

November, 1997

LIDS- P 2404

Research Supported By:

NSF grant DMI-9625489

Gradient Convergence in Gradient Methods

Bertsekas, D.P.

Tsitsiklis, J.N.

GRADIENT CONVERGENCE IN GRADIENT METHODS¹

by

Dimitri P. Bertsekas and John N. Tsitsiklis²

Abstract

For the classical gradient method $x_{t+1} = x_t - \gamma_t \nabla f(x_t)$ and several deterministic and stochastic variants, we discuss the issue of convergence of the gradient sequence $\nabla f(x_t)$ and the attendant issue of stationarity of limit points of x_t . We assume that ∇f is Lipschitz continuous, and that the stepsize γ_t diminishes to 0 and satisfies standard stochastic approximation conditions. We show that either $f(x_t) \rightarrow -\infty$ or else $f(x_t)$ converges to a finite value and $\nabla f(x_t) \rightarrow 0$ (with probability 1 in the stochastic case). Existing results assume various boundedness conditions such as boundedness from below of f , or boundedness of $\nabla f(x_t)$, or boundedness of x_t .

¹ Research supported by NSF under Grant DMI-9625489

² Dept. of Electrical Engineering and Computer Science, M.I.T., Cambridge, Mass., 02139.

1. INTRODUCTION

We consider the problem

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && x \in \mathfrak{R}^n, \end{aligned} \tag{1.1}$$

where \mathfrak{R}^n denotes the n -dimensional Euclidean space and $f : \mathfrak{R}^n \mapsto \mathfrak{R}$ is a continuously differentiable scalar function on \mathfrak{R}^n , such that for some constant L we have

$$\|\nabla f(x) - \nabla f(\bar{x})\| \leq L\|x - \bar{x}\|, \quad \forall x, \bar{x} \in \mathfrak{R}^n. \tag{1.2}$$

We focus on the gradient method

$$x_{t+1} = x_t - \gamma_t \nabla f(x_t), \tag{1.3}$$

where the positive stepsize γ_t satisfies

$$\gamma_t \rightarrow 0, \quad \sum_{t=0}^{\infty} \gamma_t = \infty. \tag{1.4}$$

The purpose of the paper is to sharpen the existing convergence theory for this classical and important method, and some of its variations involving deterministic and stochastic errors.

Our main result for the method (1.3) is that either $f(x_t) \rightarrow -\infty$ or else $f(x_t)$ converges to a finite value and $\lim_{t \rightarrow \infty} \nabla f(x_t) = 0$. Furthermore, every limit point of x_t is a stationary point of f . For the case where the stepsize γ_t is chosen by several other rules, such as the minimization and limited minimization rules, or the Armijo and Goldstein rules, these gradient convergence results are known and are relatively easy to show. However, when the stepsize is diminishing, as per Eq. (1.4), our results are stronger than those existing in the literature. This is true even for the deterministic method (1.3), but is particularly so for the case of gradient methods with errors, for which the use of a diminishing stepsize is essential for convergence.

The gradient method variants that we consider involve deterministic and stochastic errors, and scaling of the gradient direction. Such methods include among others, the standard incremental gradient/backpropagation method for neural network training, the convergence of which has been the object of much recent analysis [Luo91], [Gai94], [Gri94], [LuT94], [MaS94], [Man93], [Ber95a] (see the authors' [BeT96] for a discussion of incremental gradient methods and their application to neural network training). They also include the classical Robbins-Monro/stochastic gradient method. In particular, we consider the method

$$x_{t+1} = x_t + \gamma_t(s_t + w_t), \tag{1.5}$$

where s_t is a descent direction satisfying for some positive scalars c_1 and c_2 , and all t ,

$$c_1 \|\nabla f(x_t)\|^2 \leq -\nabla f(x_t)' s_t, \quad \|s_t\| \leq c_2 \|\nabla f(x_t)\|, \quad (1.6)$$

and w_t is an error vector satisfying for some positive scalars p and q , and all t ,

$$\|w_t\| \leq \gamma_t (q + p \|\nabla f(x_t)\|). \quad (1.7)$$

The relation (1.6) is a standard condition in gradient methods, which guarantees that the angle between $\nabla f(x_t)$ and s_t is bounded away from 90 degrees, and also provides a bound to $\|s_t\|$ that is proportional to $\|\nabla f(x_t)\|$. The relation (1.7) bounds the error w_t proportionally to the stepsize and $\|\nabla f(x_t)\|$.

We also consider stochastic variants where w_t are random errors, and the pseudogradient condition of Poljak and Tsypkin [PoT73] is satisfied; see Section 5 for a precise statement of our assumptions. Basically, the entire spectrum of unconstrained gradient methods is considered, with the only restriction being the diminishing stepsize condition (1.4) (which is essential for convergence in the case of gradient methods with errors) and the attendant Lipschitz condition (1.2) [which is necessary for showing any kind of convergence result under the stepsize condition (1.4)].

To place our analysis in perspective, we review the related results of the literature for gradient-like methods with a diminishing stepsize and in the absence of convexity. Our results relate to two types of analyses that can be found in the literature:

- (1) Results that are based on some type of deterministic or stochastic descent argument, such as the use of a Lyapounov function or a supermartingale convergence theorem. All of the results of this type known to us assume that f is bounded below, and in some cases require a boundedness assumption on the sequence $\{x_t\}$ or show only that $\liminf_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0$. By contrast, we show that $\lim_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0$ and we also deal with the case where f is unbounded below, and $\{x_t\}$ is unbounded. In fact, a principal aim of our work has been to avoid any type of boundedness assumption. For example, the classical analysis of Poljak and Tsypkin [PoT73], under essentially the same conditions as ours, shows that if f is bounded below, then $f(x_t)$ converges and $\liminf_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0$ (see Poljak [Pol87], p. 51). The analysis of Gaivoronski [Gai94], for stochastic gradient and incremental gradient methods, under similar conditions to ours shows that $\lim_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0$, but also assumes that $f(x)$ is bounded below and that $\|\nabla f(x)\|$ is bounded over \mathfrak{R}^n . The analysis of Luo and Tseng [LuT94] for the incremental gradient method shows that $\lim_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0$, but also assumes that $f(x)$ is bounded below, and makes some additional assumptions on the

stepsize γ_t . The analyses by Grippo [Gri94], and by Mangasarian and Solodov [MaS94] for the incremental gradient method (with and without a momentum term), make assumptions that are different from ours and include boundedness of the generated sequence x_t .

- (2) Results based on the so-called ODE analysis ([Lju77], [KuC78], [BMP90], [KuY97]) that relate the evolution of the algorithm to the trajectories of a differential equation $dx/dt = h(x)$. For example, if we are dealing with the stochastic steepest descent method $x_{t+1} = x_t - \gamma_t(\nabla f(x_t) + w_t)$, the corresponding ODE is $dx/dt = -\nabla f(x)$. This framework typically involves an explicit or implicit assumption that the average direction of update $h(x)$ is a well-defined function of the current iterate x . It cannot be applied, for example, to a gradient method with diagonal scaling, where the scaling may depend in a complicated way on the past history of the algorithm, unless one works with differential inclusions – rather than differential equations – for which not many results are available. For another example, an asynchronous gradient iteration that updates a single component at a time (selected by some arbitrary or hard to model mechanism) does not lead to a well-defined average direction of update $h(x)$, unless one makes some very special assumptions, e.g., the stepsize assumptions of Borkar [Bor95]. In addition to the above described difficulty, the ODE approach relies on the assumption that the sequence of iterates x_t is bounded or recurrent, something that must be independently verified. Let us also mention the more recent results by Delyon [Del96], which have some similarities with ours: they are proved using a potential function argument and can establish the convergence of $\nabla f(x_t)$ to zero. Similar to the ODE approach, these results assume a well-defined average update direction $h(x)$ and are based on boundedness or recurrence assumptions.

