{1}{\lambda n} \sum_{j=1}^n v_j x_j \right\|_2^2.$$

The stochastic dual coordinate ascent (SDCA) method [145] applied to a function of this form has an iteration similar to (7.6), except that negative gradient steps are replaced by gradient steps due to the fact that one aims to maximize the dual. At the conclusion of a run of the algorithm, the corresponding *primal* solution can be obtained as  $w \leftarrow \frac{1}{\lambda n} \sum_{j=1}^n v_j x_j$ . The per-iteration cost of this approach is on par with that of an SG method.

**Parallel CD Methods.** We close this section by noting that CD methods are also attractive when one considers the exploitation of parallel computation. For example, consider a multicore system in which the parameter vector  $w$  is stored in shared memory. Each core can then execute a CD iteration independently and in an asynchronous manner, where if  $d$  is large compared to the number of cores, then it is unlikely that two cores are attempting to update the same variable at the same time. Since, during the time it takes a core to perform an update, the parameter vector  $w$  has likely changed (due to updates produced by other cores), each update is being made based on slightly stale information. However, convergence of the method can be proved, and it improves when one can bound the degree of staleness of each update. For further information and insight into these ideas, we refer the reader to [16, 98].

**8. Methods for Regularized Models.** Our discussion of structural risk minimization (see section 2.3) highlighted the key role played by regularization functions in the formulation of optimization problems for machine learning. The optimization

methods that we presented and analyzed in the subsequent sections (sections 3–7) are all applicable when the objective function involves a smooth regularizer, such as the squared  $\ell_2$ -norm. In this section, we expand our investigation by considering optimization methods that handle the regularization as a distinct entity, in particular, when that function is *nonsmooth*. One such regularizer that deserves special attention is the  $\ell_1$ -norm, which induces sparsity in the optimal solution vector. For machine learning, sparsity can be beneficial since it leads to simpler models, and hence can be seen as a form of *feature selection*, i.e., for biasing the optimization toward solutions where only a few elements of the parameter vector are nonzero.

Broadly, this section focuses on the nonsmooth optimization problem

$$(8.1) \quad \min_{w \in \mathbb{R}^d} \Phi(w) := F(w) + \lambda \Omega(w),$$

where  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  includes the composition of a loss and prediction function,  $\lambda > 0$  is a regularization parameter, and  $\Omega : \mathbb{R}^d \rightarrow \mathbb{R}$  is a convex, nonsmooth regularization function. Specifically, we pay special attention to methods for solving the problem

$$(8.2) \quad \min_{w \in \mathbb{R}^d} \phi(w) := F(w) + \lambda \|w\|_1.$$

As discussed in section 2, it is often necessary to solve a series of such problems over a sequence of values for the parameter  $\lambda$ . For further details in terms of problem (8.2), we refer the reader to [147, 9] for examples in a variety of applications. However, in our presentation of optimization methods, we assume that  $\lambda$  has been prescribed and is fixed. We remark in passing that (8.2) has as a special case the well-known LASSO problem [151] when  $F(w) = \|Aw - b\|_2^2$  for some  $A \in \mathbb{R}^{n \times d}$  and  $b \in \mathbb{R}^n$ .

Although nondifferentiable, the regularized  $\ell_1$  problem (8.2) has a structure that can be exploited in the design of algorithms. The algorithms that have been proposed can be grouped into classes of first- or second-order methods and distinguished as those that minimize the nonsmooth objective either directly, as in a *proximal gradient* method, or by approximately minimizing a sequence of more complicated models, such as in a *proximal Newton* method.

There exist other *sparsity-inducing* regularizers besides the  $\ell_1$ -norm, including group-sparsity-inducing regularizers that combine  $\ell_1$ - and  $\ell_2$ -norms taken with respect to groups of variables [147, 9], as well as the nuclear norm for optimization over spaces of matrices [32]. While we do not discuss other such regularizers in detail, our presentation of methods for  $\ell_1$ -norm regularized problems represents how methods for alternative regularizers can be designed and characterized.

As in the previous sections, we introduce the algorithms in this section under the assumptions that  $F$  is continuously differentiable and that full, batch gradients can be computed for it in each iteration, commenting on stochastic method variants once the motivating ideas have been described.

**8.1. First-Order Methods for Generic Convex Regularizers.** The fundamental algorithm in unconstrained smooth optimization is the gradient method. For solving problem (8.1), the *proximal gradient* method represents a similar fundamental approach. Given an iterate  $w_k$ , a generic proximal gradient iteration, with  $\alpha_k > 0$ , is given by

$$(8.3) \quad w_{k+1} \leftarrow \arg \min_{w \in \mathbb{R}^d} \left( F(w_k) + \nabla F(w_k)^T (w - w_k) + \frac{1}{2\alpha_k} \|w - w_k\|_2^2 + \lambda \Omega(w) \right).$$



The term *proximal* refers to the presence of the third term in the minimization problem on the right-hand side, which encourages the new iterate to be close to  $w_k$ . Notice that if the regularization (i.e., last) term were not present, then (8.3) would exactly recover the gradient method update  $w_{k+1} \leftarrow w_k - \alpha_k \nabla F(w_k)$ ; hence, as before, we refer to  $\alpha_k$  as the stepsize parameter. On the other hand, if the regularization term is present, then, similar to the gradient method, each new iterate is found by minimizing a model formed by a first-order Taylor series expansion of the objective function plus a simple scaled quadratic. Overall, the only thing that distinguishes a proximal gradient method from the gradient method is the regularization term, which is left untouched and included explicitly in each step computation.

To show how an analysis similar to those seen in previous sections can be used to analyze (8.3), we prove the following theorem, in which we show that if  $F$  is strongly convex and its gradient function is Lipschitz continuous, then the iteration yields a global linear rate of convergence to the optimal objective value provided that the stepsizes are sufficiently small.

**THEOREM 8.1.** *Suppose that  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  is continuously differentiable, strongly convex with constant  $c > 0$ , and has a gradient that is Lipschitz continuous with constant  $L > 0$ . In addition, suppose that  $\alpha_k = \alpha \in (0, 1/L]$  for all  $k \in \mathbb{N}$ . Then, for all  $k \in \mathbb{N}$ , the iteration (8.3) yields*

$$\Phi(w_{k+1}) - \Phi(w_*) \leq (1 - \alpha c)^k (\Phi(w_1) - \Phi(w_*)),$$

where  $w_* \in \mathbb{R}^d$  is the unique global minimizer of  $\Phi$  in (8.1).

*Proof.* Since  $\alpha_k = \alpha \in (0, 1/L]$ , it follows from (4.3) that

$$\begin{aligned} \Phi(w_{k+1}) &= F(w_{k+1}) + \lambda \Omega(w_{k+1}) \\ &\leq F(w_k) + \nabla F(w_k)^T (w_{k+1} - w_k) + \frac{1}{2} L \|w_{k+1} - w_k\|_2^2 + \lambda \Omega(w_{k+1}) \\ &\leq F(w_k) + \nabla F(w_k)^T (w_{k+1} - w_k) + \frac{1}{2\alpha} \|w_{k+1} - w_k\|_2^2 + \lambda \Omega(w_{k+1}) \\ &\leq F(w_k) + \nabla F(w_k)^T (w - w_k) + \frac{1}{2\alpha} \|w - w_k\|_2^2 + \lambda \Omega(w) \quad \text{for all } w \in \mathbb{R}^d, \end{aligned}$$

where the last inequality follows since  $w_{k+1}$  is defined by (8.3). Representing  $w = w_k + d$ , we obtain

$$\begin{aligned} \Phi(w_{k+1}) &\leq F(w_k) + \nabla F(w_k)^T d + \frac{1}{2\alpha} \|d\|_2^2 + \lambda \Omega(w_k + d) \\ &\leq F(w_k) + \nabla F(w_k)^T d + \frac{1}{2} c \|d\|_2^2 - \frac{1}{2} c \|d\|_2^2 + \frac{1}{2\alpha} \|d\|_2^2 + \lambda \Omega(w_k + d) \\ &\leq F(w_k + d) + \lambda \Omega(w_k + d) - \frac{1}{2} c \|d\|_2^2 + \frac{1}{2\alpha} \|d\|_2^2 \\ &= \Phi(w_k + d) + \frac{1}{2} \left( \frac{1}{\alpha} - c \right) \|d\|_2^2, \end{aligned}$$

which for  $d = -\alpha c(w_k - w_*)$  means that

$$\begin{aligned} \Phi(w_{k+1}) &\leq \Phi(w_k - \alpha c(w_k - w_*)) + \frac{1}{2} \left( \frac{1}{\alpha} - c \right) \|\alpha c(w_k - w_*)\|_2^2 \\ (8.4) \quad &= \Phi(w_k - \alpha c(w_k - w_*)) + \frac{1}{2} \alpha c^2 (1 - \alpha c) \|w_k - w_*\|_2^2. \end{aligned}$$

On the other hand, since the  $c$ -strongly convex function  $\Phi$  satisfies (e.g., see [108, pp. 63–64])

$$\begin{aligned} (8.5) \quad \Phi(\tau w + (1 - \tau)\bar{w}) &\leq \tau \Phi(w) + (1 - \tau)\Phi(\bar{w}) - \frac{1}{2} c \tau (1 - \tau) \|w - \bar{w}\|_2^2 \\ &\quad \text{for all } (w, \bar{w}, \tau) \in \mathbb{R}^d \times \mathbb{R}^d \times [0, 1], \end{aligned}$$

we have (considering  $\bar{w} = w_k$ ,  $w = w_*$ , and  $\tau = \alpha c \in (0, 1]$  in (8.5)) that

$$\begin{aligned} \Phi(w_k - \alpha c(w_k - w_*)) &\leq \alpha c \Phi(w_*) + (1 - \alpha c) \Phi(w_k) - \frac{1}{2} c(\alpha c)(1 - \alpha c) \|w_k - w_*\|_2^2 \\ (8.6) \qquad \qquad \qquad &= \alpha c \Phi(w_*) + (1 - \alpha c) \Phi(w_k) - \frac{1}{2} \alpha c^2 (1 - \alpha c) \|w_k - w_*\|_2^2. \end{aligned}$$

Combining (8.6) with (8.4) and subtracting  $\Phi(w_*)$ , it follows that

$$\Phi(w_{k+1}) - \Phi(w_*) \leq (1 - \alpha c)(\Phi(w_k) - \Phi(w_*)).$$

The result follows by applying this inequality repeated over the first  $k \in \mathbb{N}$  iterations.  $\square$

One finds in Theorem 8.1 an identical result as for a gradient method for minimizing a smooth strongly convex function. As in such methods, the choice of the stepsize  $\alpha$  is critical in practice; the convergence guarantees demand that it be sufficiently small, but a value that is too small might unduly slow the optimization process.

