
Stochastic Polyak Step-size for SGD: An Adaptive Learning Rate for Fast Convergence

Nicolas Loizou
Mila and DIRO
Université de Montréal

Sharan Vaswani[†]
University of Alberta

Issam Laradji
McGill, Element AI

Simon Lacoste-Julien
Mila and DIRO
Université de Montréal
Canada CIFAR AI Chair

Abstract

We propose a stochastic variant of the classical Polyak step-size (Polyak, 1987) commonly used in the subgradient method. Although computing the Polyak step-size requires knowledge of the optimal function values, this information is readily available for typical modern machine learning applications. Consequently, the proposed stochastic Polyak step-size (SPS) is an attractive choice for setting the learning rate for stochastic gradient descent (SGD). We provide theoretical convergence guarantees for SGD equipped with SPS in different settings, including strongly convex, convex and non-convex functions. Furthermore, our analysis results in novel convergence guarantees for SGD with a constant step-size. We show that SPS is particularly effective when training over-parameterized models capable of interpolating the training data. In this setting, we prove that SPS enables SGD to converge to the true solution at a fast rate without requiring the knowledge of any problem-dependent constants or additional computational overhead. We experimentally validate our theoretical results via extensive experiments on synthetic and real datasets. We demonstrate the strong performance of SGD with SPS compared to state-of-the-art optimization methods when training over-parameterized models.

1 Introduction

We solve the finite-sum optimization problem:

$$\min_{x \in \mathbb{R}^d} \left[f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \right]. \quad (1)$$

This problem is prevalent in machine learning tasks where x corresponds to the model parameters, $f_i(x)$ represents the loss on the training point i and the aim is to minimize the average loss $f(x)$ across training points. We denote $\mathcal{X}^* \subset \mathbb{R}^d$ to be the set of optimal points x^* of (1) and assume that \mathcal{X}^* is not empty. We use f^* to denote the minimum value of f , obtained at a point $x^* \in \mathcal{X}^*$. For each $i \in \{1, \dots, n\}$, we denote the infimum of function f_i by $f_i^* := \inf_x f_i(x)$. Depending on the model under study, the function f can either be strongly-convex, convex, or non-convex.

1.1 Background and Main Contributions

Stochastic gradient descent (SGD) (Robbins and Monro, 1951; Nemirovski and Yudin, 1978, 1983; Shalev-Shwartz et al., 2007; Nemirovski et al., 2009; Hardt et al., 2016), is the workhorse for training supervised machine learning problems that have the generic form (1).

Step-size selection for SGD. The main parameter for guaranteeing the convergence of SGD is the *step-size* or the *learning rate*. In recent years, several ways of selecting the step-size have been proposed. Moulines and Bach (2011); Needell et al. (2016); Needell and Ward (2017); Nguyen et al. (2018); Gower et al. (2019) propose a non-asymptotic analysis of SGD with *constant step-size* for convex and strongly convex functions. For non-convex functions, such an analysis can be found in Ghadimi and Lan (2013); Bottou et al. (2018). Using a constant step-size for SGD guarantees convergence to a neighbourhood of the solution. A common technique to guarantee convergence to the

exact optimum is to use a *decreasing step-size* (Robbins and Monro, 1951; Ghadimi and Lan, 2013; Gower et al., 2019; Nemirovski et al., 2009; Karimi et al., 2016). More recently, *adaptive* methods (Duchi et al., 2011; Liu et al., 2019; Kingma and Ba, 2015; Bengio, 2015; Vaswani et al., 2019b; Li and Orabona, 2019; Ward et al., 2019) that adjust the step-size on the fly have become wide-spread and are particularly beneficial when training deep neural networks.

Contributions: Inspired by the classical Polyak step-size (Polyak, 1987) commonly used with the deterministic subgradient method (Hazan and Kakade, 2019; Boyd et al., 2003), we propose a novel adaptive learning rate for SGD. The proposed step-size is a natural extension of the Polyak step-size to the stochastic setting. We name it **stochastic Polyak step-size (SPS)**. Although computing SPS requires knowledge of the f_i^* ; we argue that this information is readily available for modern machine learning applications (for example, $f_i^* = 0$ for most standard surrogate losses), making SPS an attractive choice for SGD.

In Section 3, we provide theoretical guarantees for the convergence of SGD with SPS in different scenarios including strongly convex, convex and non-convex smooth functions. Although SPS is provably larger than the typically used constant step-size, we guarantee its convergence to a reasonable neighborhood around the optimum. We note that in the modern machine learning tasks that we consider, it is enough to converge to a small neighbourhood and not the exact minimizer to get good generalization performance. We also establish a connection between SPS and the optimal step-size used in sketch and project methods for solving linear systems. Furthermore, in Appendix C, we provide convergence guarantees for convex non-smooth functions. We also show that by progressively increasing the batch-size for computing the stochastic gradients, SGD with SPS converges to the optimum.

Technical assumptions and challenges for proving convergence. Besides smoothness and convexity, several papers (Shamir and Zhang, 2013; Recht et al., 2011; Hazan and Kale, 2014; Rakhlin et al., 2012) assume that the variance of the stochastic gradient is bounded; that is there exists a c such that $\mathbb{E}_i \|\nabla f_i(x)\|^2 \leq c$. However, in the unconstrained setting, this assumption contradicts the assumption of strong convexity (Nguyen et al., 2018; Gower et al., 2019). In another line of work, growth conditions on the stochastic gradients have been used to guarantee convergence. In particular, the weak growth condition has been used in Bertsekas and Tsitsiklis (1996); Bottou et al. (2018); Nguyen et al. (2018). It states that there exist constants ρ, δ such that

$\mathbb{E}_i \|\nabla f_i(x)\|^2 \leq \rho \mathbb{E} \|\nabla f(x)\|^2 + \delta$. Its stronger variant (strong growth condition) when $\delta = 0$ has been used in several recent papers (Schmidt and Roux, 2013; Cevher and Vü, 2019; Vaswani et al., 2019a,b). These conditions can be relaxed to the expected smoothness assumption recently used in Gower et al. (2019).

Contributions: Our analysis of SGD with SPS does not require any of these additional assumptions for guaranteeing convergence¹. We also note that our theoretical results do not require the finite-sum assumption and can be easily adapted to the streaming setting.

In addition, unlike standard analysis for constant step-size SGD, the use of SPS requires an adaptive step-size that uses the loss and stochastic gradient estimates at an iterate, resulting in correlations. One of the main technical challenges in the proofs is to carefully analyze the SGD iterates taking these correlations into account. Furthermore, since we need to be adaptive to the Lipschitz constant, we can not use the descent lemma (implied by smoothness and SGD update). This makes the convex proof more challenging than the standard analysis.