The paper is organized as follows. In the next section, we focus on the scaled gradient method $x_{t+1} = x_t + \gamma_t s_t$, which involves an error-free direction s_t that satisfies condition (1.6). The techniques used for this case are extended in Section 3 to the case where there is a nonrandom error w_t satisfying the condition (1.7). These results are then applied in Section 4 to the case of incremental gradient methods for minimizing the sum of a large number of functions. Finally, in Section 5, we focus on stochastic gradient methods.

2. GRADIENT METHODS WITHOUT ERRORS

Throughout the paper, we focus on the unconstrained minimization of a continuously differen-

table function $f : \mathfrak{R}^n \mapsto \mathfrak{R}$, satisfying for some constant L

$$\|\nabla f(x) - \nabla f(\bar{x})\| \leq L\|x - \bar{x}\|, \quad \forall x, \bar{x} \in \mathfrak{R}^n. \quad (2.1)$$

The proof of the following proposition follows standard arguments to show the known result $\liminf_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0$. The strengthened result, $\lim_{t \rightarrow \infty} \nabla f(x_t) = 0$, is then shown by arguing that if $\|\nabla f(x_t)\|$ exceeds some positive level $\epsilon > 0$ infinitely often, the corresponding cost improvement must be infinite. This line of argument, appropriately modified, is also used in the case of errors in Sections 3 and 5.

Proposition 1: Let x_t be a sequence generated by a gradient method

$$x_{t+1} = x_t + \gamma_t s_t,$$

where s_t satisfies

$$c_1 \|\nabla f(x_t)\|^2 \leq -\nabla f(x_t)' s_t, \quad \|s_t\| \leq c_2 \|\nabla f(x_t)\|, \quad (2.2)$$

for some positive scalars c_1 and c_2 , and all t . Assume that the stepsize γ_t is positive and satisfies

$$\gamma_t \rightarrow 0, \quad \sum_{t=0}^{\infty} \gamma_t = \infty.$$

Then either $f(x_t) \rightarrow -\infty$ or else $f(x_t)$ converges to a finite value and $\lim_{t \rightarrow \infty} \nabla f(x_t) = 0$. Furthermore, every limit point of x_t is a stationary point of f .

Proof: Fix two vectors x and z , let ξ be a scalar parameter, and let $g(\xi) = f(x + \xi z)$. The chain rule yields $(dg/d\xi)(\xi) = z' \nabla f(x + \xi z)$. We have

$$\begin{aligned} f(x+z) - f(x) &= g(1) - g(0) \\ &= \int_0^1 \frac{dg}{d\xi}(\xi) d\xi \\ &= \int_0^1 z' \nabla f(x + \xi z) d\xi \\ &\leq \int_0^1 z' \nabla f(x) d\xi + \left| \int_0^1 z' (\nabla f(x + \xi z) - \nabla f(x)) d\xi \right| \\ &\leq z' \nabla f(x) + \int_0^1 \|z\| \cdot \|\nabla f(x + \xi z) - \nabla f(x)\| d\xi \\ &\leq z' \nabla f(x) + \|z\| \int_0^1 L \xi \|z\| d\xi \\ &= z' \nabla f(x) + \frac{L}{2} \|z\|^2. \end{aligned} \quad (2.3)$$

Applying this relation with $z = \gamma_t s_t$ and using also Eq. (2.2), we have

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) + \gamma_t \nabla f(x_t)' s_t + \frac{\gamma_t^2 L}{2} \|s_t\|^2 \\ &\leq f(x_t) - \gamma_t \left(c_1 - \frac{\gamma_t c_2^2 L}{2} \right) \|\nabla f(x_t)\|^2. \end{aligned}$$

2. Gradient Methods Without Errors

Since $\gamma_t \rightarrow 0$, we have for some positive constant c and all t greater than some index \bar{t} ,

$$f(x_{t+1}) \leq f(x_t) - \gamma_t c \|\nabla f(x_t)\|^2. \quad (2.4)$$

From this relation, we see that for $t \geq \bar{t}$, $f(x_t)$ is monotonically nonincreasing, so either $f(x_t) \rightarrow -\infty$ or $f(x_t)$ converges to a finite value. If the former case holds we are done, so assume the latter case. By adding Eq. (2.4) over all $t \geq \bar{t}$, we obtain

$$c \sum_{t=\bar{t}}^{\infty} \gamma_t \|\nabla f(x_t)\|^2 \leq f(x_{\bar{t}}) - \lim_{t \rightarrow \infty} f(x_t) < \infty.$$

We see that there cannot exist an $\epsilon > 0$ such that $\|\nabla f(x_t)\|^2 > \epsilon$ for all t greater than some \hat{t} , since this would contradict the assumption $\sum_{t=0}^{\infty} \gamma_t = \infty$. Therefore, we must have $\liminf_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0$.

To show that $\lim_{t \rightarrow \infty} \nabla f(x_t) = 0$, assume the contrary; that is, $\limsup_{t \rightarrow \infty} \|\nabla f(x_t)\| > 0$. Then there exists an $\epsilon > 0$ such that $\|\nabla f(x_t)\| < \epsilon/2$ for infinitely many t and also $\|\nabla f(x_t)\| > \epsilon$ for infinitely many t . Therefore, there is an infinite subset of integers \mathcal{T} such that for each $t \in \mathcal{T}$, there exists an integer $i(t) > t$ such that

$$\|\nabla f(x_t)\| < \epsilon/2, \quad \|\nabla f(x_{i(t)})\| > \epsilon,$$

$$\epsilon/2 \leq \|\nabla f(x_i)\| \leq \epsilon, \quad \text{if } t < i < i(t).$$

Since

$$\begin{aligned} \|\nabla f(x_{t+1})\| - \|\nabla f(x_t)\| &\leq \|\nabla f(x_{t+1}) - \nabla f(x_t)\| \\ &\leq L \|x_{t+1} - x_t\| \\ &= \gamma_t L \|s_t\| \\ &\leq \gamma_t L c_2 \|\nabla f(x_t)\|, \end{aligned}$$

it follows that for all $t \in \mathcal{T}$ that are sufficiently large so that $\gamma_t L c_2 < 1$, we have

$$\epsilon/4 \leq \|\nabla f(x_t)\|;$$

otherwise, the condition $\epsilon/2 \leq \|\nabla f(x_{t+1})\|$ would be violated. Without loss of generality, we assume that the above relations hold for all $t \in \mathcal{T}$.