The proximal gradient iteration (8.3) is practical only when the proximal mapping

$$\text{prox}_{\lambda\Omega, \alpha_k}(\tilde{w}) := \arg \min_{w \in \mathbb{R}^n} \left( \lambda\Omega(w) + \frac{1}{2\alpha_k} \|w - \tilde{w}\|_2^2 \right)$$

can be computed efficiently. This can be seen in the fact that the iteration (8.3) can equivalently be written as  $w_{k+1} \leftarrow \text{prox}_{\lambda\Omega, \alpha_k}(w_k - \alpha_k \nabla F(w_k))$ , i.e., the iteration is equivalent to applying a proximal mapping to the result of a gradient descent step. Situations when the proximal mapping is inexpensive to compute include when  $\Omega$  is the indicator function for a simple set, such as a polyhedral set, when it is the  $\ell_1$ -norm, or, more generally, when it is separable.

A stochastic version of the proximal gradient method can be obtained, not surprisingly, by replacing  $\nabla F(w_k)$  in (8.3) by a stochastic approximation  $g(w_k, \xi_k)$ . The iteration remains cheap to perform (since  $F(w_k)$  can be ignored as it does not affect the computed step). The resulting method attains similar behavior as an SG method; analyses can be found, e.g., in [140, 8].

We now turn our attention to the most popular nonsmooth regularizer, namely, the one defined by the  $\ell_1$  norm.

**8.1.1. Iterative Soft-Thresholding Algorithm (ISTA).** In the context of solving the  $\ell_1$ -norm regularized problem (8.2), the proximal gradient method is

$$(8.7) \quad w_{k+1} \leftarrow \arg \min_{w \in \mathbb{R}^d} \left( F(w_k) + \nabla F(w_k)^T (w - w_k) + \frac{1}{2\alpha_k} \|w - w_k\|_2^2 + \lambda \|w\|_1 \right).$$

The optimization problem on the right-hand side of this expression is separable and can be solved in closed form. The solution can be written component-wise as

$$(8.8) \quad [w_{k+1}]_i \leftarrow \begin{cases} [w_k - \alpha_k \nabla F(w_k)]_i + \alpha_k \lambda & \text{if } [w_k - \alpha_k \nabla F(w_k)]_i < -\alpha_k \lambda, \\ 0 & \text{if } [w_k - \alpha_k \nabla F(w_k)]_i \in [-\alpha_k \lambda, \alpha_k \lambda], \\ [w_k - \alpha_k \nabla F(w_k)]_i - \alpha_k \lambda & \text{if } [w_k - \alpha_k \nabla F(w_k)]_i > \alpha_k \lambda. \end{cases}$$

One also finds that this iteration can be written, with  $(\cdot)_+ := \max\{\cdot, 0\}$ , as

$$(8.9) \quad w_{k+1} \leftarrow \mathcal{T}_{\alpha_k \lambda}(w_k - \alpha_k \nabla F(w_k)), \quad \text{where } [\mathcal{T}_{\alpha_k \lambda}(\tilde{w})]_i = (|\tilde{w}_i| - \alpha_k \lambda)_+ \text{sgn}(\tilde{w}_i).$$

In this form,  $\mathcal{T}_{\alpha_k \lambda}$  is referred to as the soft-thresholding operator, which leads to the name *iterative soft-thresholding algorithm* (ISTA) being used for (8.7)–(8.8) [53, 42].

It is clear from (8.8) that the ISTA iteration induces sparsity in the iterates. If the steepest descent step with respect to  $F$  yields a component with absolute value less than  $\alpha_k \lambda$ , then that component is set to zero in the subsequent iterate; otherwise, the operator still has the effect of shrinking components of the solution estimates in terms of their magnitudes. When only a stochastic estimate  $g(w_k, \xi_k)$  of the gradient is available, it can be used instead of  $\nabla F(w_k)$ .

A variant of ISTA with acceleration (recall section 7.2), known as FISTA [10], is popular in practice. We also mention that effective techniques have been developed for computing the stepsize  $\alpha_k$ , in ISTA or FISTA, based on an estimate of the Lipschitz constant of  $\nabla F$  or on curvature measured in recent iterations [10, 159, 11].

**8.1.2. Bound-Constrained Methods for  $\ell_1$ -Norm Regularized Problems.** By observing the structure created by the  $\ell_1$ -norm, one finds that an equivalent *smooth* reformulation of problem (8.2) is easily derived. In particular, by writing  $w = u - v$ , where  $u$  plays the *positive part* of  $w$  while  $v$  plays the *negative part*, problem (8.2) can equivalently be written as

$$(8.10) \quad \min_{(u,v) \in \mathbb{R}^d \times \mathbb{R}^d} \tilde{\phi}(u,v) \quad \text{s.t.} \quad (u,v) \geq 0, \quad \text{where} \quad \tilde{\phi}(u,v) = F(u-v) + \lambda \sum_{i=1}^d (u_i + v_i).$$

Now, with a bound-constrained problem in hand, one has at one's disposal a variety of optimization methods that have been developed in the optimization literature.

The fundamental iteration for solving bound-constrained optimization problems is the *gradient projection* method. In the context of (8.10), the iteration reduces to

$$(8.11) \quad \begin{bmatrix} u_{k+1} \\ v_{k+1} \end{bmatrix} \leftarrow P_+ \left( \begin{bmatrix} u_k \\ v_k \end{bmatrix} - \alpha_k \begin{bmatrix} \nabla_u \tilde{\phi}(u_k, v_k) \\ \nabla_v \tilde{\phi}(u_k, v_k) \end{bmatrix} \right) = P_+ \left( \begin{bmatrix} u_k - \alpha_k \nabla F(u_k - v_k) - \alpha_k \lambda e \\ v_k + \alpha_k \nabla F(u_k - v_k) - \alpha_k \lambda e \end{bmatrix} \right),$$

where  $P_+$  projects onto the nonnegative orthant and  $e \in \mathbb{R}^d$  is a vector of ones.

Interestingly, the gradient projection method can also be derived from the perspective of a proximal gradient method where the regularization term  $\Omega$  is chosen to be the indicator function for the feasible set (a box). In this case, the mapping  $\mathcal{T}_{\alpha_k \lambda}$  is replaced by the projection operator onto the bound constraints, causing the corresponding proximal gradient method to coincide with the gradient projection method. In the light of this observation, one should expect the iteration (8.11) to inherit the property of being globally linearly convergent when  $F$  satisfies the assumptions of Theorem 8.1. However, since the variables in (8.10) have been split into positive and negative parts, this property is maintained *only if* the iteration maintains complementarity of each iterate pair, i.e., if  $[u_k]_i [v_k]_i = 0$  for all  $k \in \mathbb{N}$  and  $i \in \{1, \dots, d\}$ . This behavior is also critical for the practical performance of the method in general, since, without it, the algorithm would not generate sparse solutions. In particular, maintaining this property allows the algorithm to be implemented in such a way that one effectively only needs  $d$  optimization variables.

A natural question that arises is whether the iteration (8.11) actually differs from an ISTA iteration, especially given that both are built upon proximal gradient ideas. In fact, the iterations can lead to the same update, but do not always do so. Consider, for example, an iterate  $w_k = u_k - v_k$  such that for  $i \in \{1, \dots, d\}$  one finds  $[w_k]_i > 0$  with  $[u_k]_i = [w_k]_i$  and  $[v_k]_i = 0$ . (A similar look, with various signs reversed, can be

taken when  $[w_k]_i < 0$ .) If  $[w_k - \alpha_k \nabla F(w_k)]_i > \alpha_k \lambda$ , then (8.8) and (8.11) yield

$$\begin{aligned} [w_{k+1}]_i &\leftarrow [w_k - \alpha_k \nabla F(w_k)]_i - \alpha_k \lambda > 0 \\ \text{and } [u_{k+1}]_i &\leftarrow [u_k - \alpha_k \nabla F(w_k)]_i - \alpha_k \lambda > 0. \end{aligned}$$

However, it is important to note the step taken in the negative part; in particular, if  $[\nabla F(w_k)]_i \leq \lambda$ , then  $[v_{k+1}]_i \leftarrow 0$ , but, if  $[\nabla F(w_k)]_i > \lambda$ , then  $[v_{k+1}]_i \leftarrow \alpha_k \nabla F(w_k) - \alpha_k \lambda$ , in which case the lack of complementarity between  $u_{k+1}$  and  $v_{k+1}$  should be rectified. A more significant difference arises when, e.g.,  $[w_k - \alpha_k \nabla F(w_k)]_i < -\alpha_k \lambda$ , in which case (8.8) and (8.11) yield

$$\begin{aligned} [w_{k+1}]_i &\leftarrow [w_k - \alpha_k \nabla F(w_k)]_i + \alpha_k \lambda < 0, \\ [u_{k+1}]_i &\leftarrow 0, \\ \text{and } [v_{k+1}]_i &\leftarrow [v_k + \alpha_k \nabla F(w_k)]_i - \alpha_k \lambda > 0. \end{aligned}$$

The pair  $([u_{k+1}]_i, [v_{k+1}]_i)$  are complementary, but  $[w_{k+1}]_i$  and  $[-v_{k+1}]_i$  differ by  $[w_k]_i > 0$ .

