# THEORY OF MACHINE LEARNING

# LECTURE 11

[Wed AM — office hours — 9:30- 11 AM.] .

(same zoom link as class.) .

# RECAP: LOSS MINIMIZATION

feature vector $\in \mathbb{R}^d$

$(\overset{\nearrow}{x_1}, y_1) \;,\; (x_2, y_2) \dots \dots \; (x_N, y_N)$

- ERM is hard, so we consider minimization of loss

label
$(\pm 1)$

$\min\limits_{h} \sum\limits_{i=1}^{N} \ell(h(x_i), y_i)$

- General problem

loss function $\ell$

parametrize hypotheses using $\underset{=}{\theta}$.

- Optimization can be hard in general, we study "easy" case of convex optimization

( hard unless $\ell$ is "nice" $\longrightarrow$ $\ell$ is convex in $\theta$, where $\theta$ are parameters that define $h$.)

- Min f(x) over D, where f is convex, domain D is convex

- Minimization is important (max can be hard)

# RECAP: CONVEX OPTIMIZATION

- **Problem.** Given a convex function defined over a convex domain, find the minimizer (or min value).
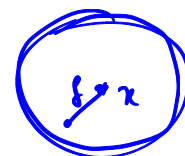
- $f(tx + (1-t)y) \leq t f(x) + (1-t)f(y)$ for all $t \in (0,1)$ and $x, y \in D$

- Local opt = global opt (just due to convexity)

- <u>Question:</u> how to find a "locally better" point? (assume f is continuous, differentiable)

- Gradient descent – inspired by Taylor approximation

$$f(x+\delta) \approx f(x) + \delta f'(x) \quad (\text{in } 1D)$$

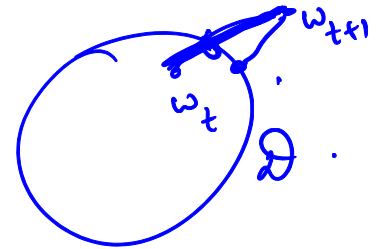$$\text{in a general,} \quad f(x+\delta) \approx f(x) + \langle \delta, \nabla f(x) \rangle.$$

$\downarrow$ vector

# GRADIENT DESCENT ALGORITHM

- Start with $w_0 = $ some feasible point; $\eta$ is some fixed param.
- For $t = 0, 1, \ldots, T-1$ : ($T$: # iterations)

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

$\rightarrow$ if $w_{t+1}$ is outside $\mathcal{D}$, set $w_{t+1}$ to be the projection.

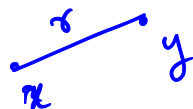- Generally applicable – even to non-convex functions

  (in which case you only find local opt)

- Choosing how much to move! (aka *learning rate*)

- Staying in the domain

  - how should we set $\eta$?

  - should it depend on $t$?

# VANILLA ANALYSIS

(As simple as possible)

$$|f(x) - f(y)| \le L \|x - y\|.$$

- Suppose f is L-Lipschitz, and domain $D = R^d$

  (no projection needed)   $\|\nabla f(x)\| \le L.$
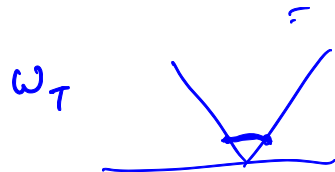
- Suppose OPT was distance B away from initial point

  $\longrightarrow \arg\min_x f(x)$

- **Theorem.**  Consider running T steps of gradient descent with a fixed

  learning rate $\eta$.  Then we have

  $$\|w_0 - w^*\| \le B$$

$w_T$

$$\frac{1}{T} \sum_{t=1}^{T} \left( f(w_t) - f(w^*) \right) \le \frac{B^2}{2\eta T} + \frac{L^2 \eta}{2}$$

$\quad e.g. \quad \sqrt{\frac{B^2 L^2}{2T}}$

$$w_t' = \frac{1}{T} \left( w_1 + \cdots + w_T \right)$$

(if $\eta$ is very small, say

$\eta = \frac{\varepsilon}{L^2}$ and T

"large"

($\varepsilon$ is some error of interest).