We have for all $t \in \mathcal{T}$, using the condition $\|s_t\| \leq c_2 \|\nabla f(x_t)\|$ and the Lipschitz condition

(2.1),

$$\begin{aligned}
\frac{\epsilon}{2} &\leq \|\nabla f(x_{i(t)})\| - \|\nabla f(x_t)\| \\
&\leq \|\nabla f(x_{i(t)}) - \nabla f(x_t)\| \\
&\leq L\|x_{i(t)} - x_t\| \\
&\leq L \sum_{i=t}^{i(t)-1} \gamma_i \|s_i\| \\
&\leq Lc_2 \sum_{i=t}^{i(t)-1} \gamma_i \|\nabla f(x_i)\| \\
&\leq Lc_2\epsilon \sum_{i=t}^{i(t)-1} \gamma_i,
\end{aligned} \tag{2.5}$$

and finally

$$\frac{1}{2Lc_2} \leq \sum_{i=t}^{i(t)-1} \gamma_i. \tag{2.6}$$

Using Eq. (2.4) for sufficiently large $t \in \mathcal{T}$, and the relation $\|\nabla f(x_i)\| \geq \epsilon/4$ for $i = t, t+1, \dots, i(t)-1$, we have

$$f(x_{i(t)}) \leq f(x_t) - c \left(\frac{\epsilon}{4}\right)^2 \sum_{i=t}^{i(t)-1} \gamma_i, \quad \forall t \in \mathcal{T}. \tag{2.7}$$

Since $f(x_t)$ converges to a finite value, the preceding relation implies that

$$\lim_{t \rightarrow \infty, t \in \mathcal{T}} \sum_{i=t}^{i(t)-1} \gamma_i = 0, \tag{2.8}$$

contradicting Eq. (2.6). Thus, $\lim_{t \rightarrow \infty} \nabla f(x_t) = 0$. Finally, if \bar{x} is a limit point of x_t , then $f(x_t)$ converges to the finite value $f(\bar{x})$. Thus we have $\nabla f(x_t) \rightarrow 0$, implying that $\nabla f(\bar{x}) = 0$.

Q.E.D.

Part of Prop. 1 can be proved if we replace the assumption $\|s_t\| \leq c_2 \|\nabla f(x_t)\|$ [cf. Eq. (2.2)] with the weaker assumption

$$\|s_t\| \leq c_2(1 + \|\nabla f(x_t)\|), \tag{2.9}$$

which allows s_t to be bounded, but not necessarily in proportion to $\|\nabla f(x_t)\|$. Under this weaker assumption, the proof of Prop. 1 can be modified to show that either $f(x_t) \rightarrow -\infty$ or else $\liminf_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0$, but in the latter case the convergence of $\nabla f(x_t)$ to 0 and of $f(x_t)$ to a finite value is unclear. To see this, note that under condition (2.9), the relation (2.4) can take the form

$$f(x_{t+1}) \leq f(x_t) - \gamma_t \xi_1 \|\nabla f(x_t)\|^2 + \gamma_t^2 \xi_2, \quad \forall t \geq \bar{t}, \tag{2.10}$$

3. Deterministic Gradient Methods With Errors

where ξ_1 and ξ_2 are some positive scalars. From this relation, we see that if there exists an $\epsilon > 0$ such that $\|\nabla f(x_t)\|^2 > \epsilon$ for all t greater than some \hat{t} , the term $\gamma_t \xi_1 \|\nabla f(x_t)\|^2$ dominates the term $\gamma_t^2 \xi_2$ in Eq. (2.10), so that the sequence $f(x_t)$ eventually becomes decreasing, leading to the conclusion that $f(x_t) \rightarrow -\infty$ or to a contradiction of the assumption $\sum_{t=0}^{\infty} \gamma_t = \infty$. Therefore, we must either have $f(x_t) \rightarrow -\infty$ or else $\liminf_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0$.

The conclusion of Prop. 1 can be proved in its entirety with the weaker assumption $\|s_t\| \leq c_2(1 + \|\nabla f(x_t)\|)$, provided we require that the stepsize γ_t satisfies in addition $\sum_{t=0}^{\infty} \gamma_t^2 < \infty$. This is shown as a special case of Props. 2 and 4 that follow, by setting $w_t \equiv 0$.

3. DETERMINISTIC GRADIENT METHODS WITH ERRORS

We now extend the results of the preceding section to cover the case where the direction contains an error w_t that is bounded by a multiple of the stepsize γ_t . We will need the following lemma, which we prove for completeness:

Lemma 1: Let Y_t , W_t , and Z_t be three sequences such that W_t is nonnegative for all t . Assume that

$$Y_{t+1} \leq Y_t - W_t + Z_t, \quad t = 0, 1, \dots,$$

and that the series $\sum_{t=0}^T Z_t$ converges as $T \rightarrow \infty$. Then either $Y_t \rightarrow -\infty$, or else Y_t converges to a finite value and $\sum_{t=0}^{\infty} W_t < \infty$.

Proof: Let \bar{t} be any nonnegative integer. By adding the relation $Y_{t+1} \leq Y_t + Z_t$ over all $t \geq \bar{t}$ and by taking the limit superior as $t \rightarrow \infty$, we obtain

$$\limsup_{t \rightarrow \infty} Y_t \leq Y_{\bar{t}} + \sum_{t=\bar{t}}^{\infty} Z_t < \infty.$$

By taking the limit inferior of the right-hand side as $\bar{t} \rightarrow \infty$ and using the fact $\lim_{\bar{t} \rightarrow \infty} \sum_{t=\bar{t}}^{\infty} Z_t = 0$, we obtain

$$\limsup_{t \rightarrow \infty} Y_t \leq \liminf_{\bar{t} \rightarrow \infty} Y_{\bar{t}} < \infty.$$

This implies that either $Y_t \rightarrow -\infty$ or else Y_t converges to a finite value. In the latter case, by adding the relation $Y_{i+1} \leq Y_i - W_i + Z_i$ from $i = 0$ to $i = t$, we obtain

$$\sum_{i=0}^t W_i \leq Y_0 + \sum_{i=0}^t Z_i - Y_{t+1}, \quad t = 0, 1, \dots,$$

which implies that $\sum_{i=0}^{\infty} W_i \leq Y_0 + \sum_{i=0}^{\infty} Z_i - \lim_{t \rightarrow \infty} Y_t < \infty$. **Q.E.D.**

We have the following result:

Proposition 2: Let x_t be a sequence generated by the method

$$x_{t+1} = x_t + \gamma_t(s_t + w_t),$$

where s_t is a descent direction satisfying for some positive scalars c_1 and c_2 , and all t ,

$$c_1 \|\nabla f(x_t)\|^2 \leq -\nabla f(x_t)'s_t, \quad \|s_t\| \leq c_2(1 + \|\nabla f(x_t)\|), \quad (3.1)$$

and w_t is an error vector satisfying for some positive scalars p and q , and all t ,

$$\|w_t\| \leq \gamma_t(q + p\|\nabla f(x_t)\|). \quad (3.2)$$

Assume that the stepsize γ_t is positive and satisfies

$$\sum_{t=0}^{\infty} \gamma_t = \infty, \quad \sum_{t=0}^{\infty} \gamma_t^2 < \infty.$$

Then either $f(x_t) \rightarrow -\infty$ or else $f(x_t)$ converges to a finite value and $\lim_{t \rightarrow \infty} \nabla f(x_t) = 0$. Furthermore, every limit point of x_t is a stationary point of f .

