# THEORY OF MACHINE LEARNING

# LECTURE 10

CONVEX OPTIMIZATION, GRADIENT DESCENT

# OPTIMIZATION

SOLVING ERM EFFICIENTLY

# RECAP: ERM IS OFTEN HARD WITH FINITE CLASSES

- Finding best linear classifier (fewest mistakes) is NP hard even to approximate!

- Common remedy: loss functions

- Many candidate loss functions

$(x_1, y_1) \quad , \quad (x_2, y_2), \cdots (x_N, y_N)$

$\downarrow \qquad \searrow$

ERM: features  label

find $h$ s.t. $\sum\limits_{i=1}^{N} \mathbb{1}\left[ h(x_i) \neq y_i \right]$ is min

Loss min:

$\sum\limits_{i=1}^{N} loss\left( h(x_i), y_i \right)$

$\downarrow$

$\sum\limits_{i=1}^{n} \left( h(x_i) - y_i \right)^2$

if $h$ has some parametric form,

e.g. $h(x) = w^T x$ (in lin. class.)

Then $\min\limits_{w} \sum\limits_{i=1}^{n} \left( w^T x_i - y_i \right)^2$ is an easy problem.

Moving to loss min converts ERM into (more tractable) optimization prob. hopefully

# RECAP: CONVEX OPTIMIZATION

$$x \in \mathbb{R}^d$$
$$\ell(x) \in \mathbb{R}.$$

$$f : \text{convex}$$
$$\mathcal{D} \subseteq \mathbb{R}^d \text{ is a convex set}$$

$$\arg\min_{x \in \mathcal{D}} f(x)$$

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

$$f : \mathbb{R}^d \to \mathbb{R}.$$

- $f(tx + (1-t)y) \le t f(x) + (1-t)f(y)$ for all $t \in (0,1)$ and $x, y \in D$
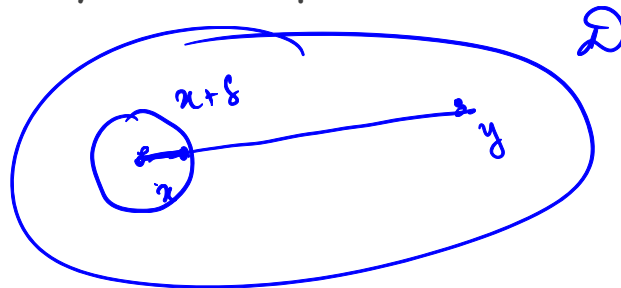
$$f\left(\frac{x+y}{2}\right) \le \frac{f(x) + f(y)}{2}.$$

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

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

$$\mathcal{D}$$

$$x + \delta$$
$$y$$

can you find some
$x + \delta$ s.t.
$$f(x + \delta) < f(x) ?$$

# TAYLOR APPROXIMATION

- Functions over $R^d$, gradients, Hessian

- First order approximation

(one dim): $f\left(x+\delta\right) = f(x) + \delta \cdot f'(x) + \frac{\delta^2}{2!} f''(x) + \cdots$

$$\nabla^2 f$$

First order approx: $\boxed{f(x+\delta) \simeq f(x) + \delta f'(x)}$

$\nabla^2 f \downarrow$ Hessian

$$f(x+\delta) = f(x) + \delta \cdot f'(x')$$

for some $x' \in (x, x+\delta)$ [Mean value theorem]

$$f''(x) \leq C$$

$$f^{(r)}(x) \leq c^r \cdot r!$$

Higher dimensions

$\forall \; x, \delta \in \mathbb{R}^d$.   ($\delta$ is "small")

$f(x) = \|x\|^2$

$(f : \mathbb{R}^d \to \mathbb{R})$
$\mathbb{R}$

$$f(x + \delta) \approx f(x) + \boxed{\langle \delta, \nabla f(x) \rangle}.$$

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}$$

$\nabla f(x)$

$x \quad \delta$

level set:
$$\{ x' : f(x) = f(x') \}$$

if we want to find $\delta$ s.t. $f(x + \delta) < f(x)$

$\langle y, y \rangle = \|y\|^2$.

idea: what if we set $\delta = -c \cdot \nabla f(x)$, for some $c > 0$?

$$\langle \delta, \underline{\nabla f(x)} \rangle = -c \|\nabla f(x)\|^2$$

This ensures that either $\underline{\nabla f^{(x)} = 0}$, or $f(x + \delta) < f(x)$.
(strictly)

(already at local opt $\longrightarrow$ global opt).