Novel analysis for constant SGD. In the existing analyses of constant step-size SGD, the neighborhood of convergence depends on the variance of the gradients at the optimum, $z^2 := \mathbb{E}_i \|\nabla f_i(x^*)\|^2$ which is assumed to be finite.

Contributions: The proposed analysis of SGD with SPS gives a novel way to analyze constant step-size SGD. In particular, we prove convergence of *constant step-size SGD* (without SPS), to a neighbourhood that depends on $\sigma^2 := f(x^*) - \mathbb{E}[f_i^*] < \infty$ (finite optimal objective difference).

Over-parametrized models and interpolation condition. Modern machine learning models such as non-parametric regression or over-parametrized deep neural networks are highly expressive and can fit or *interpolate* the training dataset completely (Zhang et al., 2017; Ma et al., 2018). In this setting, SGD with constant step-size can be shown to converge to the exact optimum at the deterministic rate (Schmidt and Roux, 2013; Ma et al., 2018; Vaswani et al., 2019a,b; Gower et al., 2019; Berrada et al., 2020).

Contributions: As a corollary of our theoretical results, we show that SPS is particularly effective under this interpolation setting. Specifically, we prove that SPS enables SGD to converge to the true solution at a fast rate matching the deterministic case. Moreover, SPS does not require the knowledge of any problem-dependent

¹Except for our analysis for non-convex smooth functions where the weak growth condition is used.

constants or additional computational overhead.

Experimental Evaluation. In Section 4, we experimentally validate our theoretical results via experiments on synthetic datasets. We also evaluate the performance of SGD equipped with SPS relative to the state-of-the-art optimization methods when training over-parameterized models for deep matrix factorization, binary classification using kernels and multi-class classification using deep neural networks. For each of these tasks, we demonstrate the superior convergence of the proposed method. The code to reproduce our results can be found at <https://github.com/IssamLaradji/spa>.

2 SGD and the Stochastic Polyak Step-size

The optimization problem (1) can be solved using SGD:

$$x^{k+1} = x^k - \gamma_k \nabla f_i(x^k),$$

where example $i \in [n]$ is chosen uniformly at random and $\gamma_k > 0$ is the step-size in iteration k .

2.1 The Polyak step-size

Before explaining the proposed stochastic Polyak step-size, we first present the deterministic variant by Polyak (Polyak, 1987). This variant is commonly used in the analysis of deterministic subgradient methods (Boyd et al., 2003; Hazan and Kakade, 2019).

The deterministic Polyak step-size. For convex functions, the deterministic Polyak step-size at iteration k is the one that minimizes an upper-bound $Q(\gamma)$ on the distance of the iterate x_{k+1} to the optimal solution: $\|x^{k+1} - x^*\|_2^2 \leq Q(\gamma)$, where $Q(\gamma) = \|x^k - x^*\|^2 - 2\gamma[f(x^k) - f^*] + \gamma^2\|g^k\|^2$. That is,

$$\gamma_k = \operatorname{argmin}_{\gamma} [Q(\gamma)] = \frac{f(x^k) - f^*}{\|g^k\|^2}.$$

Here g^k denotes a subgradient of function f at point x^k and f^* the optimum function value. For more details and a convergence analysis of the deterministic subgradient method, please check Appendix A.2. Note that the above step-size can be used only when the optimal value f^* is known, however Boyd et al. (2003) demonstrate that $f^* = 0$ for several applications (for example, finding a point in the intersection of convex sets, positive semidefinite matrix completion and solving convex inequalities).

Stochastic Polyak Step-size. It is clear that using the deterministic Polyak step-size in the update rule of

SGD is impractical. It requires the computation of the function value f and its full gradient in each iteration.

To avoid this, we propose the stochastic Polyak step-size (SPS) for SGD:

$$\text{SPS: } \gamma_k = \frac{f_i(x^k) - f_i^*}{c \|\nabla f_i(x^k)\|^2} \quad (2)$$

Note that SPS requires the evaluation of only the stochastic gradient $\nabla f_i(x^k)$ and of the function $f_i(x^k)$ at the current iterate (quantities that can be computed in the update rule of SGD without further cost). However, it requires the knowledge of f_i^* . An important quantity in the step-size is the parameter $0 < c \in \mathbb{R}$ which can be set theoretically based on the properties of the function under study. For example, for strongly convex functions, one should select $c = 1/2$ for optimal convergence.

In addition to SPS, in some of our convergence results we require its bounded variant:

$$\text{SPS}_{\max}: \gamma_k = \min \left\{ \frac{f_i(x^k) - f_i^*}{c \|\nabla f_i(x^k)\|^2}, \gamma_b \right\} \quad (3)$$

Here $\gamma_b > 0$ is a bound that restricts SPS from being very large and is essential to ensure convergence to a small neighborhood around the solution. If $\gamma_b = \infty$ then SPS_{\max} is equivalent to SPS.

Though SPS and SPS_{\max} require knowledge of f_i^* , this information is often readily available. For machine learning problems using standard *unregularized* surrogate loss functions (e.g. squared loss for regression, logistic loss for classification), $f_i^* = 0$ (Bartlett et al., 2006). In the presence of an additional regularization term (e.g. ℓ_2 regularization), f_i^* can be obtained in closed form for these standard losses. We emphasize that since $f_i^* = \inf_x f_i(x)$, the functions f_i are not required to achieve the minimum. This is important when using loss functions such as the logistic loss for which the infimum is achieved at infinity (Soudry et al., 2018). Furthermore, we note that the deterministic Polyak step-size requires knowledge of f^* which is a much stronger assumption than the knowledge of f_i^* .

Closely related work. We now briefly compare against the recently proposed stochastic variants of the Polyak step-size (Rolinek and Martius, 2018; Oberman and Prazeres, 2019; Berrada et al., 2020). In Section 3, we present a detailed comparison of the theoretical convergence rates.

In Rolinek and Martius (2018), the L_4 algorithm has been proposed showing that a stochastic variant of the Polyak step for SGD achieves good empirical results for training neural networks. However it has *no theoretical*

convergence guarantees. The step-size is very similar to SPS (2) but each update requires an online estimation of the f_i^* which does not result in robust empirical performance and *requires up to three hyper-parameters.*

Oberman and Prazeres (2019) use a different variant of the stochastic Polyak step-size: $\gamma_k = \frac{2[f(x^k) - f^*]}{\mathbb{E}_i \|\nabla f_i(x^k)\|^2}$. This step-size requires knowledge of the quantity $\mathbb{E}_i \|\nabla f_i(x^k)\|^2$ for all iterates x^k and the evaluation of $f(x^k)$ in each step, making it impractical for finite-sum problems with large n . Moreover, their theoretical results focus only on strongly convex smooth functions.