Proof: The proof is similar to the proof of Prop. 1, with the appropriate modifications to deal with the error vectors w_t . We apply Eq. (2.3) with $x = x_t$ and $z = \gamma_t(s_t + w_t)$. We obtain

$$f(x_{t+1}) \leq f(x_t) + \gamma_t \nabla f(x_t)'(s_t + w_t) + \frac{\gamma_t^2 L}{2} \|s_t + w_t\|^2.$$

Using our assumptions, we have

$$\begin{aligned} \nabla f(x_t)'(s_t + w_t) &\leq -c_1 \|\nabla f(x_t)\|^2 + \|\nabla f(x_t)\| \|w_t\| \\ &\leq -c_1 \|\nabla f(x_t)\|^2 + \gamma_t q \|\nabla f(x_t)\| + \gamma_t p \|\nabla f(x_t)\|^2. \end{aligned}$$

Furthermore, using the relations $\|s_t\|^2 \leq 2c_2^2(1 + \|\nabla f(x_t)\|^2)$ and $\|w_t\|^2 \leq 2\gamma_t^2(q^2 + p^2\|\nabla f(x_t)\|^2)$, which follow from Eqs. (3.1) and (3.2), respectively, we have

$$\begin{aligned} \|s_t + w_t\|^2 &\leq 2\|s_t\|^2 + 2\|w_t\|^2 \\ &\leq 4c_2^2(1 + \|\nabla f(x_t)\|^2) + 4\gamma_t^2(q^2 + p^2\|\nabla f(x_t)\|^2). \end{aligned}$$

Combining the above relations, we obtain

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) - \gamma_t(c_1 - \gamma_t p - 2\gamma_t c_2^2 L - 2\gamma_t^3 p^2 L) \|\nabla f(x_t)\|^2 \\ &\quad + \gamma_t^2 q \|\nabla f(x_t)\| + 2\gamma_t^2 c_2^2 L + 2\gamma_t^4 q^2 L. \end{aligned}$$

3. Deterministic Gradient Methods With Errors

Since $\gamma_t \rightarrow 0$, we have for some positive constant c and all t sufficiently large

$$f(x_{t+1}) \leq f(x_t) - \gamma_t c \|\nabla f(x_t)\|^2 + \gamma_t^2 q \|\nabla f(x_t)\| + 2\gamma_t^2 c_2^2 L + 2\gamma_t^4 q^2 L.$$

Using the inequality $\|\nabla f(x_t)\| \leq 1 + \|\nabla f(x_t)\|^2$, the above relation yields for all t

$$f(x_{t+1}) \leq f(x_t) - \gamma_t (c - \gamma_t q) \|\nabla f(x_t)\|^2 + \gamma_t^2 (q + 2c_2^2 L) + 2\gamma_t^4 q^2 L. \quad (3.3)$$

Consider Eq. (3.3) for all t sufficiently large so that $c - \gamma_t q > 0$. By using Lemma 1 and the assumption $\sum_{t=0}^{\infty} \gamma_t^2 < \infty$, we see that either $f(x_t) \rightarrow -\infty$ or else $f(x_t)$ converges and

$$\sum_{t=0}^{\infty} \gamma_t \|\nabla f(x_t)\|^2 < \infty. \quad (3.4)$$

If there existed an $\epsilon > 0$ and an integer \bar{t} such that $\|\nabla f(x_t)\| \geq \epsilon$ for all $t \geq \bar{t}$, we would have

$$\sum_{t=\bar{t}}^{\infty} \gamma_t \|\nabla f(x_t)\|^2 \geq \epsilon^2 \sum_{t=\bar{t}}^{\infty} \gamma_t = \infty,$$

which contradicts Eq. (3.4). Therefore, $\liminf_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0$.

The proof of $\lim_{t \rightarrow \infty} \nabla f(x_t) = 0$ now proceeds as in the proof of Prop. 1, by assuming that $\limsup_{t \rightarrow \infty} \|\nabla f(x_t)\| > \epsilon > 0$, in order to reach a contradiction. In particular, using the condition $\|s_t\| \leq c_2(1 + \|\nabla f(x_t)\|)$ in place of $\|s_t\| \leq c_2 \|\nabla f(x_t)\|$, Eq. (2.5) takes the form

$$\frac{\epsilon}{2} \leq Lc_2(1 + \epsilon) \sum_{i=t}^{i(t)-1} \gamma_i,$$

and Eq. (2.6) takes the form

$$\frac{\epsilon}{2Lc_2(1 + \epsilon)} \leq \sum_{i=t}^{i(t)-1} \gamma_i. \quad (3.5)$$

Using Eq. (3.3) in place of Eq. (2.4), we see that Eq. (2.7) becomes

$$f(x_{i(t)}) \leq f(x_t) - c \left(\frac{\epsilon}{4}\right)^2 \sum_{i=t}^{i(t)-1} \gamma_i + \xi \sum_{i=t}^{i(t)-1} \gamma_i^2 + \zeta \sum_{i=t}^{i(t)-1} \gamma_i^4, \quad \forall t \in \mathcal{T},$$

for appropriate positive scalars ξ and ζ . Using the already shown convergence of $f(x_t)$ and the assumption $\sum_{t=0}^{\infty} \gamma_t^2 < \infty$, this relation still implies that

$$\lim_{t \rightarrow \infty, t \in \mathcal{T}} \sum_{i=t}^{i(t)-1} \gamma_i = 0,$$

[cf. Eq. (2.8)], and contradicts Eq. (3.5). **Q.E.D.**

4. INCREMENTAL GRADIENT METHODS

In this section, we apply the results of the preceding section to the case where f has the form

$$f(x) = \sum_{i=1}^m f_i(x),$$

where $f_i : \mathfrak{R}^n \mapsto \mathfrak{R}$ is for every i , a continuously differentiable function satisfying the Lipschitz condition

$$\|\nabla f_i(x) - \nabla f_i(\bar{x})\| \leq L\|x - \bar{x}\|, \quad \forall x, \bar{x} \in \mathfrak{R}^n, \quad (4.1)$$

for some constant L .

In situations where there are many component functions f_i , it may be attractive to use an incremental method that does not wait to process the entire set of components before updating x ; instead, the method cycles through the components in sequence and updates the estimate of x after each component is processed. In particular, given x_t , we may obtain x_{t+1} as

$$x_{t+1} = \psi_m,$$

where ψ_m is obtained at the last step of the algorithm

$$\psi_i = \psi_{i-1} - \gamma_t \nabla f_i(\psi_{i-1}), \quad i = 1, \dots, m, \quad (4.2)$$

and

$$\psi_0 = x_t. \quad (4.3)$$

This method can be written as

$$x_{t+1} = x_t - \gamma_t \sum_{i=1}^m \nabla f_i(\psi_{i-1}). \quad (4.4)$$

It is referred to as the *incremental gradient method*, and it is used extensively in the training of neural networks. It should be compared with the ordinary gradient method, which is

$$x_{t+1} = x_t - \gamma_t \nabla f(x_t) = x_t - \gamma_t \sum_{i=1}^m \nabla f_i(x_t). \quad (4.5)$$

Thus, a cycle of the incremental gradient method through the components f_i differs from an ordinary gradient iteration only in that the evaluation of ∇f_i is done at the corresponding current estimates ψ_{i-1} rather than at the estimate x_t available at the start of the cycle. The advantages of incrementalism in enhancing the speed of convergence (at least in the early stages of the method) are well-known; see for example the discussion in [Ber95], [BeT96].

The main idea of the following convergence proof is that the incremental gradient method can be viewed as the regular gradient iteration where the gradient is perturbed by an error term that is proportional to the stepsize. In particular, if we compare the incremental method (4.4) with the ordinary gradient method (4.5), we see that the error term in the gradient direction is bounded by

$$\sum_{i=1}^m \|\nabla f_i(\psi_{i-1}) - \nabla f_i(x_t)\|.$$

In view of our Lipschitz assumption (4.1), this term is bounded by

$$L \sum_{i=1}^m \|\psi_{i-1} - x_t\|,$$

which from Eq. (4.2), is seen to be proportional to γ_t (a more precise argument is given below).

