



THEORY OF MACHINE LEARNING

LECTURE 13

GRADIENT DESCENT, THOUGHTS

RECAP: CONVEX OPTIMIZATION

f : convex.

- **Problem.** Given a convex function defined over a convex domain, find the minimizer (or min value).

$$\mathcal{D}.$$
$$\operatorname{argmin}_{x \in \mathcal{D}} f(x).$$

- Gradient descent - inspired by Taylor approximation

- Start with some feasible w_0
- For $t = 0, 1, \dots, T-1$, set $w_{t+1} = \underbrace{w_t}_{\text{circled}} - \eta \nabla f(w_t)$
- How do you set/"tune" the learning rate?
- Staying feasible

w_t

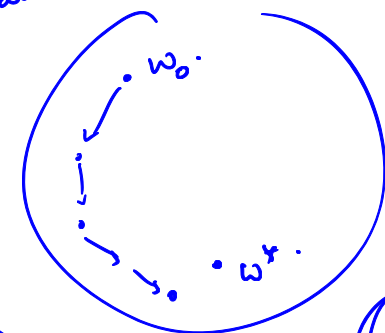
$$f(w_{t+1}) = f(w_t) - \eta \|\nabla f(w_t)\|^2.$$

keep projecting to \mathcal{D} .

BASIC THEOREM

- Assume f is L Lipschitz, domain is all of R^d , $|w_0 - w^*| \leq B$
- Theorem.** Consider running T steps of gradient descent with a fixed learning rate η . Then we have

$$\underbrace{f(\bar{w}) - f(w^*)}_{\substack{\bar{w} = \frac{w_1 + w_2 + \dots + w_T}{T}}} \leq \underbrace{\frac{1}{T} \sum_{t=1}^T \underbrace{f(w_t) - f(w^*)}_{\geq 0}}_{\substack{\text{norm of gradient} \\ \uparrow}} \leq \frac{B^2}{2\eta T} + \frac{L^2\eta}{2}$$



 RHS becomes $\left(\frac{BL}{\sqrt{T}} \right)^2$

- Same proof works if we had a constrained domain
- Proof works even if functions at different time steps were different!

$$\sum_{t=1}^T \underbrace{f_t(w_t) - f_t(w)}_{\geq 0} \leq \frac{B^2}{2\eta} + \frac{L^2\eta T}{2}$$

$$\sum_t f_t(w_t) - f_t(w)$$

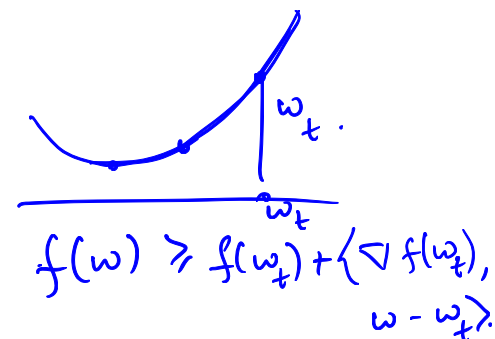
$\forall w$.

$$T = \frac{(BL)^2}{\epsilon^2}$$

ANALYSIS

- Use "basic inequality" about convex functions, for any t ,

$$f(w^*) \geq f(w_t) + \langle w^* - w_t, \nabla f(w_t) \rangle$$



- Use the potential function $\Phi_t = |w_t - w^*|^2$
- Note that $\Phi_t - \Phi_{t+1}$ (potential drop) is lower bounded by how far $f(w_t)$ is from $f(w^*)$