In the ALI-G algorithm proposed by Berrada et al. (2020), the step-size is set as: $\gamma_k = \min \left\{ \frac{f_i(x^k)}{\|\nabla f_i(x^k)\|^2 + \delta}, \eta \right\}$, where $\delta > 0$ is a positive constant. Unlike our setting, their theoretical analysis relies on an ϵ -interpolation condition. Moreover, the values of the parameter δ and η that guarantee convergence heavily depend on the smoothness parameter of the objective f , limiting the method's practical applicability. In Section 3, we show that as compared to Berrada et al. (2020), the proposed method results in both better rates and a smaller neighborhood of convergence. For the case of over-parameterized models, our step-size selection guarantees convergence to the exact solution while the step proposed in Berrada et al. (2020) finds only an approximate solution that could be δ away from the optimum. In Section 4, we also experimentally show that SPS_{max} results in better convergence than ALI-G.

2.2 Optimal Objective Difference

Unlike the typical analysis of SGD that assumes a finite gradient noise $z^2 := \mathbb{E}[\|\nabla f_i(x^*)\|^2]$, in all our results, we assume a finite optimal objective difference.

Assumption 2.1 (Finite optimal objective difference).

$$\sigma^2 := \mathbb{E}_i[f_i(x^*) - f_i^*] = f(x^*) - \mathbb{E}_i[f_i^*] < \infty \quad (4)$$

This is a very weak assumption. Moreover when (1) is the training problem of an over-parametrized model such as a deep neural network or involves solving a consistent linear system or classification on linearly separable data, each individual loss function f_i attains its minimum at x^* , and thus $f_i(x^*) - f_i^* = 0$. In this *interpolation* setting, it follows that $\sigma = 0$.

3 Convergence Analysis

In this section, we present the main convergence results. For the formal definitions and properties of functions see Appendix A.1. Proofs of all key results can be found in the Appendix B.

3.1 Upper and Lower Bounds of SPS

If a function g is μ -strongly convex and L -smooth the following bounds hold: $\frac{1}{2L} \|\nabla g(x)\|^2 \leq g(x) - \inf_x g(x) \leq \frac{1}{2\mu} \|\nabla g(x)\|^2$. Using these bounds and by assuming that the functions f_i in problem (1) are μ_i -strongly convex and L_i -smooth, it is straight forward to see that SPS can be lower and upper bounded as follows:

$$\frac{1}{2cL_{\max}} \leq \frac{1}{2cL_i} \leq \gamma_k = \frac{f_i(x^k) - f_i^*}{c \|\nabla f_i(x^k)\|^2} \leq \frac{1}{2c\mu_i}, \quad (5)$$

where $L_{\max} = \max\{L_i\}_{i=1}^n$.

3.2 Sum of convex functions: strongly convex objective

In this section, we assume that all components f_i are convex functions and that the objective function f is μ -strongly convex.

Theorem 3.1. Let f_i be L_i -smooth convex functions and assume that the objective function f is μ -strongly convex function. Then, SGD with SPS_{max} with $c \geq 1/2$ converges as:

$$\mathbb{E}\|x^k - x^*\|^2 \leq (1 - \mu\alpha)^k \|x^0 - x^*\|^2 + \frac{2\gamma_b \sigma^2}{\mu\alpha}, \quad (6)$$

where $\alpha := \min\{\frac{1}{2cL_{\max}}, \gamma_b\}$ and $L_{\max} = \max\{L_i\}_{i=1}^n$ is the maximum smoothness constant. The best convergence rate and the tightest neighborhood are obtained for $c = 1/2$.

Note that in Theorem 3.1, we do not make any assumption on the value of the upper bound γ_b . However, it is clear that for convergence to a small neighborhood of the solution x^* (unique solution for strongly convex functions) γ_b should not be very large².

Another important aspect of Theorem 3.1 is that it provides convergence guarantees without requiring strong assumptions like bounded gradients or growth conditions. We do not use these conditions because SPS provides a natural bound on the norm of the gradients. In the following corollaries we make additional assumptions to better understand the convergence of SGD with SPS_{max}.

In our first corollary, we assume that our model is able to interpolate the data (each individual loss function f_i attains its minimum at x^*). This condition is satisfied for unregularized least-squares regression on a realizable dataset, or when using the squared-hinge loss on a linearly-separable dataset. The interpolation assumption enables us to guarantee the convergence of SGD

²Note that neighborhood $\frac{2\gamma_b \sigma^2}{\mu\alpha}$ has γ_b in the numerator and for the case of large γ_b , $\alpha = \frac{1}{2cL_{\max}}$.

with SPS, without an upper-bound on the step-size ($\gamma_b = \infty$).

Corollary 3.2. Assume interpolation ($\sigma = 0$) and let all assumptions of Theorem 3.1 be satisfied. SGD with SPS with $c = 1/2$ converges as:

$$\mathbb{E}\|x^k - x^*\|^2 \leq \left(1 - \frac{\mu}{L_{\max}}\right)^k \|x^0 - x^*\|^2.$$

We compare the convergence rate in Corollary 3.2 to that of stochastic line search (SLS) proposed in Vaswani et al. (2019b). In similar setting, SLS achieves the slower linear rate $\max\left\{1 - \frac{\bar{\mu}}{L_{\max}}, 1 - \gamma_b \bar{\mu}\right\}$, where $\bar{\mu} = \sum_{i=1}^n \mu_i/n$ is the average strong-convexity of the finite sum. In particular, according to Theorem 1 of Vaswani et al. (2019b), the convergence of SLS requires that at least one of the f_i 's is μ_i -strongly convex implying that the objective function f is strongly convex. This is a stronger assumption than the one we have in Theorem 3.1. We also note that $\bar{\mu} \leq \mu$.

In Berrada et al. (2020), ALI-G is analyzed under the strong assumption that all functions f_i are μ -strongly convex and L -smooth. For detailed comparison of SPS with ALI-G, see Appendix B.1.1.

