

Lecture 14

- * Time slots for presentation review have been posted on the project homepage.

RECAP

Boosting

Weak-learning: An algorithm A is a weak learner with edge/advantage γ for class c if: for any dist. D & any target $c \in C$, given access to $\text{Ex}(c, D)$, w.p. $(1-\delta)$ A produces a hypothesis with $\text{error}(h_i, c, D) \leq \frac{1}{2} - \gamma$.

If A runs in time $\text{poly}(d, \frac{1}{\epsilon})$ & $\gamma > \frac{1}{\text{poly}(d)}$
then C is efficiently weakly-PAC learnable

Thm If C is weakly-PAC learnable (efficiently),
then C is PAC-learnable (efficiently).

TODAY

Proof AdaBoost (Freund & Schapire)
(Nobel prize '03, very influential
in practice)

Training set: $\{(x_1, y_1), \dots, (x_n, y_n)\}$
 $x_i \in \mathcal{X}$
 $y_i \in \{-1, +1\}$

By realizability,

$$\mathcal{F} \subsetneq \mathcal{C} \quad \text{s.t.} \quad y_i = c(x_i) \quad \forall i.$$

We assume \mathcal{F} weak learning algorithm (WL) for \mathcal{C} .

AdaBoost (Freund & Schapire)

$$1. D_1(i) = \frac{1}{n} \quad \forall i \in [n]$$

$$2. \text{ for } t = 1, \dots, T$$

3. Use WL with dist D_t to get h_t

4. Let $\varepsilon_t = P_{x \sim D_t} [h_t(x) \neq y] \quad \text{if } \varepsilon_t \leq \frac{1}{2} - \gamma \text{ (w.p. 1-}\delta)$

$$5. \text{ Choose } d_t = \frac{1}{2} \log \left(\frac{1 - \varepsilon_t}{\varepsilon_t} \right)$$

$$6. D_{t+1}(i) = \frac{D_t(i)}{Z_{t+1}} + \begin{cases} e^{-d_t} & \text{if } h_t(x_i) = y_i \\ e^{d_t} & \text{if } h_t(x_i) \neq y_i \end{cases}$$

$$= \frac{D_t(i)}{Z_{t+1}} \cdot e^{(-d_t h_t(x_i)) y_i}$$

(where Z_{t+1} is normalizing constant)

7. Output $\text{sign}(H(x))$ where $H(x) = \sum_{t=1}^T d_t h_t(x)$

Note:

- 1) We can assume $\varepsilon_t \leq \frac{1}{2} - \gamma$ (union bound)
- 2) Can emulate $E_t(x, D_t)$, because D_t has finite support, just reweigh.

Aside: AdaBoost fits into the "Multiplicative weight updates" framework. General framework with many algorithmic applications.

Lemma For $T = \lceil \frac{1}{2\gamma^2} \log(2n) \rceil$, training error is 0.

Proof

$$\begin{aligned}\text{Training error} &= \mathbb{P}_{D_1} [\text{sign}(H(x)) \neq y] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{1} (\text{sign}(H(x_i)) \neq y_i)\end{aligned}$$

Note that

$$\mathbb{1} (\text{sign}(H(x)) \neq y) \leq e^{-y H(x)}$$

$$\text{Define } H_t(x) = \sum_{s=t}^T \alpha_s h_s(x)$$

$$H_t(x) = d_t h_t(x) + H_{t+1}(x)$$

$$H_1(x) = H(x)$$

$$\begin{aligned}
& \frac{1}{n} \sum_{i=1}^n \mathbb{1}(\text{sign}(h(x_i)) \neq y_i) \leq \sum_{i=1}^n D_t(i) e^{-y_i h_t(x_i)} \\
& = \underbrace{\sum_{i=1}^n D_t(i)}_{= Z_2} e^{-\lambda_2 y_i h_2(x_i) - y_i H_2(x_i)} \\
& = Z_2 \sum_{i=1}^n D_2(i) e^{-y_i H_2(x_i)} \\
& = Z_2 \underbrace{\sum_{i=1}^n D_2(i)}_{= Z_3} e^{-\lambda_2 y_i h_2(x_i) - y_i H_3(x_i)} \\
& = Z_2 Z_3 \sum_{i=1}^n D_3(i) e^{-y_i H_3(x_i)} \\
& \vdots \\
& = \prod_{t=2}^{T+1} Z_t
\end{aligned}$$

What is Z_{t+1} ?