$x \longmapsto x - c \cdot \nabla f(x)$

# GRADIENT DESCENT ALGORITHM

- Generally applicable – even to non-convex functions
  (in which case you only find ~~local opt~~) some critical pt. $\left( \text{derivative} = 0 \right)$

→ Start with some pt. $w^{(0)} \in \mathcal{D}$; pick $\eta$ as some param.

→ For $T$ steps, do:    $t = 0, 1, \ldots, T-1$

$$\boxed{w^{(t+1)} = w^{(t)} - \textcircled{$\eta$} \nabla f(w^{(t)})}$$

$$\eta_t \cdot \nabla f(w^{(t)})$$

Choose $T$ s.t. $\| \nabla f(w^{(T)}) \|$ is "small enough".

# NATURAL ISSUES
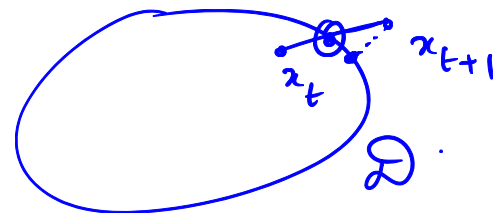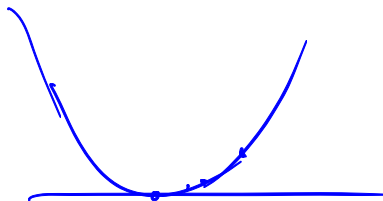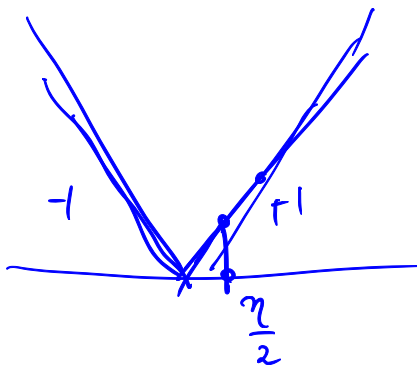
$$\eta_t = \frac{1}{t}.$$

$$\eta_t = \eta.$$

$$\boxed{\eta_t = \frac{1}{\sqrt{t}}}.$$

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

- Staying in the domain

  $\hookrightarrow$ project $x^{(t+1)}$ to $\mathcal{D}$.

$x_t$    $x_{t+1}$    $\mathcal{D}$.

projection can be tricky! (if $\mathcal{D}$ is not "simple")

$-1$    $+1$
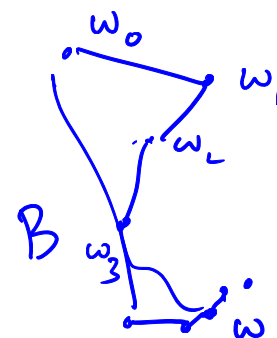
$\frac{\eta}{2}$

# GRADIENT DESCENT – VANILLA ANALYSIS

$$|f(x+\delta) - f(x)| \leq L\|\delta\|$$

L-Lipschitz: $|f(x) - f(y)| \leq L\|x-y\|_2$ $\approx \langle \delta, \nabla f(x)\rangle$

Convex,

$\implies \|\nabla f(x)\| \leq L.$

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

- Suppose OPT was distance B away from initial point

- **Theorem.** Consider running T steps of gradient descent with a fixed learning rate $\eta$. Then we have

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

w: true minimizer of f, $\bar{u}$, $w = \underset{x \in R^d}{\operatorname{argmin}} f(x)$.

$\frac{L^2 \cdot \eta}{2} \approx \eta$

Theorem tells us that if we perform "sufficiently many" steps of GD with fixed l.rate $\eta$, $\frac{1}{T} \sum_{t=1}^{T} (f(w_t) - f(w)) \leq O(\eta)$.
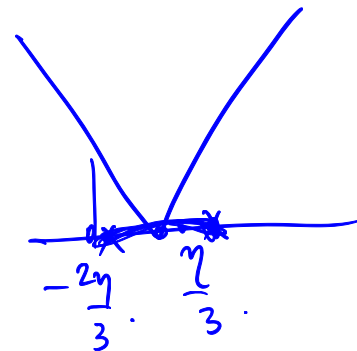
Can you find one pt $\bar{w}$ s.t. $f(\bar{w}) - f(w) \leq \boxed{O(\eta)}$?