An interesting outcome of Theorem 3.1 is a novel analysis for SGD with a constant step-size. In particular, note that if the bound in SPS_{\max} is selected to be $\gamma_b \leq \frac{1}{2cL_{\max}}$, then using the lower bound of (5), it can be easily shown that our method reduces to SGD with constant step-size $\gamma_k = \gamma = \gamma_b \leq \frac{1}{2cL_{\max}}$. In this case, we obtain the following convergence rate.

Corollary 3.3. Let all assumptions of Theorem 3.1 be satisfied. SGD with SPS_{\max} with $c = 1/2$ and $\gamma_b \leq \frac{1}{L_{\max}}$ becomes SGD with constant step-size $\gamma \leq \frac{1}{L_{\max}}$ and converges as:

$$\mathbb{E}\|x^k - x^*\|^2 \leq (1 - \mu\gamma)^k \|x^0 - x^*\|^2 + \frac{2\sigma^2}{\mu}.$$

If we further assume interpolation ($\sigma = 0$), the iterates of SGD with constant step-size $\gamma \leq \frac{1}{L_{\max}}$ satisfy:

$$\mathbb{E}\|x^k - x^*\|^2 \leq (1 - \mu\gamma)^k \|x^0 - x^*\|^2.$$

To the best of our knowledge, this is the first result that shows convergence of constant step-size SGD to a neighborhood that depends on the optimal objective difference σ^2 (4) and not on the variance $z^2 = \mathbb{E}[\|\nabla f_i(x^*)\|^2]$. Note that if we assume that all function f_i are μ -strongly convex and L -smooth functions then the two notions of variance satisfy the following connection: $\frac{1}{2L}z^2 \leq \sigma^2 \leq \frac{1}{2\mu}z^2$.

3.3 Sum of convex functions

Here, we derive the convergence rate when all component functions f_i are convex without any strong convexity and obtain the following theorem.

Theorem 3.4. Assume that f_i are convex, L_i -smooth functions. SGD with SPS_{\max} with $c = 1$ converges as:

$$\mathbb{E}[f(\bar{x}^K) - f(x^*)] \leq \frac{\|x^0 - x^*\|^2}{\alpha K} + \frac{2\sigma^2\gamma_b}{\alpha}.$$

Here $\alpha = \min\left\{\frac{1}{2cL_{\max}}, \gamma_b\right\}$ and $\bar{x}^K = \frac{1}{K} \sum_{k=0}^{K-1} x^k$.

Analogous to the strongly-convex case, the size of the neighbourhood is proportional to γ_b . When interpolation is satisfied and $\sigma = 0$, we observe that the unbounded variant of SPS with $\gamma_b = \infty$ converges to the optimum at a $O(1/K)$ rate. This rate is faster than the rates in Vaswani et al. (2019b); Berrada et al. (2020) and we refer the reader to the Appendix for a detailed comparison. As in the strongly-convex case, by setting $\gamma_b \leq \frac{1}{2cL_{\max}}$, we obtain the convergence rate obtained by constant step-size SGD.

3.4 Consistent Linear Systems

In Richtárik and Takác (2020), given the consistent linear system $\mathbf{A}x = b$, the authors provide a stochastic optimization reformulation of the form (1) which is equivalent to the linear system in the sense that their solution sets are identical. That is, the set of minimizers of the stochastic optimization problem \mathcal{X}^* is equal to the set of solutions of the stochastic linear system $\mathcal{L} := \{x : \mathbf{A}x = b\}$. An interesting property of this stochastic optimization problem is that $f_i(x) - f_i^* \stackrel{f_i^*=0}{=} f_i(x) = \frac{1}{2} \|\nabla f_i(x)\|^2 \quad \forall x \in \mathbb{R}^d$. Using the special structure of the problem, SPS (2) with $c = 1/2$ takes the following form: $\gamma_k \stackrel{(2)}{=} \frac{2[f_i(x^k) - f_i^*]}{\|\nabla f_i(x^k)\|^2} = 1$, which is the theoretically optimal constant step-size for SGD in this setting (Richtárik and Takác, 2020). This reduction implies that SPS results in an optimal convergence rate when solving consistent linear systems. We provide the convergence rate for SPS in this setting in Appendix B.

3.5 Sum of non-convex functions: PL Objective

We first focus on a special class of non-convex functions that satisfy the Polyak-Lojasiewicz (PL) condition (Polyak, 1987). The PL inequality is a generalization of strong-convexity and is satisfied for matrix

³For more details on the stochastic reformulation problem and its properties see Appendix B.3.

factorization (Sun and Luo, 2016) or when minimizing the logistic loss on a compact set (Karimi et al., 2016). In particular, we assume that function f satisfies the PL condition but do not assume convexity of the component functions f_i .

Definition 3.5 (Polyak-Lojasiewicz (PL) condition). We say that a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfies the PL condition if there exists $\mu > 0$ such that, $\forall x \in \mathbb{R}^n$:

$$\|\nabla f(x)\|^2 \geq 2\mu(f(x) - f^*) \quad (7)$$

Theorem 3.6. Assume that function f satisfies the PL condition (7), and let f and f_i be smooth functions. SGD with SPS_{max} with $c > \frac{L_{\max}}{4\mu}$ and $\gamma_b \geq \frac{1}{2cL_{\max}}$ converges as:

$$\mathbb{E}[f(x^k) - f(x^*)] \leq \nu^k [f(x^0) - f(x^*)] + \frac{L\sigma^2\gamma_b}{2(1-\nu)c}$$

where $\nu = \gamma_b \left(\frac{1}{\alpha} - 2\mu + \frac{L_{\max}}{2c} \right) \in (0, 1]$ and $\alpha = \min \left\{ \frac{1}{2cL_{\max}}, \gamma_b \right\}$.

Under the interpolation setting, $\sigma = 0$, and SPS_{max} converges to the optimal solution at a linear rate. If $\gamma_b \leq \min \left\{ \frac{1}{2cL_{\max}}, \frac{2c}{4\mu c - L_{\max}} \right\}$ using the lower bound in (5), the analyzed method is SGD with constant step-size and we obtain the following corollary.

Corollary 3.7. Assume that f satisfies the PL condition (7), and let f and f_i be smooth functions. SGD with constant step-size $\gamma_k = \gamma \leq \frac{\mu}{L_{\max}^2}$ converges as:

$$\mathbb{E}[f(x^k) - f(x^*)] \leq \nu^k [f(x^0) - f(x^*)] + \frac{L\sigma^2\gamma}{2(1-\nu)c}.$$

To the best of our knowledge this is the first result for the convergence of SGD for PL functions without assuming bounded gradient or bounded variance or interpolation (for more details see results in Karimi et al. (2016) and discussion in Gower et al. (2019)). In the interpolation case, we obtain linear convergence to the optimum with a constant step-size equal to that used in Vaswani et al. (2019a); Lei et al. (2019).

