# Gradient Descent (proofs techniques)

Saturday, March 21, 2020     5:23 PM

There are many different types of results. Here, I try to unify some of them.

**Algorithm / notation**

Algo $\begin{cases} X_1 \text{ arbitrary} \\ \text{for } t=1,2,\cdots,T \\ \quad X_t = X_{t-1} - \eta V_t \end{cases}$

TBD
ex, $V_t = \nabla f(X_t)$
or
if $V_t \in \partial f(X_t)$

**Lemma 14.1 (Shalev-Shwartz, Ben-David)**

Let $\{V_t\}_{t=1}^T$ be arbitrary, $f$ need not be convex nor smooth then Algo produces a sequence satisfying

$$\sum_{t=1}^T \langle X_t - X^*, V_t \rangle \leq \frac{\|X^* - X_1\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|V_t\|^2$$

**Corollary**

If $\|V_t\| \leq \rho \;\; \forall t$ (eg: $f$ is $\rho$-Lipschitz) and $\|X^* - X_1\| \leq B$

choosing $\eta = \sqrt{\dfrac{B^2}{\rho^2 T}}$ gives $\dfrac{1}{T} \sum_{t=1}^T \langle X_t - X^*, V_t \rangle \leq \rho \dfrac{B}{\sqrt{T}}$

($X^*$ denotes any minimizer of $f$)

proof sketch of lemma (just the good parts)

$$\sum_{t=1}^T \langle X_t - X^*, V_t \rangle = \frac{1}{2\eta} \sum \left( -\|X_{t+1} - X^*\|^2 + \|X_t - X^*\|^2 + \|V_t\|^2 \right)$$

(via completing-the-square and algebra)

$$= \frac{1}{2\eta} \left( \|X_1 - X^*\|^2 - \|X_{T+1} - X^*\|^2 \right) + \frac{\eta}{2} \sum \|V_t\|^2$$

since sum telescoped     $\underset{\|\cdot\| \geq 0}{R}$

$$\leq \frac{1}{2\eta} \|X_1 - X^*\|^2 + \eta/2 \sum \|V_t\|^2 \quad \square$$

How to use this result?

Case: f is convex (and $\rho$-Lipschitz, so corollary applies)

Choose $V_t \in \partial f(X_t)$     $\swarrow = \min_X f(x) = f(x^*)$

then by convexity,   $f(X_t) - f^* \leq \langle X_t - x^*, V_t \rangle$

So

Corollary 1: If f is convex and $\rho$-Lipschitz, running subgradient descent gives

$$\frac{1}{T}\sum_{t=1}^{T}\left(f(X_t) - f^*\right) \leq \frac{B\rho}{\sqrt{T}} \quad \begin{array}{l} \text{if } \|X_1 - x^*\| \leq B \\ \text{and stepsize } \eta = \sqrt{\frac{B^2}{\rho^2 T}} \end{array}$$

How to apply this?

If we can easily evaluate $f(X_t)$, then

let $X_{best} = \arg\min_{x \in \{X_1, \ldots, X_T\}} f(x)$

and

Corollary 1a:   $f(X_{best}) - f^* \leq \frac{B\rho}{\sqrt{T}}$

However, sometimes it's not easy to evaluate $f(x)$

ex: $f(w) = L_D(w)$ (we can only sample from it... as may be the case in SGD)

In that case, define $\bar{X} = \frac{1}{T}\sum_{t=1}^{T} X_t$

then

$f(\bar{X}) \leq \frac{1}{T}\sum_t f(X_t)$ by Jensen's ineq., hence

Corollary 1b:   $f(\bar{X}) - f^* \leq \frac{B\rho}{\sqrt{T}}$

Case: f is smooth (ie. $\nabla f$ is $\beta$-Lipschitz)

Descent lemma:   $f(X_{t+1}) \leq f(X_t) + \langle \nabla f(X_t), X_{t+1} - X_t \rangle + \frac{\beta}{2}\|X_{t+1} - X_t\|^2$

(no convexity needed just smoothness)

and if we run gradient descent w/ stepsize $\eta = \frac{1}{\beta}$ then via some algebra,

$$(*) \quad f(x_{t+1}) \leq f(x_t) - \frac{1}{2\beta} \| \underbrace{\nabla f(x_t)}_{\text{"}V_k\text{" in Algo}} \|^2$$

## Case, part 1: f is smooth but not convex

**Thm** If $f$ is $\beta$-smooth, then gradient descent w/ $\eta = \frac{1}{\beta}$ gives $\min_{t \in [T]} \| \nabla f(x_t) \|^2 \leq \frac{2\beta}{T} \left( f(x_1) - f^* \right)$

(for nonconvex, we don't show convergence to a global or even local minimizer, just $\| \nabla f(x_t) \| \to 0$, ie., a stationary pt. where $\nabla f(x) = 0$ )