Several first-order [59, 58] and second-order [138] gradient projection methods have been proposed to solve (8.2). Such algorithms should be preferred over similar techniques for general bound-constrained optimization, e.g., those in [163, 95], since such general techniques may be less effective due to not exploiting the structure of the reformulation (8.10) of (8.2).

A stochastic projected gradient method, with  $\nabla F(w_k)$  replaced by  $g(w_k, \xi_k)$ , has similar convergence properties as a standard SG method; e.g., see [105]. These properties apply in the present context, but also apply when a proximal gradient method is used to solve (8.1) when  $\Omega$  represents the indicator function of a box constraint.

**8.2. Second-Order Methods.** We now turn our attention to methods that, like Newton's method for smooth optimization, are designed to solve regularized problems through successive minimization of second-order models constructed along the iterate sequence  $\{w_k\}$ . As in a proximal gradient method, the smooth function  $F$  is approximated by a Taylor series and the regularization term is kept unchanged. We focus on two classes of methods for solving (8.2): *proximal Newton* and *orthant-based* methods.

Both classes of methods fall under the category of *active-set* methods. One could also consider the application of an *interior-point* method to solve the bound-constrained problem (8.10) [114]. This, by its nature, constitutes a second-order method that would employ Hessians of  $F$  or corresponding quasi-Newton approximations. However, a disadvantage of the interior-point approach is that, by staying away from the boundary of the feasible region, it does not promote the fast generation of sparse solutions, which is in stark contrast with the methods described below.

**8.2.1. Proximal Newton Methods.** We use the term *proximal Newton* to refer to techniques that directly minimize the nonsmooth function arising as the sum of a quadratic model of  $F$  and the regularizer. In particular, for solving problem (8.2), a proximal Newton method is one that constructs, at each  $k \in \mathbb{N}$ , a model

$$(8.12) \quad q_k(w) = F(w_k) + \nabla F(w_k)^T (w - w_k) + \frac{1}{2} (w - w_k)^T H_k (w - w_k) + \lambda \|w\|_1 \approx \phi(w),$$

where  $H_k$  represents  $\nabla^2 F(w_k)$  or a quasi-Newton approximation of it. This model has a similar form to the one in (8.7), except that the simple quadratic is replaced

by the quadratic form defined by  $H_k$ . A proximal Newton method would involve (approximately) minimizing this model to compute a trial iterate  $\tilde{w}_k$ , then a stepsize  $\alpha_k > 0$  would be taken from a predetermined sequence or chosen by a line search to ensure that the new iterate  $w_{k+1} \leftarrow w_k + \alpha_k(\tilde{w}_k - w_k)$  yields  $\Phi(w_{k+1}) < \Phi(w_k)$ .

Proximal Newton methods are more challenging to design, analyze, and implement than proximal gradient methods. That being said, they can perform better in practice once a few key challenges are addressed. The three ingredients below have proved to be essential in ensuring the practical success and scalability of a proximal Newton method. For simplicity, we assume throughout that  $H_k$  has been chosen to be positive definite.

**Choice of Subproblem Solver.** The model  $q_k$  inherits the nonsmooth structure of  $\phi$ , which has the benefit of allowing a proximal Newton method to cross manifolds of nondifferentiability while simultaneously promoting sparsity of the iterates. However, the method needs to overcome the fact that the model  $q_k$  is nonsmooth, which makes the subproblem for minimizing  $q_k$  challenging. Fortunately, with the particular structure created by a quadratic plus an  $\ell_1$ -norm term, various methods are available for minimizing such a nonsmooth function. For example, CD is particularly well suited in this context [153, 79] since the global minimizer of  $q_k$  along a CD direction can be computed analytically. Such a minimizer often occurs at a point of nondifferentiability (namely, when a component is zero), thus ensuring that the method will generate sparse iterates. Updating the gradient of the model  $q_k$  after each CD step can also be performed efficiently, even if  $H_k$  is given as a limited memory quasi-Newton approximation [137]. Numerical experiments have shown that employing a CD iteration is more efficient in certain applications than employing, say, an ISTA iteration to minimize  $q_k$ , though the latter is also a viable strategy in some applications.

**Inaccurate Subproblem Solves.** A proximal Newton method needs to overcome the fact that, in large-scale settings, it is impractical to minimize  $q_k$  accurately for all  $k \in \mathbb{N}$ . Hence, it is natural to consider the computation of an inexact minimizer of  $q_k$ . The issue then becomes: what are practical, yet theoretically sufficient termination criteria when computing an approximate minimizer of the nonsmooth function  $q_k$ ? A common suggestion in the literature has been to use the norm of the minimum-norm subgradient of  $q_k$  at a given approximate minimizer. However, this measure is not continuous, making it inadequate for these purposes.<sup>9</sup> Interestingly, the norm of an ISTA step is an appropriate measure. In particular, letting  $\text{ista}_k(w)$  represent the result of an ISTA step applied to  $q_k$  from  $w$ , the value  $\|\text{ista}_k(w) - w\|_2$  satisfies the following two key properties: (i) it equals zero if and only if  $w$  is a minimizer of  $q_k$ , and (ii) it varies continuously over  $\mathbb{R}^d$ .

Complete and sufficient termination criteria are then as follows: a trial point  $\tilde{w}_k$  represents a sufficiently accurate minimizer of  $q_k$  if, for some  $\eta \in [0, 1)$ , one finds

$$\|\text{ista}_k(\tilde{w}_k) - \tilde{w}_k\|_2 \leq \eta \|\text{ista}_k(w_k) - w_k\|_2 \quad \text{and} \quad q_k(\tilde{w}_k) < q_k(w_k).$$

The latter condition, requiring a decrease in  $q_k$ , must also be imposed since the ISTA criterion alone does not exert sufficient control to ensure convergence. Employing such criteria, it has been observed to be efficient to perform the minimization of  $q_k$

<sup>9</sup>Consider the one-dimensional case of having  $q_k(w) = |w|$ . The minimum-norm subgradient of  $q_k$  has a magnitude of 1 at all points, except at the minimizer  $w^* = 0$ ; hence, this norm does not provide a measure of proximity to  $w^*$ .

inaccurately at the start, and to increase the accuracy of the model minimizers as one approaches the solution. A superlinear rate of convergence for the overall proximal Newton iteration can be obtained by replacing  $\eta$  by  $\eta_k$ , where  $\{\eta_k\} \searrow 0$ , along with imposing a stronger descent condition on the decrease in  $q_k$  [31].

**Elimination of Variables.** Due to the structure created by the  $\ell_1$ -norm regularizer, it can be effective in some applications to first identify a set of *active* variables—i.e., variables that are predicted to be equal to zero at a minimizer for  $q_k$ —then compute an approximate minimizer of  $q_k$  over the remaining *free* variables. Specifically, supposing that a set  $\mathcal{A}_k \subseteq \{1, \dots, d\}$  of active variables has been identified, one may compute an (approximate) minimizer of  $q_k$  by (approximately) solving the reduced-space problem

$$(8.13) \quad \min_{w \in \mathbb{R}^d} q_k(w) \quad \text{s.t.} \quad [w]_i = 0, \quad i \in \mathcal{A}_k.$$

Moreover, during the minimization process for this problem, one may have reason to believe that the process may be improved by adding or removing elements from the active set estimate  $\mathcal{A}_k$ . In any case, performing the elimination of variables imposed in (8.13) has the effect of reducing the size of the subproblem, and can often lead to fewer iterations being required in the overall proximal Newton method. How should the active set  $\mathcal{A}_k$  be defined? A technique that has become popular recently is to use sensitivity information, as discussed in more detail in the next subsection.

**8.2.2. Orthant-Based Methods.** Our second class of second-order methods is based on the observation that the  $\ell_1$ -norm regularized objective  $\phi$  in problem (8.2) is smooth in any orthant in  $\mathbb{R}^d$ . Based on this observation, orthant-based methods construct, at every iteration, a smooth quadratic model of the objective, then produce a search direction by minimizing this model. After performing a line search designed to reduce the objective function, a new orthant is selected and the process is repeated. This approach can be motivated by the success of second-order gradient projection methods for bound-constrained optimization, which at every iteration employ a gradient projection search to identify an active set and perform a minimization of the objective function over a space of free variables to compute a search direction.

The selection of an orthant is typically done using sensitivity information. Specifically, with the minimum norm subgradient of  $\phi$  at  $w \in \mathbb{R}^d$ , which is given component-wise for all  $i \in \{1, \dots, d\}$  by

$$(8.14) \quad \hat{g}_i(w) = \begin{cases} [\nabla F(w)]_i + \lambda & \text{if } w_i > 0 \text{ or } \{w_i = 0 \text{ and } [\nabla F(w)]_i + \lambda < 0\}, \\ [\nabla F(w)]_i - \lambda & \text{if } w_i < 0 \text{ or } \{w_i = 0 \text{ and } [\nabla F(w)]_i - \lambda > 0\}, \\ 0 & \text{otherwise,} \end{cases}$$

the active orthant for an iterate  $w_k$  is characterized by the sign vector

$$(8.15) \quad \zeta_{k,i} = \begin{cases} \text{sgn}([w_k]_i) & \text{if } [w_k]_i \neq 0, \\ \text{sgn}(-[\hat{g}(w_k)]_i) & \text{if } [w_k]_i = 0. \end{cases}$$

Along these lines, let us also define the subsets of  $\{1, \dots, d\}$  given by

$$(8.16) \quad \mathcal{A}_k = \{i : [w_k]_i = 0 \text{ and } |[\nabla F(w_k)]_i| \leq \lambda\}$$

$$(8.17) \quad \text{and } \mathcal{F}_k = \{i : [w_k]_i \neq 0\} \cup \{i : [w_k]_i = 0 \text{ and } |[\nabla F(w_k)]_i| > \lambda\},$$

where  $\mathcal{A}_k$  represents the indices of variables that are active and kept at zero, while  $\mathcal{F}_k$  represents those that are free to move.