$$\Phi_t - \Phi_{t+1} \geq 2\eta (f(w_t) - f(w^*)) - \eta^2 |\nabla f(w_t)|^2$$

$\eta^2 L^2$

- Summing over t gives the bound

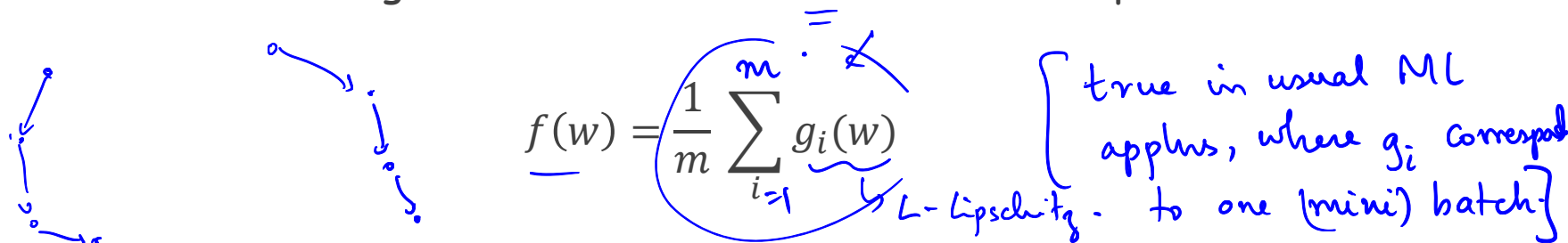
- Applications to online convex optimization, SGD

→ stochastic grad descent.

Cool thing about the analysis: can even have different f_t 's at diff times t .

STOCHASTIC GRADIENT DESCENT

- Consider the setting where the function f can be decomposed as



$$f(w) = \frac{1}{m} \sum_{i=1}^m g_i(w)$$

true in usual ML applies, where g_i corresponds to one (mini) batch.

L -Lipschitz

- In iteration t , pick index i_t uniformly at random and take a gradient step,

i.e., $w_{t+1} = w_t - \eta \nabla g_{i_t}(w_t)$ not the same as $w_t - \eta \nabla f(w_t)$

- Now w_t is a random variable, and we need to argue about $E[f(w_t)]$

- We still have that the expected potential drop at step t (conditioned on any trajectory so far) is $2\eta(f(w_t) - f(w^*)) - \eta^2 L^2$

- Earlier bound holds in expectation

$$E[\Phi_t - \Phi_{t+1} | w_t]$$

NOISY GRADIENT DESCENT

$$w \rightarrow w - \eta \nabla f(w)$$

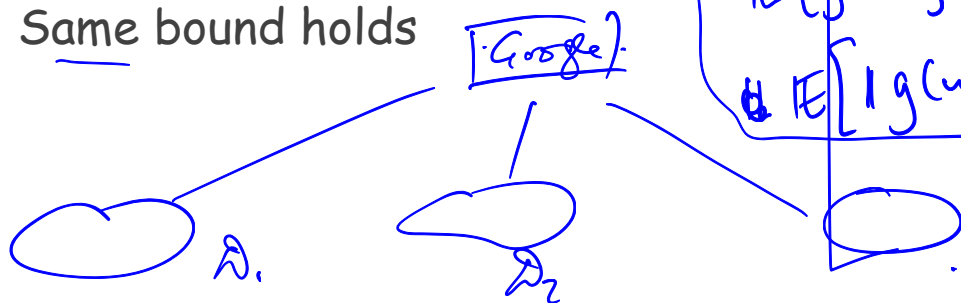
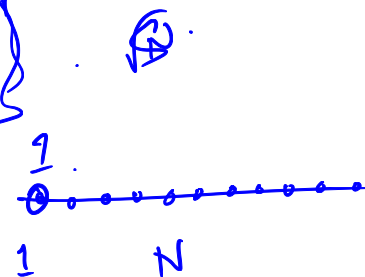
- Consider the setting where we perform gradient descent on the function f using a "noisy gradient oracle"

- Given a point w , suppose we get "noisy gradient"

Oracle that gives $g(w) \sim$ noisy gradient; with the props that

- Same bound holds

$$\begin{aligned} \mathbb{E}[g(w)] &= \nabla f(w) \\ \mathbb{E}[\|g(w)\|^2] &\leq L^2 \end{aligned}$$



f is said to be

~~quadratic~~ ax^2+bx+c
 $-\frac{b}{2a}$

ADDITIONAL STRUCTURE ON FUNCTIONS

$$f(y) \leq f(x) + \langle \nabla f(x), y-x \rangle + M \|y-x\|^2$$

(earlier: $|f(x) - f(y)| \leq L \|x-y\|$; $\|\nabla f(x) - \nabla f(y)\| \leq M \|x-y\|$)
 $\| \nabla f(x) \| \leq L$ \Leftrightarrow $\| \nabla^2 f(x) \| \leq M$