Proposition 3: Let x_t be a sequence generated by the incremental gradient method (4.2)-(4.4). Assume that for some positive constants C , and D , and all $i = 1, \dots, m$, we have

$$\|\nabla f_i(x)\| \leq C + D\|\nabla f(x)\|, \quad \forall x \in \mathfrak{R}^n. \quad (4.6)$$

Assume also that

$$\sum_{t=0}^{\infty} \gamma_t = \infty, \quad \sum_{t=0}^{\infty} \gamma_t^2 < \infty.$$

Then either $f(x_t) \rightarrow -\infty$ or else $f(x_t)$ converges to a finite value and $\lim_{t \rightarrow \infty} \nabla f(x_t) = 0$. Furthermore, every limit point of x_t is a stationary point of f .

Proof: We formulate the incremental gradient method as a gradient method with errors that are proportional to the stepsize, and then apply Prop. 2. For simplicity we will assume that there are only two functions f_i ; that is, $m = 2$. The proof is similar when $m > 2$. We have

$$\psi_1 = x_t - \gamma_t \nabla f_1(x_t),$$

$$x_{t+1} = \psi_1 - \gamma_t \nabla f_2(\psi_1).$$

By adding these two relations, we obtain

$$x_{t+1} = x_t + \gamma_t(-\nabla f(x_t) + w_t),$$

where

$$w_t = \nabla f_2(x_t) - \nabla f_2(\psi_1).$$

We have

$$\|w_t\| \leq L\|x_t - \psi_1\| = \gamma_t L \|\nabla f_1(x_t)\| \leq \gamma_t(LC + LD\|\nabla f(x_t)\|).$$

Thus Prop. 3 applies. **Q.E.D.**

Condition (4.6) is guaranteed to hold if each f_k is of the form

$$f_k(x) = x'Q_kx + g'_kx + h_k,$$

where each Q_k is a positive semidefinite matrix, each g_k is a vector, and each h_k is a scalar. (This is the generic situation encountered in linear least squares problems.) If $\sum_{k=1}^K Q_k$ is positive definite, there exists a unique minimum to which the algorithm must converge. In the absence of positive definiteness, we obtain $\nabla f(x_t) \rightarrow 0$ if the optimal cost is finite. If on the other hand the optimal cost is $-\infty$, it can be shown that $\|\nabla f(x)\| \geq \alpha$ for some $\alpha > 0$ and for all x . This implies that $f(x) \rightarrow -\infty$ and that $\|x\| \rightarrow \infty$.

5. STOCHASTIC GRADIENT METHODS

In this section, we study stochastic gradient methods. Our main result is similar to Proposition 2, except that we let the noise term w_t be of a stochastic nature. Once more, we will prove that $f(x_t)$ converges and, if the limit is finite, $\nabla f(x_t)$ converges to 0. We comment on the technical issues that arise in establishing such a result. The sequence $f(x_t)$ can be shown to be approximately a supermartingale. However, the variance of the underlying noise is allowed to grow with $\|\nabla f(x_t)\|$ and can therefore be unbounded. Furthermore, since no lower bound on $f(x_t)$ is assumed, the supermartingale convergence theorem or its variants cannot be used in a simple manner. Our approach is to show that whenever $\|\nabla f(x_t)\|$ is large, it remains so for a sufficiently long time interval, guaranteeing a decrease in the value of $f(x_t)$ which is significant and dominates the noise effects.

Proposition 4: Let x_t be a sequence generated by the method

$$x_{t+1} = x_t + \gamma_t(s_t + w_t),$$

where γ_t is a deterministic positive stepsize, s_t is a descent direction, and w_t is a random noise term. Let \mathcal{F}_t be an increasing sequence of σ -fields. We assume the following:

- (a) x_t and s_t are \mathcal{F}_t -measurable.
- (b) There exist positive scalars c_1 and c_2 such that

$$c_1\|\nabla f(x_t)\|^2 \leq -\nabla f(x_t)'s_t, \quad \|s_t\| \leq c_2(1 + \|\nabla f(x_t)\|), \quad \forall t. \quad (5.1)$$

(c) We have, for all t , and with probability one,

$$E[w_t \mid \mathcal{F}_t] = 0, \quad (5.2)$$

$$E[\|w_t\|^2 \mid \mathcal{F}_t] \leq A(1 + \|\nabla f(x_t)\|^2), \quad (5.3)$$

where A is a positive deterministic constant.

(d) We have

$$\sum_{t=0}^{\infty} \gamma_t = \infty, \quad \sum_{t=0}^{\infty} \gamma_t^2 < \infty.$$

Then, either $f(x_t) \rightarrow -\infty$ or else $f(x_t)$ converges to a finite value and $\lim_{t \rightarrow \infty} \nabla f(x_t) = 0$. Furthermore, every limit point of x_t is a stationary point of f .

Remarks:

- (a) The σ -field \mathcal{F}_t should be interpreted as the history of the algorithm up to time t , just before w_t is generated. In particular, conditioning on \mathcal{F}_t can be thought of as conditioning on $x_0, s_0, w_0, \dots, x_{t-1}, s_{t-1}, w_{t-1}, x_t, s_t$.
- (b) Strictly speaking, the conclusions of the proposition only hold “with probability 1.” For simplicity, an explicit statement of this qualification will often be omitted.
- (c) Our assumptions on w_t are of the same type as those considered in [PoT73].

Proof: We apply Eq. (2.3) with $x = x_t$ and $z = \gamma_t(s_t + w_t)$. We obtain

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) + \gamma_t \nabla f(x_t)'(s_t + w_t) + \frac{\gamma_t^2 L}{2} \|s_t + w_t\|^2 \\ &\leq f(x_t) - \gamma_t c_1 \|\nabla f(x_t)\|^2 + \gamma_t \nabla f(x_t)' w_t + \gamma_t^2 L (\|s_t\|^2 + \|w_t\|^2) \\ &\leq f(x_t) - \gamma_t c_1 \|\nabla f(x_t)\|^2 + \gamma_t \nabla f(x_t)' w_t + \gamma_t^2 2Lc_2^2 \\ &\quad + \gamma_t^2 2Lc_2^2 \|\nabla f(x_t)\|^2 + \gamma_t^2 L \|w_t\|^2 \\ &\leq f(x_t) - \gamma_t \frac{c_1}{2} \|\nabla f(x_t)\|^2 + \gamma_t \nabla f(x_t)' w_t + \gamma_t^2 2Lc_2^2 + \gamma_t^2 L \|w_t\|^2, \end{aligned} \quad (5.4)$$

where the last inequality is only valid when t is large enough so that $\gamma_t 2Lc_2^2 \leq c_1/2$. Without loss of generality, we will assume that this is the case for all $t \geq 0$.

Let $\delta > 0$ be an arbitrary positive number that will be kept constant until the very end of this proof. Let η be a positive constant defined, in terms of δ , by

$$\eta c_2 \left(\frac{1}{\delta} + 2 \right) + \eta = \frac{1}{2L}. \quad (5.5)$$

We will partition the set of all times t (the nonnegative integers) into a set S of times at which $\|\nabla f(x_t)\|$ is “small” and intervals $I_k = \{\tau_k, \tau_k + 1, \dots, \tau'_k\}$ during which $\|\nabla f(x_t)\|$ stays “large.”