Given these quantities, an orthant-based method proceeds as follows. First, one computes the (approximate) solution  $d_k$  of the (smooth) quadratic problem

$$\begin{aligned} \min_{d \in \mathbb{R}^n} \quad & \hat{g}(w_k)^T d + \frac{1}{2} d^T H_k d \\ \text{s.t.} \quad & d_i = 0, \quad i \in \mathcal{A}_k, \end{aligned}$$

where  $H_k$  represents  $\nabla^2 F(x^k)$  or an approximation of it. Then, the algorithm performs a line search—over a path contained in the current orthant—to compute the next iterate. For example, one option is a projected backtracking line search along  $d_k$ , which involves computing the largest  $\alpha_k$  in a decreasing geometric sequence so

$$F(P_k(w_k + \alpha_k d_k)) < F(w_k).$$

Here,  $P_k(w)$  projects  $w \in \mathbb{R}^d$  onto the orthant defined by  $\zeta^k$ , i.e.,

$$(8.18) \quad [P_k(w)]_i = \begin{cases} w_i & \text{if } \text{sgn}(w_i) = \zeta_{k,i}, \\ 0 & \text{otherwise.} \end{cases}$$

In this way, the initial and final points of an iteration lie in the same orthant. Orthant-based methods have proved to be quite useful in practice; e.g., see [7, 29].

**Commentary.** Proximal Newton and orthant-based methods represent two efficient classes of second-order active-set methods for solving the  $\ell_1$ -norm regularized problem (8.2). The proximal Newton method is reminiscent of the sequential quadratic programming method (SQP) for constrained optimization; they both solve a complex subproblem that yields a useful estimate of the optimal active set. Although solving the piecewise quadratic model (8.12) is very expensive in general, the CD method has proven to be well suited for this task and allows the proximal Newton method to be applied to very large problems [78]. Orthant-based methods have been shown to be equally useful, but in a more heuristic way, since some popular implementations lack convergence guarantees [58, 29]. Stochastic variants of both proximal Newton and orthant-based schemes can be devised in natural ways and generally inherit the properties of stochastic proximal gradient methods as long as the Hessian approximations are forced to possess eigenvalues within a positive interval.

**9. Summary and Perspectives.** Mathematical optimization is one of the foundations of machine learning, touching almost every aspect of the discipline. In particular, numerical optimization algorithms, the main subject of this paper, have played an integral role in the transformational progress that machine learning has experienced over the past two decades. In our study, we highlight the dominant role played by the stochastic gradient (SG) method of Robbins and Monro [130], whose success derives from its superior work complexity guarantees. A concise, yet broadly applicable convergence and complexity theory for the SG method is presented here, providing insight into how these guarantees have translated into practical gains.

Although the title of this paper suggests that our treatment of optimization methods for machine learning is comprehensive, much more could be said about this rapidly evolving field. Perhaps most importantly, we have neither discussed nor analyzed at length the opportunities offered by parallel and distributed computing, which may

alter our perspectives in the years to come. In fact, it has already been widely acknowledged that the SG method, despite its other benefits, may not be the best suited method for emerging computer architectures.

This leads to our discussion of a spectrum of methods that have the potential to surpass the SG method in the next generation of optimization methods for machine learning. These methods, such as those built on noise reduction and second-order techniques, offer the ability to attain improved convergence rates, overcome the adverse effects of high nonlinearity and ill-conditioning, and exploit parallelism and distributed architectures in new ways. There are important methods that are not included in our presentation—such as the alternating direction method of multipliers (ADMM) [57, 64, 67] and the expectation-maximization (EM) method and its variants [48, 160]—but our study covers many of the core algorithmic frameworks in optimization for machine learning, with emphasis on methods and theoretical guarantees that have the largest impact on practical performance.

With the great strides that have been made and the various avenues for continued contributions, numerical optimization promises to continue to have a profound impact on the rapidly growing field of machine learning.

**Appendix A. Convexity and Analyses of SG.** Our analyses of the SG method in section 4 can be characterized as relying primarily on smoothness in the sense of Assumption 4.1. This has advantages and disadvantages. On the positive side, it allows us to prove convergence results that apply equally for the minimization of convex and nonconvex functions, the latter of which has been rising in importance in machine learning; recall section 2.2. It also allows us to present results that apply equally to situations in which the stochastic vectors are unbiased estimators of the gradient of the objective, or when such estimators are scaled by a symmetric positive definite matrix; recall (4.2). A downside, however, is that it requires us to handle the minimization of nonsmooth models separately, which we do in section 8.

As an alternative, a common tactic employed by many authors is to leverage convexity instead of smoothness, allowing for the establishment of guarantees that can be applied in smooth and nonsmooth settings. For example, a typical approach for analyzing SG-based methods is to commence with the following fundamental equations related to squared distances to the optimum:

$$\begin{aligned} \|w_{k+1} - w_*\|_2^2 - \|w_k - w_*\|_2^2 &= 2(w_{k+1} - w_k)^T(w_k - w_*) + \|w_{k+1} - w_k\|_2^2 \\ (A.1) \qquad \qquad \qquad &= -2\alpha_k g(w_k, \xi_k)^T(w_k - w_*) + \alpha_k^2 \|g(w_k, \xi_k)\|_2^2. \end{aligned}$$

Assuming that  $\mathbb{E}_{\xi_k}[g(w_k, \xi_k)] = \hat{g}(w_k) \in \partial F(w_k)$ , one then obtains

$$\begin{aligned} (A.2) \qquad \mathbb{E}_{\xi_k}[\|w_{k+1} - w_*\|_2^2] - \|w_k - w_*\|_2^2 \\ = -2\alpha_k \hat{g}(w_k)^T(w_k - w_*) + \alpha_k^2 \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2], \end{aligned}$$

which has certain similarities with (4.10a). One can now introduce an assumption of convexity to bound the first term on the right-hand side in this expression; in particular, convexity offers the subgradient inequality

$$\hat{g}(w_k)^T(w_k - w_*) \geq F(w_k) - F(w_*) \geq 0,$$

while strong convexity offers the stronger condition (4.11). Combined with a suitable assumption on the second moment of the stochastic subgradients to bound the second term in the expression, the entire right-hand side can be adequately controlled through



judicious stepsize choices. The resulting analysis then has many similarities with that presented in section 4, especially if one introduces an assumption about Lipschitz continuity of the gradients of  $F$  in order to translate results on decreases in the distance to the solution in terms of decreases in  $F$  itself. The interested reader will find a clear exposition of such results in [105].

Note, however, that one can see in (A.2) that analyses based on distances to the solution do not carry over easily to nonconvex settings or when (quasi-)Newton-type steps are employed. In such situations, without explicit knowledge of  $w_*$ , one cannot easily ensure that the first term on the right-hand side can be bounded appropriately.

## Appendix B. Proofs.

*Inequality (4.3).* Under Assumption 4.1, one obtains

$$\begin{aligned} F(w) &= F(\bar{w}) + \int_0^1 \frac{\partial F(\bar{w} + t(w - \bar{w}))}{\partial t} dt \\ &= F(\bar{w}) + \int_0^1 \nabla F(\bar{w} + t(w - \bar{w}))^T (w - \bar{w}) dt \\ &= F(\bar{w}) + \nabla F(\bar{w})^T (w - \bar{w}) + \int_0^1 [\nabla F(\bar{w} + t(w - \bar{w})) - \nabla F(\bar{w})]^T (w - \bar{w}) dt \\ &\leq F(\bar{w}) + \nabla F(\bar{w})^T (w - \bar{w}) + \int_0^1 L \|t(w - \bar{w})\|_2 \|w - \bar{w}\|_2 dt, \end{aligned}$$

from which the desired result follows.  $\square$

*Inequality (4.12).* Given  $w \in \mathbb{R}^d$ , the quadratic model

$$q(\bar{w}) := F(w) + \nabla F(w)^T (\bar{w} - w) + \frac{1}{2} c \|\bar{w} - w\|_2^2$$

has the unique minimizer  $\bar{w}_* := w - \frac{1}{c} \nabla F(w)$  with  $q(\bar{w}_*) = F(w) - \frac{1}{2c} \|\nabla F(w)\|_2^2$ . Hence, the inequality (4.11) with  $\bar{w} = w_*$  and any  $w \in \mathbb{R}^d$  yields

$$F_* \geq F(w) + \nabla F(w)^T (w_* - w) + \frac{1}{2} c \|w_* - w\|_2^2 \geq F(w) - \frac{1}{2c} \|\nabla F(w)\|_2^2,$$

from which the desired result follows.  $\square$

*Corollary 4.12.* Define  $G(w) := \|\nabla F(w)\|_2^2$  and let  $L_G$  be the Lipschitz constant of  $\nabla G(w) = 2\nabla^2 F(w)\nabla F(w)$ . Then,

$$\begin{aligned} G(w_{k+1}) - G(w_k) &\leq \nabla G(w_k)^T (w_{k+1} - w_k) + \frac{1}{2} L_G \|w_{k+1} - w_k\|_2^2 \\ &\leq -\alpha_k \nabla G(w_k)^T g(w_k, \xi_k) + \frac{1}{2} \alpha_k^2 L_G \|g(w_k, \xi_k)\|_2^2. \end{aligned}$$

Taking the expectation with respect to the distribution of  $\xi_k$ , one obtains from Assumptions 4.1 and 4.3 and inequality (4.9) that

$$\begin{aligned} &\mathbb{E}_{\xi_k} [G(w_{k+1}) - G(w_k)] \\ &\leq -2\alpha_k \nabla F(w_k)^T \nabla^2 F(w_k)^T \mathbb{E}_{\xi_k} [g(w_k, \xi_k)] + \frac{1}{2} \alpha_k^2 L_G \mathbb{E}_{\xi_k} [\|g(w_k, \xi_k)\|_2^2] \\ &\leq 2\alpha_k \|\nabla F(w_k)\|_2 \|\nabla^2 F(w_k)\|_2 \mathbb{E}_{\xi_k} [g(w_k, \xi_k)] + \frac{1}{2} \alpha_k^2 L_G \mathbb{E}_{\xi_k} [\|g(w_k, \xi_k)\|_2^2] \\ &\leq 2\alpha_k L \mu_G \|\nabla F(w_k)\|_2^2 + \frac{1}{2} \alpha_k^2 L_G (M + M_V \|\nabla F(w_k)\|_2^2). \end{aligned}$$