"Last iterate convergence"

$\underset{t}{\text{argmin}} \ f(w_t)$

$\bar{w} = \frac{1}{T} \sum_{t} w_t$

$f(\bar{w}) \leq \frac{1}{t} \sum_{t} f(w_t)$

$-\frac{2\eta}{3}$      $\frac{\eta}{3}$
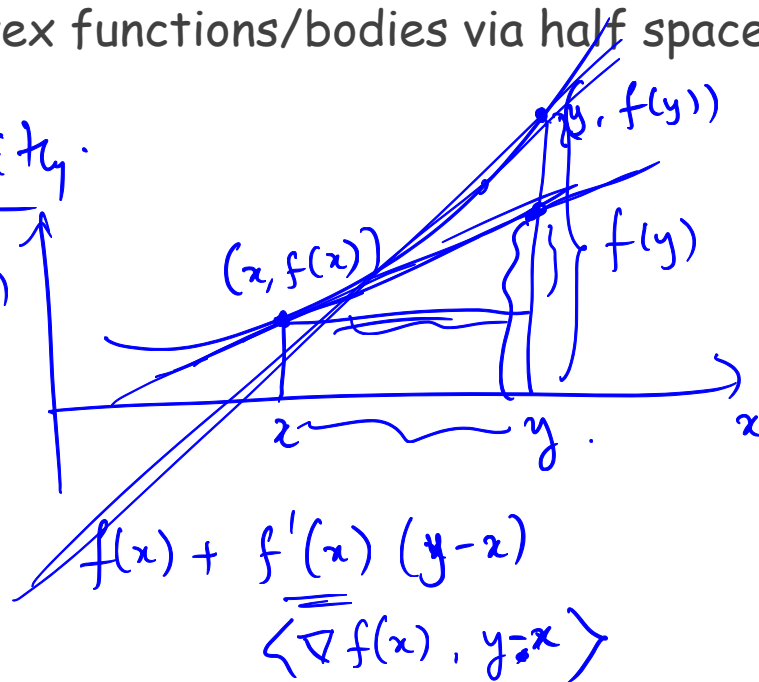
# ALTERNATE DEFINITION OF CONVEXITY

- Function lies "above" the tangent plane *at any point* !

- Related to the definition of convex functions/bodies via half spaces
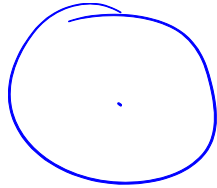
"Basic inequality" in convexity.

$$f(y) - f(x) \geq \langle \nabla f(x), y-x \rangle$$

$$f\left(\frac{x+y}{2}\right) \leq \frac{f(x) + f(y)}{2}.$$

(Homework 2).

$f(x)$

$(y, f(y))$

$(x, f(x))$

$f(y)$

$z$    $y$ .    $x$

$$f(x) + f'(x)(y-x)$$
$$\underbrace{\langle \nabla f(x), y-x \rangle}$$
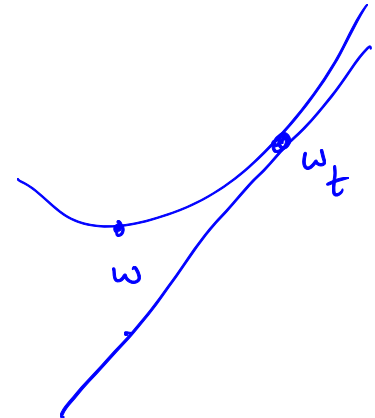
# GRADIENT DESCENT – VANILLA ANALYSIS

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

- Suppose OPT was distance B away from initial point

- **Theorem.** Consider running T steps of gradient descent with a fixed learning rate $\eta$. Then we have

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

$$f(w) - f(w_t) \geq \cdots.$$

$$f(w_t) - f(w) \leq \cdots$$

Basic ineq of convexity:

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

$$\Longleftrightarrow \boxed{f(w_t) - f(w)} \leq \boxed{\langle \nabla f(w_t), w_t - w \rangle}$$

How to connect to grad. descent?

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

(Algebraic magic)

$$\Phi_t = \| w_t - w \|^2$$

$$\| z - \eta \nabla f(w_t) \|^2 - \| z \|^2$$
$$\| z \|^2 + \eta^2 \| \nabla f(w_t) \|^2 - 2 \langle z, \eta \nabla f. \rangle$$

$$\boxed{\overline{\Phi}_{t+1} - \overline{\Phi}_t} = \| w_t - \eta \nabla f(w_t) - w \|^2 - \| \overset{z}{\overbrace{w_t - w}} \|^2$$

$$= -2\eta \boxed{\langle \nabla f(w_t), w_t - w \rangle} + \eta^2 \| \nabla f(w_t) \|^2$$

$$\Phi_{t+1} - \cancel{\Phi}_t + \cancel{\Phi}_t - \cancel{\Phi}_{t-1} + \ldots \Big/ + \cancel{\Phi}_1 - \overline{\Phi}_0$$

## ANALYSIS VIA POTENTIAL FUNCTIONS

→ drop in potential is "related" to ~~the~~ change in function value...

# DEALING WITH THE DOMAIN – PROJECTED GD