3.6 General Non-Convex Functions

In this section, we assume a common condition used to prove convergence of SGD in the non-convex setting (Bottou et al., 2018).

$$\mathbb{E}[\|\nabla f_i(x)\|^2] \leq \rho \|\nabla f(x)\|^2 + \delta \quad (8)$$

where $\rho, \delta > 0$ constants.

Theorem 3.8. Let f and f_i be smooth functions and assume that there exist $\rho, \delta > 0$ such that the condition (8) is satisfied. SGD with SPS_{max} with $c > \frac{\rho L}{4L_{\max}}$ and $\gamma_b < \max \left\{ \frac{2}{L\rho}, \bar{\gamma}_b \right\}$ converges as:

$$\min_{k \in [K]} \mathbb{E} \|\nabla f(x^k)\|^2 \leq \frac{2}{\zeta K} (f(x^0) - f(x^*)) + \frac{(\gamma_b - \alpha + L\gamma_b^2) \delta}{\zeta},$$

where $\alpha = \min \left\{ \frac{1}{2cL_{\max}}, \gamma_b \right\}$, $\zeta = (\gamma_b + \alpha) - \rho(\gamma_b - \alpha + L\gamma_b^2)$ and

$$\bar{\gamma}_b := \frac{-(\rho - 1) + \sqrt{(\rho - 1)^2 + \frac{4L\rho(\rho + 1)}{2cL_{\max}}}}{2L\rho}.$$

From the above theorem, we observe that SGD with SPS results in $O(1/K)$ convergence to a neighborhood governed by δ . For the case that $\delta = 0$, condition (8) reduces to the strong growth condition (SGC) used in several recent papers (Schmidt and Roux, 2013; Vaswani et al., 2019b,a). It can be easily shown that functions that satisfy the SGC condition necessarily satisfy the interpolation property (Vaswani et al., 2019a). In the special case of interpolation, SGD with SPS is able to find a first-order stationary point as efficiently as deterministic gradient descent. Moreover, for $c \in \left(\frac{\rho L}{4L_{\max}}, \frac{\rho L}{2L_{\max}} \right]$, the lower bound $\frac{1}{2cL_{\max}}$ of SPS lies in the range $\left[\frac{1}{\rho L}, \frac{2}{\rho L} \right)$ and thus the step-size is larger than $\frac{1}{\rho L}$, the best constant step-size analyzed in this setting (Vaswani et al., 2019a).

3.7 Additional Convergence Results

In Appendix C, we prove a $O(1/\sqrt{K})$ convergence rate for non-smooth convex functions. Furthermore, similar to Schmidt et al. (2011), we propose a method to increase the mini-batch size for evaluating the stochastic gradient and guarantee convergence to the optimal solution without interpolation.

4 Experimental Evaluation

We validate our theoretical results using synthetic experiments in Section 4.1. In Section 4.2, we evaluate the performance of SGD with SPS when training over-parametrized models. In particular, we compare against state-of-the-art optimization methods for deep matrix factorization, binary classification using kernel methods and multi-class classification using standard deep neural network models.

4.1 Synthetic experiments

We use a synthetic dataset to validate our theoretical results. Following the procedure outlined in Nutini et al. (2017), we generate a sparse dataset for binary classification with the number of examples $n = 1k$ and dimension $d = 100$. We use the logistic loss with and without ℓ_2 regularization. The data is generated to ensure that the function f is strongly convex in both cases. We evaluate the performance of SPS_{\max} and set $c = 1/2$ as suggested by Theorem 3.1. We experiment with three values of $\gamma_b = \{1, 5, 100\}$. In the regularized case, f_i^* can be pre-computed in closed form for each i using the Lambert W function (Corless et al., 1996) (see Appendix D); while f_i^* is simply zero in the unregularized case. A similar observation has been used to construct a “truncated” model for improving the robustness of gradient descent in Asi and Duchi (2019). In both cases, we benchmark the performance of SPS against constant step-size SGD with $\gamma = \{0.1, 0.01\}$. From Figure 1, we observe that constant step-size SGD is not robust to the step-size; it has good convergence with step-size 0.1, slow convergence when using a step-size of 0.01 and we observe divergence for larger step-sizes. In contrast, all the variants of SPS converge to a neighbourhood of the optimum and the size of the neighbourhood increases as γ_b increases as predicted by the theory.

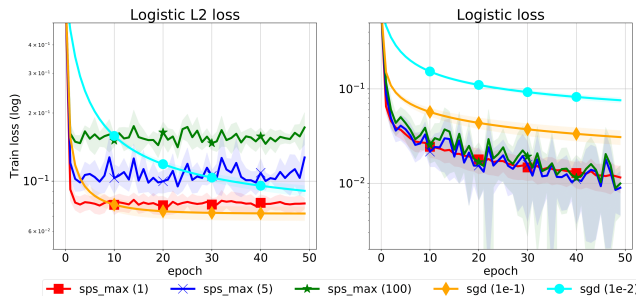


Figure 1: Synthetic experiment to benchmark SPS against constant step-size SGD for binary classification using the (left) regularized and (right) unregularized logistic loss.

4.2 Experiments for over-parametrized models

In this section, we consider training over-parameterized models that (approximately) satisfy the interpolation condition. Following the logic of the previous section, we evaluate the performance of both the SPS and SPS_{\max} variants with $f_i^* = 0$. Throughout our experiments, we found that SPS without an upper-bound on the step-size is not robust to the misspecification of interpolation and results in large fluctuations when inter-

polation is not exactly satisfied. For SPS_{\max} , the value of γ_b that results in good convergence depends on the problem and requires careful parameter tuning. This is also evidenced by the highly variable performance of ALI-G (Berrada et al., 2020) that uses a constant upper-bound on the step-size. To alleviate this problem, we use a *smoothing* procedure that prevents large fluctuations in the step-size across iterations. This can be viewed as using an adaptive iteration-dependent upper-bound γ_b^k where $\gamma_b^k = \tau^{b/n} \gamma^{k-1}$. Here, τ is a tunable hyper-parameter set to 2 in all our experiments, b is the batch-size and n is the number of examples. We note that using an adaptive γ_b can be easily handled by our theoretical results. A similar smoothing procedure has been used to control the magnitude of the step-sizes when using the Barzilai-Borwein step-size selection procedure for SGD (Tan et al., 2016) and is related to the “reset” option for using larger step-sizes in Vaswani et al. (2019b). We set $c = 1/2$ for binary classification using kernels (convex case) and deep matrix factorization (non-convex PL case). For multi-class classification using deep networks, we empirically find that any value of $c \geq 0.2$ results in convergence. In this case, we observed that across models and datasets, the fastest convergence is obtained with $c = 0.2$ and use this value.