Taking the total expectation simply yields

$$(B.1) \quad \begin{aligned} & \mathbb{E}[G(w_{k+1})] - \mathbb{E}[G(w_k)] \\ & \leq 2\alpha_k L\mu_G \mathbb{E}[\|\nabla F(w_k)\|_2^2] + \frac{1}{2}\alpha_k^2 L_G(M + M_V \mathbb{E}[\|\nabla F(w_k)\|_2^2]). \end{aligned}$$

Recall that Theorem 4.10 establishes that the first component of this bound is the term of a convergent sum. The second component of this bound is also the term of a convergent sum since  $\sum_{k=1}^{\infty} \alpha_k^2$  converges and since  $\alpha_k^2 \leq \alpha_k$  for sufficiently large  $k \in \mathbb{N}$ , meaning that again the result of Theorem 4.10 can be applied. Therefore, the right-hand side of (B.1) is the term of a convergent sum. Let us now define

$$\begin{aligned} S_K^+ &= \sum_{k=1}^K \max(0, \mathbb{E}[G(w_{k+1})] - \mathbb{E}[G(w_k)]) \\ \text{and } S_K^- &= \sum_{k=1}^K \max(0, \mathbb{E}[G(w_k)] - \mathbb{E}[G(w_{k+1})]). \end{aligned}$$

Since the bound (B.1) is positive and forms a convergent sum, the nondecreasing sequence  $S_K^+$  is upper bounded and therefore converges. Since, for any  $K \in \mathbb{N}$ , one has  $G(w_K) = G(w_0) + S_K^+ - S_K^- \geq 0$ , the nondecreasing sequence  $S_K^-$  is upper bounded and therefore also converges. Therefore  $G(w_K)$  converges and Theorem 4.9 means that this limit must be zero.  $\square$

## REFERENCES

- [1] A. AGARWAL, P. L. BARTLETT, P. RAVIKUMAR, AND M. J. WAINWRIGHT, *Information-theoretic lower bounds on the oracle complexity of stochastic convex optimization*, IEEE Trans. Inform. Theory, 58 (2012), pp. 3235–3249. (Cited on pp. 258, 259, 270, 276, 290)
- [2] N. AGARWAL, B. BULLINS, AND E. HAZAN, *Second-order stochastic optimization for machine learning in linear time*, J. Mach. Learn. Res., 18 (2017), pp. 4148–4187. (Cited on p. 273)
- [3] Z. ALLEN ZHU AND E. HAZAN, *Optimal black-box reductions between optimization objectives*, in Advances in Neural Information Processing Systems 29, D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, eds., 2016, pp. 1606–1614. (Cited on p. 247)
- [4] S.-I. AMARI, *A theory of adaptive pattern classifiers*, IEEE Trans. Electronic Computers, EC-16 (1967), pp. 299–307. (Cited on p. 238)
- [5] S.-I. AMARI, *Natural gradient works efficiently in learning*, Neural Comput., 10 (1998), pp. 251–276. (Cited on pp. 242, 280, 282, 283)
- [6] S.-I. AMARI AND H. NAGAOKA, *Methods of Information Geometry*, AMS, Providence, RI, 1997. (Cited on pp. 280, 281, 282)
- [7] G. ANDREW AND J. GAO, *Scalable training of  $L^1$ -regularized log-linear models*, in Proceedings of the 24th International Conference on Machine Learning, ACM, 2007, pp. 33–40. (Cited on p. 301)
- [8] Y. F. ATCHADE, G. FORT, AND E. MOULINES, *On Stochastic Proximal Gradient Algorithms*, preprint, <https://arxiv.org/abs/1402.2365v1>, 2014. (Cited on p. 296)
- [9] F. BACH, R. JENATTON, J. MAIRAL, AND G. OBOZINSKI, *Optimization with sparsity-inducing penalties*, Found. Trends Mach. Learning, 4 (2012), pp. 1–106. (Cited on p. 294)
- [10] A. BECK AND M. TEOULLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*, SIAM J. Imaging Sci., 2 (2009), pp. 183–202, <https://doi.org/10.1137/080716542>. (Cited on p. 297)
- [11] S. BECKER, E. J. CANDÈS, AND M. GRANT, *Templates for convex cone problems with applications to sparse signal recovery*, Math. Program. Comput., 3 (2011), pp. 165–218. (Cited on p. 297)
- [12] S. BECKER AND Y. LECUN, *Improving the convergence of back-propagation learning with second-order methods*, in Proceedings of the 1988 Connectionist Models Summer School, D. Touretzky, G. Hinton, and T. Sejnowski, eds., Morgan Kaufman, San Mateo, CA, 1989, pp. 29–37. (Cited on pp. 275, 284, 285)

- [13] D. P. BERTSEKAS, *Nonlinear Programming*, Athena Scientific, Belmont, Massachusetts, 1995. (Cited on p. 266)
- [14] D. P. BERTSEKAS, *Incremental least squares methods and the extended Kalman filter*, SIAM J. Optim., 6 (1996), pp. 807–822, <https://doi.org/10.1137/S1052623494268522>. (Cited on p. 242)
- [15] D. P. BERTSEKAS, *Convex Optimization Algorithms*, Athena Scientific, Nashua, NH, 2015. (Cited on pp. 239, 243)
- [16] D. P. BERTSEKAS AND J. N. TSITSIKLIS, *Parallel and Distributed Computation: Numerical Methods*, Prentice-Hall, Englewood Cliffs, NJ, 1989. (Cited on p. 293)
- [17] D. BLATT, A. O. HERO, AND H. GAUCHMAN, *A convergent incremental gradient method with a constant step size*, SIAM J. Optim., 18 (2007), pp. 29–51, <https://doi.org/10.1137/040615961>. (Cited on p. 268)
- [18] A. BORDES, L. BOTTOU, AND P. GALLINARI, *SGD-QN: Careful Quasi-Newton stochastic gradient descent*, J. Mach. Learn. Res., 10 (2009), pp. 1737–1754. (Cited on pp. 277, 285)
- [19] A. BORDES, L. BOTTOU, P. GALLINARI, J. CHANG, AND S. A. SMITH, *Erratum: SGDQN is less careful than expected*, J. Mach. Learn. Res., 11 (2010), pp. 2229–2240. (Cited on p. 285)
- [20] L. BOTTOU, *Stochastic gradient learning in neural networks*, in Proceedings of Neuro-Nîmes 91, Nîmes, France, 1991, EC2. (Cited on p. 238)
- [21] L. BOTTOU, *Online algorithms and stochastic approximations*, in Online Learning and Neural Networks, D. Saad, ed., Cambridge University Press, Cambridge, UK, 1998. (Cited on p. 247)
- [22] L. BOTTOU, *Large-scale machine learning with stochastic gradient descent*, in Proceedings of the 19th International Conference on Computational Statistics (COMPSTAT'2010), Y. Lechevallier and G. Saporta, eds., Paris, 2010, Springer, pp. 177–187. (Cited on pp. 255, 269)
- [23] L. BOTTOU AND O. BOUSQUET, *The tradeoffs of large scale learning*, in Advances in Neural Information Processing Systems 20, J. C. Platt, D. Koller, Y. Singer, and S. T. Roweis, eds., Curran Associates, Inc., 2008, pp. 161–168. (Cited on p. 255)
- [24] L. BOTTOU, F. FOGELMAN SOULIÉ, P. BLANCHET, AND J. S. LIENARD, *Experiments with time delay networks and dynamic time warping for speaker independent isolated digit recognition*, in Proceedings of EuroSpeech 89, Vol. 2, Paris, 1989, pp. 537–540. (Cited on p. 229)
- [25] L. BOTTOU AND Y. LE CUN, *On-line learning for very large datasets*, Appl. Stoch. Models Business Indust., 21 (2005), pp. 137–151. (Cited on pp. 269, 276)
- [26] O. BOUSQUET, *Concentration Inequalities and Empirical Processes Theory Applied to the Analysis of Learning Algorithms*, Ph.D. thesis, Ecole Polytechnique, 2002. (Cited on p. 233)
- [27] C. G. BROYDEN, *The convergence of a class of double-rank minimization algorithms*, J. Institute Math. Appl., 6 (1970), pp. 76–90. (Cited on p. 275)
- [28] R. H. BYRD, G. M. CHIN, W. NEVEITT, AND J. NOCEDAL, *On the use of stochastic Hessian information in optimization methods for machine learning*, SIAM J. Optim., 21 (2011), pp. 977–995, <https://doi.org/10.1137/10079923X>. (Cited on p. 273)
- [29] R. H. BYRD, G. M. CHIN, J. NOCEDAL, AND F. OZTOPRAK, *A family of second-order methods for convex  $\ell_1$ -regularized optimization*, Math. Program. Ser. A, 159 (2016), pp. 435–467. (Cited on p. 301)
- [30] R. H. BYRD, G. M. CHIN, J. NOCEDAL, AND Y. WU, *Sample size selection in optimization methods for machine learning*, Math. Program. Ser. B, 134 (2012), pp. 127–155. (Cited on p. 265)
- [31] R. H. BYRD, J. NOCEDAL, AND F. OZTOPRAK, *An inexact successive quadratic approximation method for convex  $L_1$  regularized optimization*, Math. Program., 157 (2017), pp. 375–396. (Cited on p. 300)
- [32] E. J. CANDÈS AND B. RECHT, *Exact matrix completion via convex optimization*, Found. Comput. Math., 9 (2009), pp. 717–772. (Cited on p. 294)
- [33] J.-F. CARDOSO, *Blind signal separation: Statistical principles*, Proc. IEEE, 9 (1998), pp. 2009–2025. (Cited on pp. 281, 283)
- [34] N. CESA-BIANCHI, A. CONCONI, AND C. GENTILE, *On the generalization ability of on-line learning algorithms*, IEEE Trans. Inform. Theory, 50 (2004), pp. 2050–2057. (Cited on p. 260)
- [35] A. CHOROMANSKA, M. HENAFF, M. MATHIEU, G. B. AROUS, AND Y. LECUN, *The loss surfaces of multilayer networks*, in Proceedings of the 18th International Conference on Artificial Intelligence and Statistics (AISTATS), 2015, pp. 192–204. (Cited on p. 275)