The definition of the times τ_k and τ'_k is recursive and is initialized by letting $\tau'_0 = -1$. We then let, for $k = 1, 2, \dots$,

$$\tau_k = \min \{t > \tau'_{k-1} \mid \|\nabla f(x_t)\| \geq \delta\}.$$

(We leave τ_k undefined if $\|\nabla f(x_t)\| < \delta$ for all $t > \tau'_{k-1}$.) We also let

$$\tau'_k = \max \left\{ t \geq \tau_k \mid \sum_{i=\tau_k}^t \gamma_i \leq \eta, \text{ and } \frac{\|\nabla f(x_{\tau_k})\|}{2} \leq \|\nabla f(x_t)\| \leq 2\|\nabla f(x_{\tau_k})\| \right\}.$$

We say that the interval I_k is *full* if $\sum_{t=\tau_k}^{\tau'_k+1} \gamma_t > \eta$. Let S be the set of all times that do not belong to any of the intervals I_k .

We define a sequence G_t , used to scale the noise terms w_t , by

$$G_t = \begin{cases} \delta, & \text{if } t \in S, \\ \|\nabla f(x_{\tau_k})\| = H_k, & \text{if } t \in I_k, \end{cases}$$

where the last equality should be taken as the definition of H_k . In particular, G_t is constant during an interval I_t . Note that $G_t \geq \delta$ for all t .

We now collect a few observations that are direct consequences of our definitions.

(P1) For all $t \in S$, we have $\|\nabla f(x_t)\| < \delta = G_t$.

(P2) For all $t \in I_k$, we have

$$\frac{G_t}{2} = \frac{H_k}{2} \leq \|\nabla f(x_t)\| \leq 2H_k = 2G_t.$$

Combining this with (P1), we also see that the ratio $\|\nabla f(x_t)\|/G_t$ is bounded above by 2.

(P3) If τ_k is defined and I_k is a full interval, then

$$\frac{\eta}{2} \leq \eta - \gamma_{\tau'_k+1} < \sum_{t=\tau_k}^{\tau'_k} \gamma_t \leq \eta, \quad (5.6)$$

where the leftmost inequality holds when k is large enough so that $\gamma_{\tau'_k+1} \leq \eta/2$. Without loss of generality, we will assume that this condition actually holds for all k .

(P4) The value of G_t is completely determined by x_0, x_1, \dots, x_t and is therefore \mathcal{F}_t -measurable. Similarly, the indicator function

$$\chi_t = \begin{cases} 1, & \text{if } t \in S, \\ 0, & \text{otherwise,} \end{cases}$$

is also \mathcal{F}_t -measurable.

Lemma 2: Let r_t be a sequence of random variables with each r_t being \mathcal{F}_{t+1} -measurable, and suppose that $E[r_t | \mathcal{F}_t] = 0$ and $E[\|r_t\|^2 | \mathcal{F}_t] \leq B$, where B is some deterministic constant. Then, the sequences

$$\sum_{t=0}^T \gamma_t r_t \quad \text{and} \quad \sum_{t=0}^T \gamma_t^2 \|r_t\|^2, \quad T = 0, 1, \dots,$$

converge to finite limits (with probability 1).

Proof: It is seen that $\sum_{t=0}^T \gamma_t r_t$ is a martingale whose variance is bounded by $B \sum_{t=0}^{\infty} \gamma_t^2$. It must therefore converge, by the martingale convergence theorem. Furthermore,

$$E \left[\sum_{t=0}^{\infty} \gamma_t^2 \|r_t\|^2 \right] \leq B \sum_{t=0}^{\infty} \gamma_t^2 < \infty,$$

which shows that $\sum_{t=0}^{\infty} \gamma_t^2 \|r_t\|^2$ is finite with probability 1. This establishes convergence of the second sequence. **Q.E.D.**

Using Lemma 2, we obtain the following:

Lemma 3: The following sequences converge (with probability 1):

- (a) $\sum_{t=0}^T \chi_t \gamma_t \nabla f(x_t)' w_t$;
- (b) $\sum_{t=0}^T \gamma_t \frac{w_t}{G_t}$;
- (c) $\sum_{t=0}^T \gamma_t \frac{\nabla f(x_t)' w_t}{G_t^2}$;
- (d) $\sum_{t=0}^T \gamma_t^2 \frac{\|w_t\|^2}{G_t^2}$;
- (e) $\sum_{t=0}^T \gamma_t^2 \chi_t \|w_t\|^2$.

Proof: (a) Let $r_t = \chi_t \nabla f(x_t)' w_t$. Since χ_t and $\nabla f(x_t)$ are \mathcal{F}_t -measurable and $E[w_t | \mathcal{F}_t] = 0$, we obtain $E[r_t | \mathcal{F}_t] = 0$. Whenever $\chi_t = 1$, we have $\|\nabla f(x_t)\| \leq \delta$ and $E[\|w_t\|^2 | \mathcal{F}_t] \leq A(1 + \delta^2)$. It follows easily that $E[\|r_t\|^2 | \mathcal{F}_t]$ is bounded. The result follows from Lemma 2.

(b) Let $r_t = w_t/G_t$. Since G_t is \mathcal{F}_t -measurable and $E[w_t | \mathcal{F}_t] = 0$, we obtain $E[r_t | \mathcal{F}_t] = 0$. Furthermore,

$$E[\|r_t\|^2 | \mathcal{F}_t] \leq \frac{A(1 + \|\nabla f(x_t)\|^2)}{G_t^2}.$$

Since the ratio $\|\nabla f(x_t)\|/G_t$ is bounded above [cf. observation (P2)], Lemma 2 applies and establishes the desired convergence result.

(c) Let $r_t = \nabla f(x_t)'w_t/G_t^2$. Note that

$$\frac{\nabla f(x_t)'w_t}{G_t^2} \leq \frac{\|\nabla f(x_t)\| \cdot \|w_t\|}{G_t^2} \leq 2 \frac{\|w_t\|}{G_t}.$$

The ratio in the left-hand side has bounded conditional second moment, by the same argument as in the proof of part (b). The desired result follows from Lemma 2.

(d) This follows again from Lemma 2. The needed assumptions have already been verified while proving part (b).

(e) This follows from Lemma 2 because $\chi_t w_t$ has bounded conditional second moment, by an argument similar to the one used in the proof of part (a). **Q.E.D.**

We now assume that we have removed the zero probability set of sample paths for which the series in Lemma 3 do not converge. For the remainder of the proof, we will concentrate on a single sample path outside this zero probability set. Let ϵ be a positive constant that satisfies

$$\epsilon \leq \eta, \quad 2\epsilon + 2L\epsilon \leq \frac{c_1\eta}{48}, \quad 4Lc_2^2\epsilon \leq \frac{c_1\delta^2\eta}{48}. \quad (5.7)$$

Let us choose some t_0 after which all of the series in Lemma 3, as well as the series $\sum_{t=0}^T \gamma_t^2$, stay within ϵ from their limits.

Lemma 4: Let t_0 be as above. If τ_k is defined and is larger than t_0 , then the interval I_k is full.