We compare our methods against Adam (Kingma and Ba, 2015), which is the most common adaptive method, and other recent methods that report better performance than Adam: (i) stochastic line-search (SLS) (Vaswani et al., 2019b) (ii) ALI-G (Berrada et al., 2020)⁴ (iii) rectified Adam (RADAM) (Liu et al., 2019) (iv) Look-ahead optimizer (Zhang et al., 2019). We use the default learning rates and momentum (non-zero) parameters and the publicly available code for the competing methods. All our results are averaged across 5 independent runs.

Deep matrix factorization. In the first experiment, we use deep matrix factorization to examine the effect of over-parametrization for the different optimizers. In particular, we solve the non-convex regression problem: $\min_{W_1, W_2} \mathbb{E}_{x \sim N(0, I)} \|W_2 W_1 x - Ax\|^2$ and use the experimental setup in Rolinek and Martius (2018); Vaswani et al. (2019b); Rahimi and Recht (2017). We choose $A \in \mathbb{R}^{10 \times 6}$ with condition number $\kappa(A) = 10^{10}$ and generate a fixed dataset of 1000 samples. We control the degree of over-parametrization via the rank k of the matrix factors $W_1 \in \mathbb{R}^{k \times 6}$ and

⁴With ALI-G we refer to the method analyzed in Berrada et al. (2020). This is SGD with step-size the one described in Section 2. We highlight that the experiments in Berrada et al. (2020) used momentum on top of the analyzed method but without any convergence guarantees. To ensure a fair comparison with SPS, we do not use such momentum.

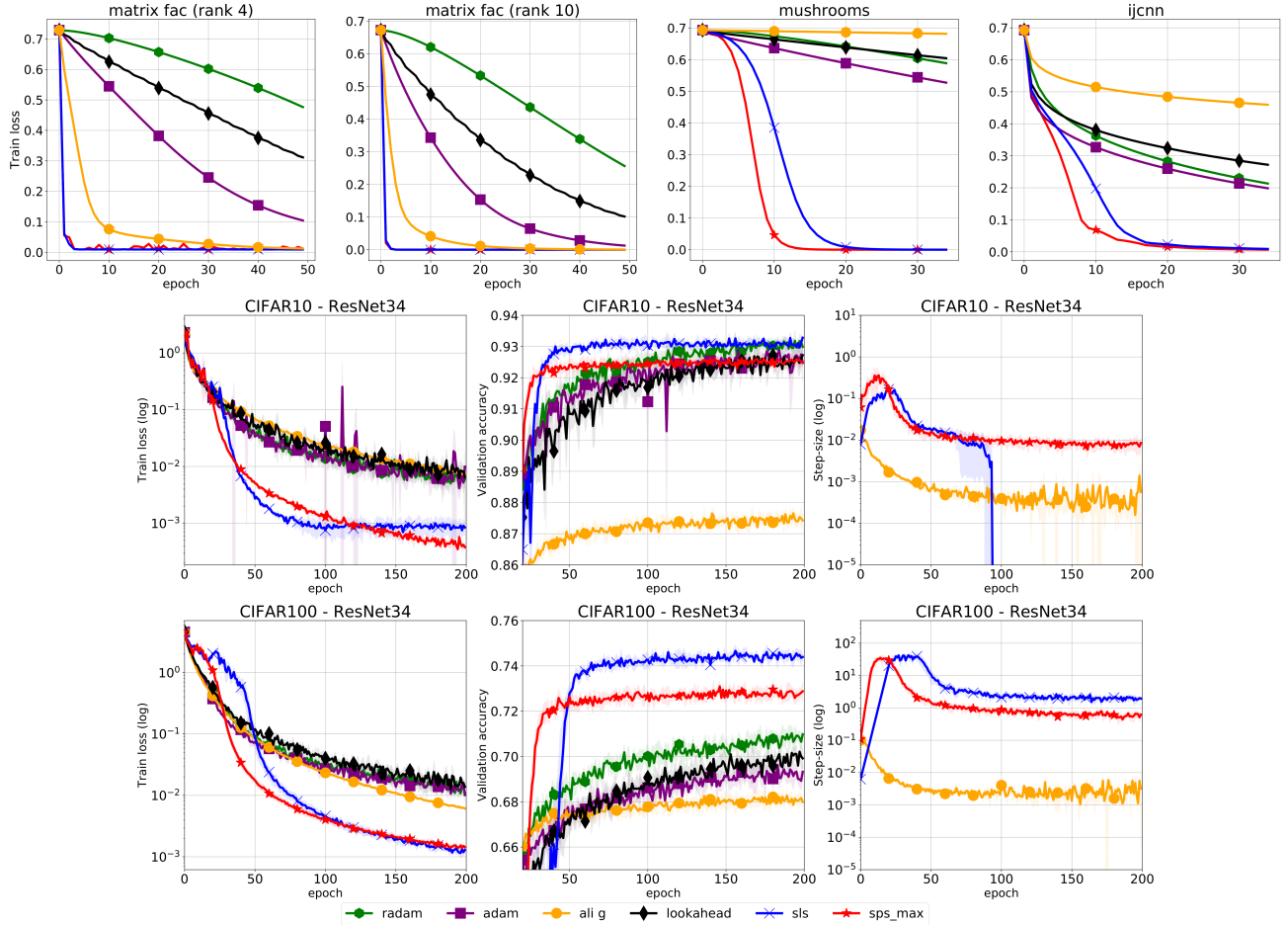


Figure 2: Comparing the performance of optimizers on deep matrix factorization (top left) and binary classification using kernels (top right) and multi-class classification on CIFAR-10 and CIFAR-100 with ResNet34.

$W_2 \in \mathbb{R}^{10 \times k}$. In Figure 2, we show the training loss as we vary the rank $k \in \{4, 10\}$ (additional experiments are in Appendix E). For $k = 4$, the interpolation condition is *not* satisfied, whereas it is exactly satisfied for $k = 10$. We observe that (i) SPS is robust to the degree of over-parametrization and (ii) has performance equal to that of SLS. However, note that SPS does not require the expensive back-tracking procedure of SLS and is arguably simpler to implement.