- [36] K. L. CHUNG, *On a stochastic approximation method*, Ann. Math. Statist., 25 (1954), pp. 463–483. (Cited on p. 251)
- [37] A. R. CONN, N. I. M. GOULD, AND P. L. TOINT, *Trust Region Methods*, SIAM, Philadelphia, 2000, <http://doi.org/10.1137/1.9780898719857>. (Cited on p. 274)
- [38] C. CORTES AND V. N. VAPNIK, *Support-Vector Networks*, Machine Learning, 20 (1995), pp. 273–297. (Cited on pp. 228, 233)
- [39] R. COURANT AND H. ROBBINS., *What Is Mathematics?*, 1st ed., Oxford University Press, 1941. (Cited on p. 238)
- [40] R. COURANT AND H. ROBBINS, *What Is Mathematics?*, 2nd ed., Oxford University Press, 1996; revised by I. Stewart. (Cited on p. 238)
- [41] T. M. COVER, *Universal portfolios*, Math. Finance, 1 (1991), pp. 1–29. (Cited on p. 260)
- [42] I. DAUBECHIES, M. DEFRISE, AND C. DE MOL, *An iterative thresholding algorithm for linear inverse problems with a sparsity constraint*, Comm. Pure Appl. Math., 58 (2004), pp. 1413–1457. (Cited on p. 296)
- [43] Y. N. DAUPHIN, H. DE VRIES, AND Y. BENGIO, *Equilibrated adaptive learning rates for non-convex optimization*, in Advances in Neural Information Processing Systems 28, 2015, pp. 1504–1512. (Cited on p. 286)
- [44] Y. N. DAUPHIN, R. PASCANU, Ç. GÜLÇEHRE, K. CHO, S. GANGULI, AND Y. BENGIO, *Identifying and attacking the saddle point problem in high-dimensional non-convex optimization*, in Advances in Neural Information Processing Systems 27, 2014, pp. 2933–2941. (Cited on p. 275)
- [45] J. DEAN, G. CORRADO, R. MONGA, K. CHEN, M. DEVIN, Q. V. LE, M. Z. MAO, M. RANZATO, A. W. SENIOR, P. A. TUCKER, K. YANG, AND A. Y. NG, *Large scale distributed deep networks*, in Advances in Neural Information Processing Systems 25, 2012, pp. 1232–1240. (Cited on p. 259)
- [46] A. DEFAZIO, F. BACH, AND S. LACOSTE-JULIEN, *SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives*, in Advances in Neural Information Processing Systems 27, 2014, pp. 1646–1654. (Cited on pp. 267, 268)
- [47] R. S. DEMBO, S. C. EISENSTAT, AND T. STEIHAUG, *Inexact Newton methods*, SIAM J. Numer. Anal., 19 (1982), pp. 400–408, <https://doi.org/10.1137/0719025>. (Cited on pp. 272, 274)
- [48] A. P. DEMPSTER, N. M. LAIRD, AND D. B. RUBIN, *Maximum likelihood from incomplete data via the EM algorithm*, J. Roy. Statist. Soc. Ser. B, 39 (1977), pp. 1–38, <https://www.jstor.org/stable/2984875>. (Cited on p. 302)
- [49] J. DENG, N. DING, Y. JIA, A. FROME, K. MURPHY, S. BENGIO, Y. LI, H. NEVEN, AND H. ADAM, *Large-scale object classification using label relation graphs*, in Proceedings of the 13th European Conference on Computer Vision (ECCV), Part I, 2014, pp. 48–64. (Cited on p. 228)
- [50] L. DENG, G. E. HINTON, AND B. KINGSBURY, *New types of deep neural network learning for speech recognition and related applications: An overview*, in IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2013), Vancouver, BC, Canada, 2013, pp. 8599–8603. (Cited on p. 229)
- [51] J. E. DENNIS, JR., AND J. J. MORÉ, *A characterization of superlinear convergence and its application to quasi-Newton methods*, Math. Comp., 28 (1974), pp. 549–560. (Cited on p. 276)
- [52] J. E. DENNIS, JR., AND R. B. SCHNABEL, *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*, Classics Appl. Math. 16, SIAM, Philadelphia, 1996, <https://doi.org/10.1137/1.9781611971200>. (Cited on p. 270)
- [53] D. DONOHO, *De-noising by soft-thresholding*, IEEE Trans. Inform. Theory, 41 (1995), pp. 613–627. (Cited on p. 296)
- [54] J. DUCHI, E. HAZAN, AND Y. SINGER, *Adaptive subgradient methods for online learning and stochastic optimization*, J. Mach. Learn. Res., 12 (2011), pp. 2121–2159. (Cited on pp. 242, 260, 286)
- [55] R. M. DUDLEY, *Uniform Central Limit Theorems*, Cambridge Stud. Adv. Math. 63, Cambridge University Press, Cambridge, UK, 1999. (Cited on p. 233)
- [56] S. T. DUMAIS, J. C. PLATT, D. HECHERMAN, AND M. SAHAMI, *Inductive learning algorithms and representations for text categorization*, in Proceedings of the 1998 ACM CIKM International Conference on Information and Knowledge Management, Bethesda, MD, 1998, pp. 148–155. (Cited on pp. 226, 227)
- [57] J. ECKSTEIN AND D. P. BERTSEKAS, *On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators*, Math. Program., 55 (1992), pp. 293–318. (Cited on p. 302)

- [58] R.-E. FAN, K.-W. CHANG, C.-J. HSIEH, X.-R. WANG, AND C.-J. LIN, *Liblinear: A library for large linear classification*, J. Mach. Learn. Res., 9 (2008), pp. 1871–1874. (Cited on pp. 298, 301)
- [59] M. A. T. FIGUEIREDO, R. D. NOWAK, AND S. J. WRIGHT, *Gradient projection for sparse reconstruction: Application to compressed sensing and other inverse problems*, IEEE J. Selected Topics Signal Process., 1 (2007), pp. 586–597. (Cited on p. 298)
- [60] R. FLETCHER, *A new approach to variable metric algorithms*, Comput. J., 13 (1970), pp. 317–322. (Cited on p. 275)
- [61] J. E. FREUND, *Mathematical Statistics*, Prentice-Hall, Englewood Cliffs, NJ, 1962. (Cited on p. 263)
- [62] M. P. FRIEDLANDER AND M. SCHMIDT, *Hybrid deterministic-stochastic methods for data fitting*, SIAM J. Sci. Comput., 34 (2012), pp. A1380–A1405, <https://doi.org/10.1137/110830629>. (Cited on p. 264)
- [63] M. C. FU, *Optimization for simulation: Theory vs. practice*, INFORMS J. Comput., 14 (2002), pp. 192–215. (Cited on p. 237)
- [64] D. GABAY AND B. MERCIER, *A dual algorithm for the solution of nonlinear variational problems via finite element approximations*, Comput. Math. Appl., 2 (1976), pp. 17–40. (Cited on p. 302)
- [65] G. GASSO, A. PAPPASIOANNOU, M. SPIVAK, AND L. BOTTOU, *Batch and online learning algorithms for nonconvex Neyman-Pearson classification*, ACM Trans. Intelligent System Technol., 3 (2011), article 28. (Cited on p. 228)
- [66] E. G. GLADYSHEV, *On stochastic approximations*, Theory Probab. Appl., 10 (1965), pp. 275–278. (Cited on p. 247)
- [67] R. GLOWINSKI AND A. MARROCCO, *Sur l'approximation, par elements finis d'ordre un, et la resolution, par penalisation-dualité, d'une classe de problèmes de Dirichlet non lineaires*, Rev. Française Automat. Informat. Recherche Opérationnelle Sér. Rouge Anal. Numér., 9 (1975), pp. 41–76. (Cited on p. 302)
- [68] D. GOLDFARB, *A family of variable metric updates derived by variational means*, Math. Comp., 24 (1970), pp. 23–26. (Cited on p. 275)
- [69] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, 4th ed., Johns Hopkins University Press, Baltimore, MD, 2012. (Cited on pp. 272, 289)
- [70] I. J. GOODFELLOW, O. VINYALS, AND A. M. SAKE, *Qualitatively Characterizing Neural Network Optimization Problems*, CoRR preprint, <https://arxiv.org/abs/1412.6544>, 2014. (Cited on p. 275)
- [71] A. GRIEWANK, *Automatic differentiation*, in Princeton Companion to Applied Mathematics, N. Higham, ed., Princeton University Press, Princeton, NJ, 2014, pp. 749–752. (Cited on p. 229)
- [72] M. GÜRBÜZBALABAN, A. OZDAGLAR, AND P. A. PARRILO, *On the convergence rate of incremental aggregated gradient algorithms*, SIAM J. Optim., 27 (2017), pp. 1035–1048, <https://doi.org/10.1137/15M1049695>. (Cited on pp. 243, 268)
- [73] F. S. HASHEMI, S. GHOSH, AND R. PASUPATHY, *On adaptive sampling rules for stochastic recursions*, in Proceedings of the 2014 Winter Simulation Conference (WSC '14), 2014, pp. 3959–3970. (Cited on p. 265)
- [74] K. HE, X. ZHANG, S. REN, AND J. SUN, *Deep residual learning for image recognition*, in Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 770–778. (Cited on p. 287)
- [75] W. HOEFFDING, *Probability inequalities for sums of bounded random variables*, J. Amer. Statist. Assoc., 58 (1963), pp. 13–30. (Cited on p. 232)
- [76] M. D. HOFFMAN, D. M. BLEI, C. WANG, AND J. W. PAISLEY, *Stochastic variational inference*, J. Mach. Learn. Res., 14 (2013), pp. 1303–1347. (Cited on pp. 281, 283)
- [77] T. HOMEM-DE MELLO AND A. SHAPIRO, *A simulation-based approach to two-stage stochastic programming with recourse*, Math. Program., 81 (1998), pp. 301–325. (Cited on p. 264)
- [78] C.-J. HSIEH, M. A. SUSTIK, I. S. DHILLON, P. K. RAVIKUMAR, AND R. POLDRACK, *BIG & QUIC: Sparse inverse covariance estimation for a million variables*, in Advances in Neural Information Processing Systems, 2013, pp. 3165–3173. (Cited on p. 301)
- [79] C. J. HSIEH, M. A. SUSTIK, P. RAVIKUMAR, AND I. S. DHILLON, *Sparse inverse covariance matrix estimation using quadratic approximation*, in Advances in Neural Information Processing Systems 24, 2011, pp. 2330–2338. (Cited on p. 299)
- [80] S. IOFFE AND C. SZEGEDY, *Batch normalization: Accelerating deep network training by reducing internal covariate shift*, in Proceedings of the 32nd International Conference on Machine Learning (ICML 2015), Lille, France, 2015, pp. 448–456. (Cited on p. 287)