- Smoothness - function f is M smooth if gradient is M -Lipschitz

- Key observation: in this case, every iteration yields drop in function value!

$$w_{t+1} = w_t - \eta \nabla f(w_t)$$

$$f(w_{t+1}) \leq f(w_t) + \langle \nabla f(w_t), w_{t+1} - w_t \rangle + M \|w_{t+1} - w_t\|^2$$

$$= f(w_t) - \eta \|\nabla f(w_t)\|^2 + M \eta^2 \|\nabla f(w_t)\|^2$$

$$= f(w_t) - \frac{1}{2M} G^2 + \frac{1}{4M} G^2 = f(w_t) - \frac{\eta}{2} \|\nabla f(w_t)\|^2$$

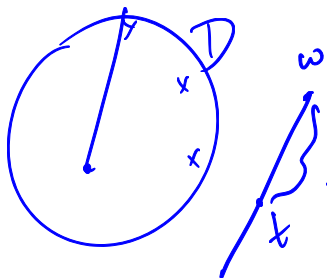
$$\eta \approx \frac{1}{2M}$$

- After T steps, $\sum_t \|\nabla f(w_t)\|^2$ is bounded by $4M (f(w_0) - f(w^*))$

- Convergence rate of $1/T$

$$f(w_{t+1}) \leq f(w_t) - \frac{\eta}{2} \|\nabla f(w_t)\|^2$$

$$\Leftrightarrow \|\nabla f(w_t)\|^2 \leq \frac{2}{\eta} (f(w_t) - f(w_{t+1}))$$



Earlier analysis:

$$\Phi_t = \|\omega_t - \omega^*\|^2$$

Saw that: ~~for~~ $\Phi_t - \Phi_{t+1} \geq 2\eta(f(\omega_t) - f(\omega^*)) - \eta^2 \|\nabla f(\omega_t)\|^2$

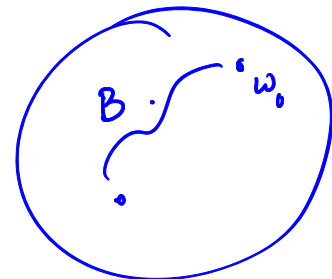
$$\Leftrightarrow f(\omega_t) - f(\omega^*) \leq \frac{\Phi_t - \Phi_{t+1}}{2\eta} + \frac{\eta}{2} \underbrace{\|\nabla f(\omega_t)\|^2}$$

$$\leq \frac{\Phi_t - \Phi_{t+1}}{2\eta} + (f(\omega_t) - f(\omega_{t+1}))$$

$$\frac{1}{T} \sum_{t=1}^T f(\omega_t) - f(\omega^*) \leq \frac{2M(\Phi_0 - \dots)}{2\eta = 1/2M} + \frac{f(\omega_0) - f(\omega_T)}{T} \leq C$$

error after T steps
 $\approx \frac{1}{T}$

$$\frac{2MB^2}{T} + \frac{C}{T}$$



NONCONVEX (SMOOTH) FUNCTIONS

→ If f is non-convex ~~and~~ but M -smooth, then analysis before (with $\eta = \frac{1}{2M}$) implies that $\|\nabla f(w_t)\|^2 \leq 4M(f(w_t) - f(w_{t+1}))$

$$\Rightarrow \frac{1}{T} \sum_{t=1}^T \|\nabla f(w_t)\|^2 \leq \frac{4M(f(w_0) - f(w^*))}{T}$$

→

$$\exists t \text{ s.t. } \|\nabla f(w_t)\|^2 \leq \frac{4MC}{T}.$$

↓
approximate singular pt.

ADDITIONAL STRUCTURE ON FUNCTIONS

- **Smoothness:** function is M smooth if *gradient* is M -Lipschitz
- **Strongly convex:** function is m -strongly convex if we have a “lower bound” via a parabola

IMPROVEMENTS, GENERALIZATIONS

- **Polyak-Lojasiewicz inequality:** suppose f satisfies:
 $|\nabla f(w)|^2 \geq c(f(w) - f(w^*))$ for all w
- “Global” condition, but can be satisfied for non-convex f
- Polyak’s “heavy ball” method (momentum)
- AdaGrad and related methods
- Second order (Newton) methods
- ...