Proof: Recall that for $t \in I_k = \{\tau_k, \dots, \tau'_k\}$ we have $G_t = H_k = \|\nabla f(x_{\tau_k})\| \geq \delta$ and $\|s_t\| \leq c_2(1 + \|\nabla f(x_t)\|) \leq c_2(1 + 2H_k)$. Therefore,

$$\begin{aligned} \|x_{\tau'_k+1} - x_{\tau_k}\| &\leq \sum_{t=\tau_k}^{\tau'_k} \gamma_t \|s_t\| + \left\| \sum_{t=\tau_k}^{\tau'_k} \gamma_t w_t \right\| \\ &= \sum_{t=\tau_k}^{\tau'_k} \gamma_t \|s_t\| + H_k \left\| \sum_{t=\tau_k}^{\tau'_k} \gamma_t \frac{w_t}{G_t} \right\| \\ &\leq \eta c_2(1 + 2H_k) + H_k \epsilon \\ &\leq \eta c_2 H_k \left(\frac{1}{\delta} + 2 \right) + \eta H_k \\ &= \frac{H_k}{2L}, \end{aligned}$$

where the last equality follows from our choice of η [cf. Eq. (5.5)]. Thus,

$$\|\nabla f(x_{\tau'_k+1}) - \nabla f(x_{\tau_k})\| \leq L \|x_{\tau'_k+1} - x_{\tau_k}\| \leq \frac{H_k}{2} = \frac{\|\nabla f(x_{\tau_k})\|}{2},$$

which implies that

$$\frac{1}{2}\|\nabla f(x_{\tau_k})\| \leq \|\nabla f(x_{\tau'_k+1})\| \leq 2\|\nabla f(x_{\tau_k})\|.$$

If we also had $\sum_{t=\tau_k}^{\tau'_k+1} \gamma_t \leq \eta$, then $\tau'_k + 1$ should be an element of I_k , which it isn't. This shows that $\sum_{t=\tau_k}^{\tau'_k+1} \gamma_t > \eta$, and I_k is a full interval. **Q.E.D.**

Our next lemma shows that after a certain time, $f(x_t)$ is guaranteed to decrease by at least a constant amount during full intervals.

Lemma 5: Let t_0 be the same as earlier. If τ_k is defined and larger than t_0 , then

$$f(x_{\tau'_k+1}) \leq f(x_{\tau_k}) - h,$$

where h is a positive constant that only depends on δ .

Proof: Note that I_k is a full interval, by Lemma 4. Using Eq. (5.4), we have

$$f(x_{t+1}) - f(x_t) \leq -\gamma_t \frac{c_1}{2} \|\nabla f(x_t)\|^2 + \gamma_t \nabla f(x_t)' w_t + \gamma_t^2 2Lc_2^2 + \gamma_t^2 L \|w_t\|^2.$$

We will sum (from τ_k to τ'_k) the terms in the right-hand side of the above inequality, and provide suitable upper bounds. Recall that for $t \in I_k$, we have $\|\nabla f(x_t)\| \geq H_k/2$. Thus, using also Eq. (5.6),

$$-\sum_{t=\tau_k}^{\tau'_k} \gamma_t \frac{c_1}{2} \|\nabla f(x_t)\|^2 \leq -\frac{c_1 H_k^2}{8} \sum_{t=\tau_k}^{\tau'_k} \gamma_t \leq -\frac{c_1 H_k^2 \eta}{16}. \quad (5.8)$$

Furthermore,

$$\sum_{t=\tau_k}^{\tau'_k} \gamma_t \nabla f(x_t)' w_t \leq 2H_k^2 \epsilon, \quad (5.9)$$

which follows from the convergence of the series in Lemma 3(c) and the assumption that after time t_0 the series is within ϵ of its limit. By a similar argument based on Lemma 3(d), we also have

$$L \sum_{t=\tau_k}^{\tau'_k} \gamma_t^2 \|w_t\|^2 \leq 2LH_k^2 \epsilon. \quad (5.10)$$

Finally,

$$2Lc_2^2 \sum_{t=\tau_k}^{\tau'_k} \gamma_t^2 \leq 4Lc_2^2 \epsilon. \quad (5.11)$$

We add Eqs. (5.8)-(5.11) and obtain

$$\begin{aligned} f(x_{\tau'_k+1}) &\leq f(x_{\tau_k}) - \frac{c_1 \eta H_k^2}{16} + (2\epsilon + 2L\epsilon)H_k^2 + 4Lc_2^2 \epsilon \\ &\leq f(x_{\tau_k}) - \frac{2c_1 \eta H_k^2}{48} + \frac{c_1 \eta \delta^2}{48} \\ &\leq f(x_{\tau_k}) - \frac{c_1 \eta \delta^2}{48}. \end{aligned}$$

The second inequality made use of (5.7); the third made use of $H_k \geq \delta$. **Q.E.D.**

Lemma 6: For almost every sample path, $f(x_t)$ converges to a finite value or to $-\infty$. If $\lim_{t \rightarrow \infty} f(x_t) \neq -\infty$, then $\limsup_{t \rightarrow \infty} \|\nabla f(x_t)\| \leq \delta$.

Proof: Suppose that there are only finitely many intervals I_k and, in particular,

$$\limsup_{t \rightarrow \infty} \|\nabla f(x_t)\| \leq \delta.$$

Let t^* be some time such that $t \in S$ for all $t \geq t^*$. We then have $\chi_t = 1$ for all $t \geq t^*$. We use Eq. (5.4) to obtain

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) + \gamma_t \chi_t \nabla f(x_t)' w_t + \gamma_t^2 2Lc_2^2 + \chi_t \gamma_t^2 L \|w_t\|^2 \\ &= f(x_t) + Z_t, \quad \text{for } t \geq t^*, \end{aligned}$$

where the last equality can be taken as the definition of Z_t . Using parts (a) and (e) of Lemma 3, the series $\sum_t Z_t$ converges. Lemma 1 then implies that $f(x_t)$ converges to a finite value or to $-\infty$. This proves Lemma 6 for the case where there are finitely many intervals.

We consider next the case where there are infinitely many intervals. We will prove that $f(x_t)$ converges to $-\infty$. We first establish such convergence along a particular subsequence. Let $\mathcal{T} = S \cup \{\tau_1, \tau_2, \dots\}$. We will show that the sequence $\{f(x_t)\}_{t \in \mathcal{T}}$ converges to $-\infty$. To see why this must be the case, notice that whenever $t \in S$, we have $f(x_{t+1}) \leq f(x_t) + Z_t$, where Z_t is as in the preceding paragraph and is summable. Also, whenever $t \in \mathcal{T}$ but $t \notin S$, then $t = \tau_k$, for some k , and the next element of \mathcal{T} is the time $\tau'_k + 1$. Using Lemma 5, $f(x_t)$ decreases by at least h during this interval (for k large enough). We are now in the situation captured by Lemma 1, with $W_t = h$ whenever $t = \tau_k$. The convergence of the subsequence $\{f(x_t)\}_{t \in \mathcal{T}}$ follows. Furthermore, since $W_t = h$ infinitely often, the limit can only be $-\infty$.

Having shown that $f(x_{\tau_k})$ converges to $-\infty$, it now remains to show that the fluctuations of $f(x_t)$ during intervals I_k cannot be too large. Because the technical steps involved here are very similar to those given earlier, we only provide an outline. In order to carry out this argument, we consider the events that immediately precede an interval I_k .

Let us first consider the case where I_k is preceded by an element of S , i.e., $\tau_k - 1 \in S$. By replicating the first half of the proof of Lemma 4, we can show that $x_t - x_{\tau_k - 1}$, for $t \in I_k$, is bounded by a constant multiple of δ (for k large enough). Since $\|\nabla f(x_{\tau_k - 1})\| \leq \delta$, this leads to a $c\delta^2$ bound on the difference $f(x_t) - f(x_{\tau_k - 1})$, where c is some absolute constant. Since $f(x_{\tau_k - 1}) \rightarrow -\infty$, the same must be true for $f(x_t)$, $t \in I_k$.