- [81] T. JOACHIMS, *Text categorization with support vector machines: Learning with many relevant features*, in Proceedings of the 10th European Conference on Machine Learning (ECML'98), Chemnitz, Germany, Springer, 1998, pp. 137–142. (Cited on pp. 227, 228)
- [82] R. JOHNSON AND T. ZHANG, *Accelerating stochastic gradient descent using predictive variance reduction*, in Advances in Neural Information Processing Systems 26, 2013, pp. 315–323. (Cited on p. 266)
- [83] R. E. KALMAN, *A new approach to linear filtering and prediction problems*, Trans. ASME J. Basic Engrg. Ser. D, 82 (1960), pp. 35–45. (Cited on p. 238)
- [84] D. P. KINGMA AND J. BA, *Adam: A Method for Stochastic Optimization*, CoRR preprint, <https://arxiv.org/abs/1412.6980>, 2014. (Cited on p. 286)
- [85] A. KRIZHEVSKY, I. SUTSKEVER, AND G. E. HINTON, *ImageNet Classification with Deep Convolutional Neural Networks*, in Advances in Neural Information Processing Systems 25, 2012, pp. 84–90. (Cited on pp. 229, 230)
- [86] J. LAFOND, N. VASILACHE, AND L. BOTTOU, *Manuscript in preparation*, 2016. (Cited on p. 287)
- [87] Y. LE CUN, B. E. BOSER, J. S. DENKER, D. HENDERSON, R. E. HOWARD, W. E. HUBBARD, AND L. D. JACKEL, *Handwritten digit recognition with a back-propagation network*, in Advances in Neural Information Processing Systems 2, 1989, pp. 396–404. (Cited on p. 229)
- [88] Y. LE CUN, L. BOTTOU, Y. BENGIO, AND P. HAFFNER, *Gradient based learning applied to document recognition*, Proc. IEEE, 86 (1998), pp. 2278–2324. (Cited on p. 229)
- [89] Y. LE CUN, L. BOTTOU, G. B. ORR, AND K.-R. MÜLLER, *Efficient backprop*, in Neural Networks: Tricks of the Trade, Lecture Notes in Comput. Sci. 1524, Springer-Verlag, 1998, pp. 9–50. (Cited on pp. 231, 285)
- [90] N. LE ROUX, M. SCHMIDT, AND F. R. BACH, *A stochastic gradient method with an exponential convergence rate for finite training sets*, in Advances in Neural Information Processing Systems 25, 2012, pp. 2663–2671. (Cited on p. 268)
- [91] W. S. LEE, P. L. BARTLETT, AND R. C. WILLIAMSON, *The importance of convexity in learning with squared loss*, IEEE Trans. Inform. Theory, 44 (1998), pp. 1974–1980. (Cited on p. 257)
- [92] T. K. LEEN AND G. B. ORR, *Optimal stochastic search and adaptive momentum*, in Advances in Neural Information Processing Systems 6, 1993, pp. 477–484. (Cited on p. 289)
- [93] D. LEVENTHAL AND A. S. LEWIS, *Randomized methods for linear constraints: Convergence rates and conditioning*, Math. Oper. Res., 35 (2010), pp. 641–654. (Cited on p. 292)
- [94] D. D. LEWIS, Y. YANG, T. G. ROSE, AND F. LI, *RCV1: A new benchmark collection of text categorization research*, J. Mach. Learn. Res., 5 (2004), pp. 361–397. (Cited on p. 227)
- [95] C.-J. LIN AND J. J. MORE, *Newton's method for large bound-constrained optimization problems*, SIAM J. Optim., 9 (1999), pp. 1100–1127, <https://doi.org/10.1137/S1052623498345075>. (Cited on p. 298)
- [96] H. LIN, J. MAIRAL, AND Z. HARCHAoui, *A universal catalyst for first-order optimization*, Advances in Neural Information Processing Systems 28, 2015, pp. 3384–3392. (Cited on p. 290)
- [97] D. C. LIU AND J. NOCEDAL, *On the limited memory BFGS method for large scale optimization*, Math. Program., 45 (1989), pp. 503–528. (Cited on pp. 238, 276)
- [98] J. LIU, S. J. WRIGHT, C. RÉ, V. BITTORF, AND S. SRIDHAR, *An asynchronous parallel stochastic coordinate descent algorithm*, J. Mach. Learn. Res., 16 (2015), pp. 285–322. (Cited on p. 293)
- [99] G. MARCEAU-CARON AND Y. OLLIVIER, *Practical Riemannian Neural Networks*, CoRR preprint, <https://arxiv.org/abs/1602.08007>, 2016. (Cited on pp. 281, 283, 284)
- [100] J. MARTENS, *New Insights and Perspectives on the Natural Gradient Method*, preprint, <https://arxiv.org/abs/1412.1193>, 2014. (Cited on p. 281)
- [101] J. MARTENS AND I. SUTSKEVER, *Learning recurrent neural networks with Hessian-free optimization*, in 28th International Conference on Machine Learning (ICML), 2011. (Cited on p. 275)
- [102] P. MASSART, *Some applications of concentration inequalities to statistics*, Ann. Fac. Sci. Toulouse Math. (6), 9 (2000), pp. 245–303. (Cited on p. 233)
- [103] A. MOKHTARI AND A. RIBEIRO, *RES: Regularized stochastic BFGS algorithm*, IEEE Trans. Signal Process., 62 (2014), pp. 6089–6104. (Cited on p. 277)
- [104] N. MURATA, *A statistical study of on-line learning*, in On-line Learning in Neural Networks, D. Saad, ed., Cambridge University Press, New York, 1998, pp. 63–92. (Cited on p. 270)
- [105] A. NEMIROVSKI, A. JUDITSKY, G. LAN, AND A. SHAPIRO, *Robust stochastic approximation approach to stochastic programming*, SIAM J. Optim., 19 (2009), pp. 1574–1609, <https://doi.org/10.1137/070704277>. (Cited on pp. 250, 269, 298, 303)

- [106] A. S. NEMIROVSKI AND D. B. YUDIN, *On Cesaro's convergence of the steepest descent method for approximating saddle point of convex-concave functions*, Soviet Math. Dokl., 19 (1978). (Cited on pp. 258, 269)
- [107] Y. NESTEROV, *A method of solving a convex programming problem with convergence rate  $\mathcal{O}(1/k^2)$* , Soviet Math. Dokl., 27 (1983), pp. 372–376. (Cited on pp. 242, 289)
- [108] Y. NESTEROV, *Introductory Lectures on Convex Optimization: A Basic Course*, Kluwer Academic, 2004. (Cited on pp. 242, 295)
- [109] Y. NESTEROV, *Primal-dual subgradient methods for convex problems*, Math. Program. Ser. B, 120 (2009), pp. 221–259. (Cited on p. 269)
- [110] Y. NESTEROV, *Efficiency of coordinate descent methods on huge-scale optimization problems*, SIAM J. Optim., 22 (2012), pp. 341–362, <https://doi.org/10.1137/100802001>. (Cited on p. 292)
- [111] NIPS FOUNDATION, *Advances in Neural Information Processing Systems (29 volumes)*, 1987–2016. Volumes 0–28, <http://papers.nips.cc>. (Cited on p. 228)
- [112] F. NIU, B. RECHT, C. RE, AND S. J. WRIGHT, *Hogwild!: A lock-free approach to parallelizing stochastic gradient descent*, in Advances in Neural Information Processing Systems 24, 2011, pp. 693–701. (Cited on p. 259)
- [113] J. NOCEDAL, *Updating quasi-Newton matrices with limited storage*, Math. Comp., 35 (1980), pp. 773–782. (Cited on pp. 238, 276, 278)
- [114] J. NOCEDAL AND S. J. WRIGHT, *Numerical Optimization*, 2nd ed., Springer, New York, 2006. (Cited on pp. 237, 254, 273, 274, 278, 298)
- [115] A. B. J. NOVIKOFF, *On convergence proofs on perceptrons*, in Proceedings of the Symposium on the Mathematical Theory of Automata, Vol. 12, 1962, pp. 615–622. (Cited on p. 260)
- [116] J. NUTINI, M. SCHMIDT, I. H. LARADJI, M. FRIEDLANDER, AND H. KOEPKE, *Coordinate descent converges faster with the Gauss-Southwell rule than random selection*, in 32nd International Conference on Machine Learning (ICML), 2015, pp. 1632–1641. (Cited on p. 292)
- [117] J. M. ORTEGA AND W. C. RHEINBOLDT, *Iterative Solution of Nonlinear Equations in Several Variables*, SIAM, Philadelphia, 2000, <https://doi.org/10.1137/1.9780898719468>. (Cited on p. 240)
- [118] M. R. OSBORNE, *Fisher's method of scoring*, Internat. Statist. Rev., 60 (1992), pp. 99–117. (Cited on p. 283)
- [119] H. PARK, S. AMARI, AND K. FUKUMIZU, *Adaptive natural gradient learning algorithms for various stochastic models*, Neural Networks, 13 (2000), pp. 755–764. (Cited on pp. 281, 283, 284)
- [120] R. PASUPATHY, P. W. GLYNN, S. GHOSH, AND F. S. HASHEMI, *On sampling rates in stochastic recursions*, submitted for publication (2015). (Cited on p. 264)
- [121] B. A. PEARLMUTTER, *Fast exact multiplication by the Hessian*, Neural Comput., 6 (1994), pp. 147–160. (Cited on p. 272)
- [122] M. PILANCI AND M. J. WAINWRIGHT, *Newton sketch: A linear-time optimization algorithm with linear-quadratic convergence*, SIAM J. Optim., 27 (2017), pp. 205–245, <https://doi.org/10.1137/15M1021106>. (Cited on p. 273)
- [123] B. T. POLYAK, *Some methods of speeding up the convergence of iteration methods*, USSR Comput. Math. Math. Phys., 4 (1964), pp. 1–17. (Cited on pp. 242, 288)
- [124] B. T. POLYAK, *Comparison of the convergence rates for single-step and multi-step optimization algorithms in the presence of noise*, Engrg. Cybernet., 15 (1977), pp. 6–10. (Cited on p. 269)
- [125] B. T. POLYAK, *New method of stochastic approximation type*, Automat. Remote Control, 51 (1991), pp. 937–946. (Cited on p. 269)
- [126] B. T. POLYAK AND A. B. JUDITSKY, *Acceleration of stochastic approximation by averaging*, SIAM J. Control Optim., 30 (1992), pp. 838–855, <https://doi.org/10.1137/0330046>. (Cited on p. 269)
- [127] M. J. D. POWELL, *On search directions for minimization algorithms*, Math. Program., 4 (1973), pp. 193–201. (Cited on p. 291)
- [128] M. RAGINSKY AND A. RAKHLIN, *Information-based complexity, feedback and dynamics in convex programming*, IEEE Trans. Inform. Theory, 57 (2011), pp. 7036–7056. (Cited on pp. 258, 259)
- [129] A. RAZAVIAN, H. AZIZPOUR, J. SULLIVAN, AND S. CARLSSON, *CNN features off-the-shelf: An astounding baseline for recognition*, in 2014 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), 2014, pp. 512–519. (Cited on pp. 230, 231)
- [130] H. ROBBINS AND S. MONRO, *A stochastic approximation method*, Ann. Math. Statist., 22 (1951), pp. 400–407. (Cited on pp. 225, 236, 238, 241, 243, 249, 301)