Binary classification using kernels. Next, we compare the optimizers’ performance in the convex, interpolation regime. We consider binary classification using RBF kernels, using the logistic loss without regularization. The bandwidths for the RBF kernels are set according to the validation procedure described in Vaswani et al. (2019b). We experiment with four standard datasets: mushrooms, rcv1, ijcnn, and w8a from LIBSVM (Chang and Lin, 2011). Figure 2 shows the training loss on the mushrooms and ijcnn for the different optimizers. Again, we observe the strong

performance of SPS compared to the other optimizers.

Multi-class classification using deep networks.

We benchmark the convergence rate and generalization performance of SPS methods on standard deep learning experiments. We consider non-convex minimization for multi-class classification using deep network models on the CIFAR10 and CIFAR100 datasets. Our experimental choices follow the setup in Luo et al. (2019). For CIFAR10 and CIFAR100, we experiment with the standard image-classification architectures: ResNet-34 (He et al., 2016) and DenseNet-121 (Huang et al., 2017). For space concerns, we report only the ResNet experiments in the main paper and relegate the DenseNet and MNIST experiments to Appendix E. From Figure 2, we observe that SPS results in the best training loss across models and datasets. For CIFAR-10, SPS results in competitive generalization performance compared to the other optimizers, whereas for CIFAR-100, its generalization performance is better than all optimizers except SLS. Note that ALI-G, the closest related

optimizer results in worse generalization performance in all cases. We note that SPS is able to match the performance of SLS, but does not require an expensive back-tracking line-search or additional tricks.

For this set of experiments, we also plot how the step-size varies across iterations for SLS, SPS and ALI-G. Interestingly, for both CIFAR-10 and CIFAR-100, we find that step-size for both SPS and SLS follows a cyclic behaviour - a warm-up period where the step-size first increases and then decreases to a constant value. Such a step-size schedule has been empirically found to result in good training and generalization performance (Loshchilov and Hutter, 2017) and our results show that SPS is able to simulate this behaviour.

5 Conclusion

We proposed and theoretically analyzed a stochastic variant of the classical the Polyak step-size. We quantified the convergence rate of SPS in numerous settings and used our analysis techniques to prove new results for constant step-size SGD. Furthermore, via experiments on a variety of tasks we showed the strong performance of SGD with SPS as compared to state-of-the-art optimization methods. There are many possible interesting extensions of our work: using SPS with accelerated methods, studying the effect of mini-batching and non-uniform sampling techniques and extensions to the distributed and decentralized settings.

Acknowledgements

Nicolas Loizou and Sharan Vaswani acknowledge support by the IVADO Postdoctoral Funding Program. Issam Laradji is funded by the UBC Four-Year Doctoral Fellowships (4YF). This research was partially supported by the Canada CIFAR AI Chair Program and by a Google Focused Research award. Simon Lacoste-Julien is a CIFAR Associate Fellow in the Learning in Machines & Brains program.

The authors would like to thank Frederik Kunstner for help with the convex proofs, and Aaron Defazio for fruitful discussions and feedback on the manuscript.

References

- Asi, H. and Duchi, J. C. (2019). The importance of better models in stochastic optimization. *Proceedings of the National Academy of Sciences*, 116(46):22924–22930.
- Bartlett, P. L., Jordan, M. I., and McAuliffe, J. D. (2006). Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156.
- Bengio, Y. (2015). Rmsprop and equilibrated adaptive learning rates for nonconvex optimization. *corr abs/1502.04390*.
- Berrada, L., Zisserman, A., and Kumar, M. P. (2020). Training neural networks for and by interpolation. In *ICML*.
- Bertsekas, D. P. and Tsitsiklis, J. N. (1996). *Neuro-Dynamic Programming*. Athena Scientific, 1st edition.
- Bottou, L., Curtis, F. E., and Nocedal, J. (2018). Optimization methods for large-scale machine learning. *SIAM Review*, 60(2):223–311.
- Boyd, S., Xiao, L., and Mutapcic, A. (2003). Subgradient methods. *lecture notes of EE392o, Stanford University, Autumn Quarter*, 2004:2004–2005.
- Cevher, V. and Vũ, B. C. (2019). On the linear convergence of the stochastic gradient method with constant step-size. *Optimization Letters*, 13(5):1177–1187.
- Chang, C.-C. and Lin, C.-J. (2011). LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- Corless, R. M., Gonnet, G. H., Hare, D. E., Jeffrey, D. J., and Knuth, D. E. (1996). On the Lambert W function. *Advances in Computational mathematics*, 5(1):329–359.
- Duchi, J., Hazan, E., and Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(Jul):2121–2159.
- Ghadimi, S. and Lan, G. (2013). Stochastic first- and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368.
- Gower, R. and Richtárik, P. (2015). Randomized iterative methods for linear systems. *SIAM Journal on Matrix Analysis and Applications*, 36(4):1660–1690.
- Gower, R. M., Loizou, N., Qian, X., Sailanbayev, A., Shulgin, E., and Richtárik, P. (2019). SGD: General analysis and improved rates. In *ICML*.
- Hardt, M., Recht, B., and Singer, Y. (2016). Train faster, generalize better: stability of stochastic gradient descent. In *ICML*.
- Harikandeh, R., Ahmed, M. O., Virani, A., Schmidt, M., Konečný, J., and Sallinen, S. (2015). Stop wasting my gradients: Practical SVRG. In *NeurIPS*.
- Hazan, E. and Kakade, S. (2019). Revisiting the polyak step size. *arXiv preprint arXiv:1905.00313*.
- Hazan, E. and Kale, S. (2014). Beyond the regret minimization barrier: optimal algorithms for stochastic