Proof uses "basic inequality" of convexity

$$f(w_{\frac{1}{T}}') \le \frac{1}{T} \sum_t f(w_t) \quad \text{(by convexity)}$$

$$f(w_t') - f(w^*) \le \dots$$

$$\forall y, \quad f(y) \geq f(x) + \langle y - x, \nabla f(x) \rangle.$$

$$[\iff \text{old definition of convexity.})$$

$$f(w^*) \geq f(w_t) + \boxed{\langle w^* - w_t, \nabla f(w_t) \rangle} . \longrightarrow \ast$$

(intuition) Tells us that keeping track of $\| w^* - w_t \|$ may be useful.

$$\Phi_t := \| w^* - w_t \|^2 \qquad (\text{hope: this reduces with time..})$$

$$\| a + b \|^2 = \| a \|^2 + \| b \|^2 + 2 \langle a, b \rangle.$$

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

$$= -2\eta \langle w^* - w_t, \nabla f(w_t) \rangle - \eta^2 \| \nabla f(w_t) \|^2$$

$$\| w^* - w_{t+1} \|^2 = \\ \leq \| w^* - w_{t+\frac{1}{2}} \|^2$$

(use $\ast$)

$\rightarrow$ f-divergences

($\ast$ HW)

$$\geq -2\eta\left[f(w^*) - f(w_t)\right] - \eta^2 \|\nabla f(w_t)\|^2$$

$$\underbrace{\phantom{-2\eta\left[f(w^*) - f(w_t)\right]}}_{= \, 2\eta\left[f(w_t) - f(w^*)\right]}$$

("local" analysis)

$$\therefore \quad 2\eta\left[f(w_t) - f(w^*)\right] \leq \frac{\Phi_t - \Phi_{t+1}}{2\eta} + \frac{\eta^2}{2}\|\nabla f(w_t)\|^2$$

(single t.)

$$\bigvee \qquad \sum_{t=0}^{T-1}\left(f(w_t) - f(w^*)\right) \leq \frac{\Phi_0 - \Phi_T}{2\eta} + \frac{\eta \cdot L^2 \cdot T}{2} \qquad \left| \begin{array}{l} \sum_t A_t - A_{t+1} \\[4pt] A_1 - A_2 + A_2 - A_3 + \\[4pt] \qquad A_3 - A_4 + \cdots \\[4pt] = A_0 - A_T \end{array}\right.$$

What is the best choice of $\eta$ for a given $T$?   Must set $\eta$ so

$$\leq \frac{B^2}{2\eta} + \frac{\eta \cdot L^2 T}{2}$$

that $\dfrac{B^2}{2\eta} = \dfrac{\eta \cdot L^2 T}{2}$ i.e, $\eta = \sqrt{\dfrac{B^2}{L^2 T}} := \dfrac{B}{L\sqrt{T}}$

(in theorem)

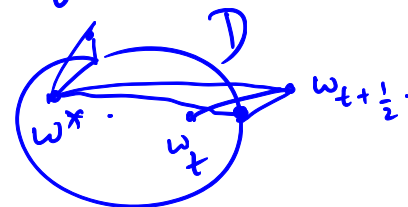RHS becomes $\dfrac{BL}{\sqrt{T}}$.

$\rightarrow$ Entire analysis goes through with projections!

$\rightarrow$ Formally when $D$ is bounded, $GD$ is the foll:

Say we have $\omega_t \in D$

$$\omega_{t+\frac{1}{2}} = \omega_t - \eta \cdot \nabla f(\omega_t)$$

$$\omega_{t+1} = \prod_D \left(\omega_{t+\frac{1}{2}}\right) = \underset{x \in \partial}{\text{argmin}} \left\| \omega_{t+\frac{1}{2}} \rightarrow x \right\|$$

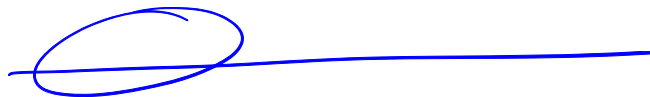Obsn: $\left\| \omega^* - \omega_{t+1} \right\|^2 \leq \left\| \omega^* - \omega_{t+\frac{1}{2}} \right\|^2 . \longrightarrow$ lets the original analysis go through!

(angle with separating plane is $90°$)

$\therefore$ angle with any pt on the other side of the plane is $\geqslant 90°$

# EXTENSIONS

- What if function is "smooth"? Get improved 'rate'

- What if function is "strongly convex"?

- What if functions at different steps are different? (!)

What if at time $t$, $f_t(x) \rightarrow$ convex, $\underline{L}$- Lipschitz, varies with time.

$\forall \underline{w^* \in D},$   $f_t(w^*) \geq f_t(w_t) + \langle w^* - w_t, \nabla f_t(w_t) \rangle$

$$\frac{1}{T} \sum_t (f_t(w_t) - f_t(w^*)) \leq \frac{B^2}{2\eta T} + \frac{L^2 \eta}{2} \cdot \quad \left[ \begin{array}{c} \text{Online Convex} \\ \text{opt.} \end{array} \right].$$