- [131] H. ROBBINS AND D. SIEGMUND, *A convergence theorem for non negative almost supermartingales and some applications*, in *Optimizing Methods in Statistics*, J. S. Rustagi, ed., Academic Press, 1971, pp. 233–257. (Cited on p. 247)
- [132] F. ROOSTA-KHORASANI AND M. W. MAHONEY, *Sub-sampled Newton Methods II: Local Convergence Rates*, preprint, <https://arxiv.org/abs/1601.04738>, 2016. (Cited on p. 273)
- [133] S. ROSS, P. MINEIRO, AND J. LANGFORD, *Normalized online learning*, in *Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence (UAI)*, 2013. (Cited on p. 284)
- [134] D. E. RUMELHART, G. E. HINTON, AND R. J. WILLIAMS, *Learning internal representations by error propagation*, in *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, Vol. 1, MIT Press, Cambridge, MA, 1986, pp. 318–362. (Cited on pp. 228, 229, 238, 280)
- [135] D. E. RUMELHART, G. E. HINTON, AND R. J. WILLIAMS, *Learning representations by back-propagating errors*, *Nature*, 323 (1986), pp. 533–536. (Cited on pp. 228, 229, 238, 280)
- [136] D. RUPPERT, *Efficient Estimations from a Slowly Convergent Robbins-Monro Process*, Tech. report 781, Cornell University Operations Research and Industrial Engineering, 1988. (Cited on p. 269)
- [137] K. SCHEINBERG AND X. TANG, *Practical inexact proximal quasi-Newton method with global complexity analysis*, *Math. Program.*, 160 (2016), pp. 495–529. (Cited on p. 299)
- [138] M. SCHMIDT, *Graphical Model Structure Learning with  $\ell_1$ -Regularization*, Ph.D. thesis, University of British Columbia, 2010. (Cited on p. 298)
- [139] M. SCHMIDT, N. LE ROUX, AND F. BACH, *Minimizing finite sums with the stochastic average gradient*, *Math. Program.*, 162 (2017), pp. 83–112. (Cited on p. 268)
- [140] M. SCHMIDT, N. L. ROUX, AND F. R. BACH, *Convergence rates of inexact proximal-gradient methods for convex optimization*, in *Advances in Neural Information Processing Systems* 24, 2011, pp. 1458–1466. (Cited on p. 296)
- [141] N. N. SCHRAUDOLPH, *Fast curvature matrix-vector products*, in *Artificial Neural Networks, ICANN 2001*, G. Dorffner, H. Bischof, and K. Hornik, eds., *Lecture Notes in Comput. Sci.* 2130, Springer-Verlag, Berlin, Heidelberg, 2001, pp. 19–26. (Cited on pp. 242, 280)
- [142] N. N. SCHRAUDOLPH, J. YU, AND S. GÜNTHER, *A stochastic quasi-Newton method for online convex optimization*, in *International Conference on Artificial Intelligence and Statistics, Society for Artificial Intelligence and Statistics*, 2007, pp. 436–443. (Cited on p. 277)
- [143] G. SHAFER AND V. VOVK, *Probability and Finance: It's Only a Game!*, John Wiley & Sons, 2005. (Cited on p. 260)
- [144] S. SHALEV-SHWARTZ, Y. SINGER, N. SREBRO, AND A. COTTER, *Pegasos: Primal estimated sub-gradient solver for SVM*, *Math. Program.*, 127 (2011), pp. 3–30. (Cited on p. 260)
- [145] S. SHALEV-SHWARTZ AND T. ZHANG, *Stochastic dual coordinate ascent methods for regularized loss*, *J. Mach. Learn. Res.*, 14 (2013), pp. 567–599. (Cited on p. 293)
- [146] D. F. SHANNO, *Conditioning of quasi-Newton methods for function minimization*, *Math. Comp.*, 24 (1970), pp. 647–656. (Cited on p. 275)
- [147] S. SRA, S. NOWOZIN, AND S. WRIGHT, *Optimization for Machine Learning*, MIT Press, Cambridge, MA, 2011. (Cited on p. 294)
- [148] STANFORD VISION LAB, *ImageNet Large Scale Visual Recognition Challenge (ILSVRC)*, <http://image-net.org/challenges/LSVRC/>, 2015 (accessed October 25, 2015). (Cited on p. 230)
- [149] T. STEihaug, *The conjugate gradient method and trust regions in large scale optimization*, *SIAM J. Numer. Anal.*, 20 (1983), pp. 626–637, <https://doi.org/10.1137/0720042>. (Cited on p. 274)
- [150] I. SUTSKEVER, J. MARTENS, G. DAHL, AND G. HINTON, *On the importance of initialization and momentum in deep learning*, in *30th International Conference on Machine Learning (ICML)*, 2013. (Cited on p. 289)
- [151] R. TIBSHIRANI, *Regression shrinkage and selection via the lasso*, *J. Roy. Statist. Soc. Ser. B*, 58 (1996), pp. 267–288. (Cited on p. 294)
- [152] T. TIELEMAN AND G. HINTON, *Lecture 6.5. RMSPROP: Divide the gradient by a running average of its recent magnitude*. COURSERA: Neural Networks for Machine Learning, 2012. (Cited on pp. 242, 286)
- [153] P. TSENG AND S. YUN, *A coordinate gradient descent method for nonsmooth separable minimization*, *Math. Program.*, 117 (2009), pp. 387–423. (Cited on p. 299)
- [154] A. B. TSYBAKOV, *Optimal aggregation of classifiers in statistical learning*, *Ann. Statist.*, 32 (2004), pp. 135–166. (Cited on p. 257)
- [155] V. N. VAPNIK, *Estimation of Dependences Based on Empirical Data*, Springer, 1983. (Cited on pp. 231, 256)



- [156] V. N. VAPNIK, *Statistical Learning Theory*, Wiley, New York, 1998. (Cited on p. 233)
- [157] V. N. VAPNIK AND A. Y. CHERVONENKIS, *On the uniform convergence of relative frequencies of events to their probabilities*, Proc. USSR Acad. Sci., 181 (1968), pp. 781–783; English translation: Soviet Math. Dokl., 9 (1968), pp. 915–918. (Cited on pp. 231, 232)
- [158] V. N. VAPNIK AND A. Y. CHERVONENKIS, *Theory of Pattern Recognition*, Nauka, Moscow, 1974; German Translation: *Theorie der Zeichenerkennung*, Akademie-Verlag, Berlin, 1979. (Cited on p. 233)
- [159] S. J. WRIGHT, R. D. NOWAK, AND M. A. T. FIGUEIREDO, *Sparse reconstruction by separable approximation*, IEEE Trans. Signal Process., 57 (2009), pp. 2479–2493. (Cited on p. 297)
- [160] C. F. J. WU, *On the convergence properties of the EM algorithm*, Ann. Statist., 11 (1983), pp. 95–103, <https://doi.org/10.1214/aos/1176346060>. (Cited on p. 302)
- [161] W. XU, *Towards Optimal One Pass Large Scale Learning with Averaged Stochastic Gradient Descent*, preprint, <https://arxiv.org/abs/1107.2490>, 2011. (Cited on p. 269)
- [162] M. D. ZEILER, *ADADELTA: An Adaptive Learning Rate Method*, CoRR preprint, <http://arxiv.org/abs/1212.5701>, 2012. (Cited on p. 286)
- [163] C. ZHU, R. H. BYRD, P. LU, AND J. NOCEDAL, *Algorithm 78: L-BFGS-B: Fortran subroutines for large-scale bound constrained optimization*, ACM Trans. Math. Softw., 23 (1997), pp. 550–560. (Cited on p. 298)
- [164] M. ZINKEVICH, *Online Convex Programming and Generalized Infinitesimal Gradient Ascent*, in Proceedings of the 20th International Conference on Machine Learning (ICML), 2003, pp. 928–936. (Cited on p. 260)