- strongly-convex optimization. *The Journal of Machine Learning Research*, 15(1):2489–2512.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *CVPR*.
- Huang, G., Liu, Z., Van Der Maaten, L., and Weinberger, K. Q. (2017). Densely connected convolutional networks. In *CVPR*.
- Kaczmarz, S. (1937). Angenäherte auflösung von systemen linearer gleichungen. *Bulletin International de l'Academie Polonaise des Sciences et des Lettres*, 35:355–357.
- Karimi, H., Nutini, J., and Schmidt, M. (2016). Linear convergence of gradient and proximal-gradient methods under the Polyak-Lojasiewicz condition. In *ECML-PKDD*.
- Kingma, D. and Ba, J. (2015). Adam: A method for stochastic optimization. In *ICLR*.
- Lei, Y., Hu, T., Li, G., and Tang, K. (2019). Stochastic gradient descent for nonconvex learning without bounded gradient assumptions. *IEEE transactions on neural networks and learning systems*, 31(10):4394–4400.
- Li, X. and Orabona, F. (2019). On the convergence of stochastic gradient descent with adaptive stepsizes. In *AISTATS*.
- Liu, L., Jiang, H., He, P., Chen, W., Liu, X., Gao, J., and Han, J. (2019). On the variance of the adaptive learning rate and beyond. *arXiv preprint arXiv:1908.03265*.
- Lohr, S. L. (2019). *Sampling: Design and Analysis: Design and Analysis*. Chapman and Hall/CRC.
- Loizou, N. (2019). *Randomized iterative methods for linear systems: momentum, inexactness and gossip*. PhD thesis, University of Edinburgh.
- Loizou, N. and Richtárik, P. (2019). Revisiting randomized gossip algorithms: General framework, convergence rates and novel block and accelerated protocols. *arXiv preprint arXiv:1905.08645*.
- Loizou, N. and Richtárik, P. (2020a). Convergence analysis of inexact randomized iterative methods. *SIAM Journal on Scientific Computing*, 42(6):A3979–A4016.
- Loizou, N. and Richtárik, P. (2020b). Momentum and stochastic momentum for stochastic gradient, newton, proximal point and subspace descent methods. *Computational Optimization and Applications*, 77(3):653–710.
- Loshchilov, I. and Hutter, F. (2017). SGDR: Stochastic gradient descent with warm restarts. *ICLR*.
- Luo, L., Xiong, Y., Liu, Y., and Sun, X. (2019). Adaptive gradient methods with dynamic bound of learning rate. In *ICLR*.
- Ma, S., Bassily, R., and Belkin, M. (2018). The power of interpolation: Understanding the effectiveness of sgd in modern over-parametrized learning. In *ICML*.
- Mező, I. and Baricz, Á. (2017). On the generalization of the Lambert W function. *Transactions of the American Mathematical Society*, 369(11):7917–7934.
- Moulines, E. and Bach, F. R. (2011). Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *NeurIPS*.
- Needell, D., Srebro, N., and Ward, R. (2016). Stochastic gradient descent, weighted sampling, and the randomized kaczmarz algorithm. *Mathematical Programming, Series A*, 155(1):549–573.
- Needell, D. and Ward, R. (2017). Batched stochastic gradient descent with weighted sampling. In *Approximation Theory XV, Springer*, volume 204 of *Springer Proceedings in Mathematics & Statistics*, pages 279 – 306.
- Nemirovski, A., Juditsky, A., Lan, G., and Shapiro, A. (2009). Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609.
- Nemirovski, A. and Yudin, D. B. (1978). On Cezari’s convergence of the steepest descent method for approximating saddle point of convex-concave functions. *Soviet Mathematics Doklady*, 19.
- Nemirovski, A. and Yudin, D. B. (1983). *Problem complexity and method efficiency in optimization*. Wiley Interscience.
- Nguyen, L., Nguyen, P. H., van Dijk, M., Richtárik, P., Scheinberg, K., and Takáč, M. (2018). SGD and hogwild! Convergence without the bounded gradients assumption. In *ICML*.
- Nutini, J., Laradji, I., and Schmidt, M. (2017). Let’s make block coordinate descent go fast: Faster greedy rules, message-passing, active-set complexity, and superlinear convergence. *arXiv preprint arXiv:1712.08859*.
- Oberman, A. M. and Prazeres, M. (2019). Stochastic gradient descent with polyak’s learning rate. *arXiv preprint arXiv:1903.08688*.
- Polyak, B. (1987). Introduction to optimization. translations series in mathematics and engineering. *Optimization Software*.
- Rahimi, A. and Recht, B. (2017). Reflections on random kitchen sinks - arg min blog.
- Rakhlin, A., Shamir, O., and Sridharan, K. (2012). Making gradient descent optimal for strongly convex stochastic optimization. In *ICML*.
- Recht, B., Re, C., Wright, S., and Niu, F. (2011). Hogwild: A lock-free approach to parallelizing stochastic gradient descent. In *NeurIPS*.

- Richtárik, P. and Takác, M. (2020). Stochastic reformulations of linear systems: algorithms and convergence theory. *SIAM Journal on Matrix Analysis and Applications*, 41(2):487–524.
- Robbins, H. and Monro, S. (1951). A stochastic approximation method. *The Annals of Mathematical Statistics*, pages 400–407.
- Rolinek, M. and Martius, G. (2018). L4: Practical loss-based stepsize adaptation for deep learning. In *NeurIPS*.
- Schmidt, M., Kim, D., and Sra, S. (2011). 11 projected Newton-type methods in machine learning. *Optimization for Machine Learning*, page 305.
- Schmidt, M. and Roux, N. (2013). Fast convergence of stochastic gradient descent under a strong growth condition. *arXiv preprint arXiv:1308.6370*.
- Shalev-Shwartz, S., Singer, Y., and Srebro, N. (2007). Pegasos: primal estimated subgradient solver for SVM. In *ICML*.
- Shamir, O. and Zhang, T. (2013). Stochastic gradient descent for non-smooth optimization: Convergence results and optimal averaging schemes. In *ICML*.
- Soudry, D., Hoffer, E., Nacson, M. S., Gunasekar, S., and Srebro, N. (2018). The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878.
- Strohmer, T. and Vershynin, R. (2009). A randomized Kaczmarz algorithm with exponential convergence. *J. Fourier Anal. Appl.*, 15(2):262–278.
- Sun, R. and Luo, Z.-Q. (2016). Guaranteed matrix completion via non-convex factorization. *IEEE Transactions on Information Theory*, 62(11):6535–6579.
- Tan, C., Ma, S., Dai, Y.-H., and Qian, Y. (2016). Barzilai-borwein step size for stochastic gradient descent. In *NeurIPS*.
- Vaswani, S., Bach, F., and Schmidt, M. (2019a). Fast and faster convergence of sgd for over-parameterized models and an accelerated perceptron. In *AISTATS*.
- Vaswani, S., Mishkin, A., Laradji, I., Schmidt, M., Gidel, G., and Lacoste-Julien, S. (2019b). Painless stochastic gradient: Interpolation, line-search, and convergence rates. In *NeurIPS*.
- Ward, R., Wu, X., and Bottou, L. (2019). Adagrad stepsizes: Sharp convergence over nonconvex landscapes. In *ICML*.
- Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. (2017). Understanding deep learning requires rethinking generalization. *ICLR*.
- Zhang, M., Lucas, J., Ba, J., and Hinton, G. E. (2019). Lookahead optimizer: k steps forward, 1 step back. In *NeurIPS*.