**proof** Sum $(*)$ from $t = 1, \dots T$

$$\frac{1}{2\beta} \sum_{t=1}^{T} \| \nabla f(x_t) \|^2 \leq \sum_{t=1}^{T} f(x_t) - f(x_{t+1})$$
$$\text{(telescopes)}$$
$$= f(x_1) - f(x_{T+1}) \leq f(x_1) - f^*$$

and

$$\min_{t \in [T]} \| \nabla f(x_t) \|^2 \leq \frac{1}{T} \sum_{t=1}^{T} \| \nabla f(x_t) \|^2 \quad (\min \leq \text{avg}) \quad \square$$

## Case, part 2: f is $\beta$-smooth and convex

As we already saw above, when $f$ is convex, combined w/ Lemma 14.1 (but don't use the corollary yet)

$$\sum_{t=1}^{T} f(x_t) - f^* \leq \frac{\| x_1 - x^* \|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^{T} \| \underbrace{\nabla f(x_t)}_{V_t} \|^2$$

and by smoothness, descent lemma gives

$$f(x_{t+1}) + \frac{1}{2\beta} \| \nabla f(x_t) \|^2 \leq f(x_t). \quad \text{Choose } \eta = \frac{1}{\beta}$$

then

$$\sum_{t=1}^{T} \left( f(x_{t+1}) + \frac{\eta}{2} \| \nabla f(x_t) \|^2 - f^* \right) \leq \sum_{t=1}^{T} \left( f(x_t) - f^* \right)$$

$$\leq \| x_1 - x^* \|^2 \cdot \eta \, \overline{I} \dots / \dots^2$$

So (using $\eta = 1/\beta$)                                    $\frac{1}{2\eta}$    $+\frac{1}{2}\sum_{t=1}^{T} \|\nabla f(x_t)\|$

Thm: $f(x_T) \leq \frac{1}{T}\sum_{t=1}^{T}\left(f(x_t) - f^*\right) \leq \frac{1}{T}\frac{\beta}{2}\|x_1 - x^*\|^2$

↑

since $X_T = X_{best}$ in this case (descent lemma $\Rightarrow$
$f(x_{t+1}) \leq f(x_t)$ )

Case, part 3: f is $\beta$ smooth and $\lambda$-strongly convex

First, define

"Polyak-Łojasiewicz Inequality" or just "PL"

$\forall x, \frac{1}{2}\|\nabla f(x)\|^2 \geq \lambda(f(x) - f^*)$

then f $\lambda$-strongly convex $\Rightarrow$ f is $\lambda$-PL

Then to analyze, start at descent lemma again

$f(x_{t+1}) - f(x_t) \leq \frac{-1}{2\beta}\|\nabla f(x_t)\|^2$ (by $\beta$-smoothness)

$\leq \frac{-\lambda}{\beta}\left(f(x_t) - f^*\right)$ (by PL)

So re-arrange

$f(x_{t+1}) - f^* \leq \left(1 - \frac{\lambda}{\beta}\right)\left(f(x_t) - f^*\right)$

$\leq \left(1 - \frac{\lambda}{\beta}\right)^t \left(f(x_0) - f^*\right)$

So
Thm   $f(x_{t+1}) - f^* \leq c^t\left(f(x_0) - f^*\right)$

$c = \left(1 - \frac{\lambda}{\beta}\right) < 1$  "linear convergence"

Discussion of convergence rates

Error $e_t$. How many more iterations needed to go
from $\varepsilon = 1$ accuracy to $\varepsilon = 0.01$ accuracy?

| Rate | Iter. | Examples |
|---|---|---|
| 1. $e_T \propto \frac{1}{\sqrt{T}}$ $(T = O(\varepsilon^{-2}))$ | 10,000 times more | subgradient or gradient descent (not smooth); SGD |
| 2. $e_T \propto \frac{1}{T}$ $(T = O(\varepsilon^{-1}))$ | 100 times more | gradient descent (smooth) |

-linear

3. $e_T \propto \frac{1}{T^2}$ $\left(T = O(\varepsilon^{-1/2})\right)$ 10 ~~more~~ times more    accelerated gradient descent

$c < 1$

4. $e_T \propto c^T$ $\left(T = \log(\varepsilon^{-1})\right)$    2·const more    gradient descent (smooth and strongly convex)

$e_{t+1} \leq c \cdot e_t$

5. $e_{t+1} \leq c_t e_t$ , $\begin{array}{l} c_t < 1 \\ c_t \to 0 \end{array}$

6. $e_{t+1} \leq c e_t^2$ $\left(T = \log(\log(\varepsilon^{-1}))\right)$    1 more    Newton's method near a solution

$c > 0$

## Geometric Picture



quadratic upper bound (if smooth)

minimizer of upper bound

what we want to bound

$\frac{1}{2\beta}\nabla f(x_{t-1})$

$\geq 0$

← Value unknown

$\geq 0$

← Value (sort of) known

linear lower bound (if convex)

$f$

$X_{t-1}$  $X_t$  $X^{*}$

true minimizer