Let us now consider the case where I_k is immediately preceded by an interval I_{k-1} . By replicating the proof of Lemma 5 (with a somewhat smaller choice of ϵ), we can show that (for k

6. The Incremental Gradient Method Revisited

large enough) we will have $f(x_t) \leq f(x_{\tau_{k-1}})$ for all $t \in I_k$. Once more, since $f(x_{\tau_{k-1}})$ converges to $-\infty$, the same must be true for $f(x_t)$, $t \in I_k$. **Q.E.D.**

According to Lemma 6, $f(x_t)$ converges and if

$$\lim_{t \rightarrow \infty} f(x_t) \neq -\infty,$$

then $\limsup_{t \rightarrow \infty} \|\nabla f(x_t)\| \leq \delta$. Since this has been proved for an arbitrary $\delta > 0$, we conclude that if $\lim_{t \rightarrow \infty} f(x_t) \neq -\infty$, then $\limsup_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0$, that is, $\nabla f(x_t) \rightarrow 0$.

Finally, if x^* is a limit point of x_t , this implies that $f(x_t)$ has a subsequence that converges to $f(x^*)$. Therefore, the limit of the entire sequence $f(x_t)$, which we have shown to exist, must be finite and equal to $f(x^*)$. We have shown that in this case $\nabla f(x_t)$ converges to zero. By taking the limit of $\nabla f(x_t)$ along a sequence of times such that x_t converges to x^* , we conclude that $\nabla f(x^*) = 0$. **Q.E.D.**

6. THE INCREMENTAL GRADIENT METHOD REVISITED

We now provide an alternative view of the incremental gradient method that was discussed in Section 4.

Consider again a cost function f of the form

$$f(x) = \frac{1}{m} \sum_{i=1}^m f_i(x),$$

where each f_i is a function from \mathfrak{R}^n into \mathfrak{R} that satisfies the Lipschitz condition (4.1). In contrast to the setting of Section 4, we now assume that each update is based on a single component function f_i , chosen at random. More specifically, let $k(t)$, $t = 1, 2, \dots$, be a sequence of independent random variables, each distributed uniformly over the set $\{1, \dots, m\}$. The algorithm under consideration is

$$x_{t+1} = x_t - \gamma_t \nabla f_{k(t)}(x_t), \tag{6.1}$$

where γ_t is a nonnegative scalar stepsize. We claim that this is a special case of the stochastic gradient algorithm. Indeed, the algorithm (6.1) can be rewritten as

$$x_{t+1} = x_t - \frac{\gamma_t}{m} \sum_{i=1}^m \nabla f_i(x_t) - \gamma_t \left(\nabla f_{k(t)}(x_t) - \frac{1}{m} \sum_{i=1}^m \nabla f_i(x_t) \right),$$

which is of the form

$$x_{t+1} = x_t - \gamma_t \nabla f(x_t) - \gamma_t w_t,$$

where

$$w_t = \nabla f_{k(t)}(x_t) - \frac{1}{m} \sum_{i=1}^m \nabla f_i(x_t).$$

We now verify that w_t satisfies the assumptions of Proposition 4. Due to the way that $k(t)$ is chosen, we have

$$E[\nabla f_{k(t)}(x_t) \mid \mathcal{F}_t] = \frac{1}{m} \sum_{i=1}^m \nabla f_i(x_t),$$

from which it follows that $E[w_t \mid \mathcal{F}_t] = 0$. We also have

$$\begin{aligned} E[\|w_t\|^2 \mid \mathcal{F}_t] &= E[\|\nabla f_{k(t)}(x_t)\|^2 \mid \mathcal{F}_t] - \|E[\nabla f_{k(t)}(x_t) \mid \mathcal{F}_t]\|^2 \\ &\leq E[\|\nabla f_{k(t)}(x_t)\|^2 \mid \mathcal{F}_t], \end{aligned}$$

which yields

$$E[\|w_t\|^2 \mid \mathcal{F}_t] \leq \max_k \|\nabla f_k(x_t)\|^2.$$

Let us assume that there exist constants C and D such that

$$\|\nabla f_i(x)\| \leq C + D\|\nabla f(x)\|, \quad \forall i, x, \quad (6.2)$$

(cf. the assumption of Prop. 3). It follows that

$$E[\|w_t\|^2 \mid \mathcal{F}_t] \leq 2C^2 + 2D^2\|\nabla f(x_t)\|^2,$$

so that condition (5.3) is satisfied and the assertion of Prop. 4 holds.

REFERENCES

- [BMP90] Benveniste, A., Metivier, M., and Priouret, P., 1990. Adaptive Algorithms and Stochastic Approximations, Springer-Verlag, N. Y.
- [BeT89] Bertsekas, D. P., and Tsitsiklis, J. N., 1989. "Parallel and Distributed Computation: Numerical Methods," Prentice-Hall, Englewood Cliffs.
- [BeT96] Bertsekas, D. P., and Tsitsiklis, J. N., 1996. "Neuro-Dynamic Programming," Athena Scientific, Belmont, MA.
- [Ber95a] Bertsekas, D. P., 1995. "Nonlinear Programming," Athena Scientific, Belmont, MA.

- [Ber95b] Bertsekas, D. P., 1995. "A Hybrid Incremental Gradient Method for Least Squares Problems," Lab. for Info. and Decision Systems Report LIDS-P-2257, Massachusetts Institute of Technology, Cambridge, MA; to appear in SIAM J. on Optimization.
- [Bor95] Borkar, V. S., 1995. "Asynchronous Stochastic Approximations," to appear in SIAM J. on Control and Optimization.
- [Del96] Delyon, B., 1996. "General Results on the Convergence of Stochastic Algorithms," IEEE Transactions on Aut. Control, Vol. 41, pp. 1245-1255.
- [Gai94] Gaivoronski, A. A., 1994. "Convergence Analysis of Parallel Backpropagation Algorithm for Neural Networks," Optimization Methods and Software, Vol. 4, pp. 117-134.
- [Gri94] Grippo, L., 1994. "A Class of Unconstrained Minimization Methods for Neural Network Training," Optimization Methods and Software, Vol. 4, pp. 135-150.
- [KuC78] Kushner, H. J., and Clark, D. S., 1978. "Stochastic Approximation Methods for Constrained and Unconstrained Systems," Springer-Verlag, NY.
- [KuY97] Kushner, H. J., and Yin, G., 1996. "Stochastic Approximation Methods," Springer-Verlag, NY.
- [Lju77] Ljung, L., 1977. "Analysis of Recursive Stochastic Algorithms," IEEE Trans. on Automatic Control, Vol. 22, pp. 551-575.
- [LuT94] Luo, Z. Q., and Tseng, P., 1994. "Analysis of an Approximate Gradient Projection Method with Applications to the Backpropagation Algorithm," Optimization Methods and Software, Vol. 4, pp. 85-101.
- [Luo91] Luo, Z. Q., 1991. "On the Convergence of the LMS Algorithm with Adaptive Learning Rate for Linear Feedforward Networks," Neural Computation, Vol. 3, pp. 226-245.
- [MaS94] Mangasarian, O. L., and Solodov, M. V., 1994. "Serial and Parallel Backpropagation Convergence Via Nonmonotone Perturbed Minimization," Optimization Methods and Software, Vol. 4, 1994, pp. 103-116.
- [PoT73] Poljak, B. T., and Tsytkin, Y. Z., 1973. "Pseudogradient Adaptation and Training Algorithms," Automation and Remote Control, Vol. 12, pp. 83-94.
- [Pol87] Poljak, B. T., 1987. "Introduction to Optimization," Optimization Software Inc., N.Y.
- [TBA86] Tsitsiklis, J. N., Bertsekas, D. P., and Athans, M., 1986. "Distributed Asynchronous

References

Deterministic and Stochastic Gradient Optimization Algorithms," IEEE Trans. on Aut. Control, Vol. AC-31, pp. 803-812.