$$\begin{aligned}
Z_{t+1} &= \sum_{i=1}^n D_t(i) e^{-\lambda_t y_i h_t(x_i)} \\
&= \sum_{i: y_i = h_t(x_i)} D_t(i) e^{-\lambda_t} + \sum_{i: y_i \neq h_t(x_i)} D_t(i) e^{\lambda_t} \\
&= (1 - \epsilon_t) e^{-\lambda_t} + \epsilon_t \cdot e^{\lambda_t}
\end{aligned}$$

$$\text{Recall } \lambda_t = \frac{1}{2} \log \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

$$\begin{aligned}\therefore Z_{t+1} &= (1 - \varepsilon_t) \cdot \sqrt{\frac{\varepsilon_t}{1 - \varepsilon_t}} + \varepsilon_t \sqrt{\frac{1 - \varepsilon_t}{\varepsilon_t}} \\ &= 2 \sqrt{\varepsilon_t(1 - \varepsilon_t)}\end{aligned}$$

Define $r_t = \frac{1}{2} - \varepsilon_t$ Note that $r_t \geq r$

$$\begin{aligned}Z_{t+1} &= 2 \sqrt{\left(\frac{1}{2} - r_t\right)\left(\frac{1}{2} + r_t\right)} \\ &= \sqrt{1 - 4r_t^2} \\ &\leq \left(e^{-4r_t^2}\right)^{1/2} \\ &= e^{-2r_t^2} \\ &\leq e^{-2r^2}\end{aligned}$$

$$\begin{aligned}\therefore \prod_{t=1}^{T+1} z_t &\leq e^{-2Tr^2} < \frac{1}{2n} \\ \text{if } T &\geq \frac{1}{2r^2} \log(2n).\end{aligned}$$

♦

What about test error?

Suppose that WL (weak learning algo.) always outputs a hypothesis from some class H whose VC-dim is d .

$$LC(\mathcal{H}, T) = \left\{ \text{sign} \left(\sum_{i=1}^T d_i h_i(x) \right); h_i \in \mathcal{H}, d_i \in \mathbb{R} \right\}$$

Exercise: VC-dim($LC(\mathcal{H}, T)$) $\leq c \cdot T \cdot d \cdot \log T$ for some constant c .

Hint: As in HW1, compute growth function & use Sauer's Lemma.

\therefore Due to VC Theorem,

$$\text{if } n \geq \frac{c}{\varepsilon} \left(\log \left(\frac{1}{\delta} \right) + T \cdot d \log T \log \left(\frac{1}{\varepsilon} \right) \right)$$

then the hypothesis produced by Ada Boost has error $\leq \varepsilon$, w.p. $1 - \delta$.

Putting in the bound of T ,

$$\text{if } n \geq \frac{c}{\varepsilon} \left(\log \left(\frac{1}{\delta} \right) + \frac{1}{2\gamma^2} \log(2n) d \log \left(\frac{1}{2\gamma^2} \log(2n) \right) \cdot \log \left(\frac{1}{\varepsilon} \right) \right)$$

we will get error $\leq \varepsilon$ w.p. $1 - \delta$.

This is satisfied for

$$n \geq \frac{c}{\varepsilon} \left(\frac{d \operatorname{polylog}(d, \frac{1}{r^2}, \frac{1}{\varepsilon})}{r^2} + \log(\frac{1}{\delta}) \right)$$

$$\left[n \geq c \log n \text{ is satisfied for } n \geq c \log \alpha \right]$$

Convex Optimization

Basic properties,

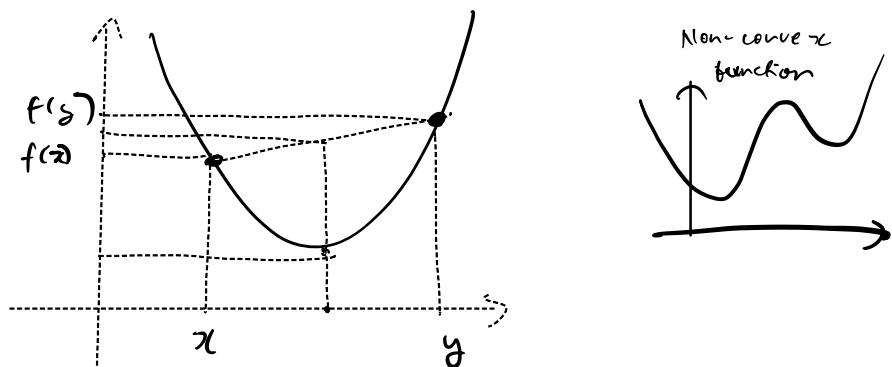
Def (convex set) A set $C \subseteq \mathbb{R}^d$ is convex if $x, y \in C \Rightarrow tx + (1-t)y \in C$ for all $0 \leq t \leq 1$.



Def (convex functions) A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if $\text{domain}(f) \subseteq \mathbb{R}^d$ is convex &

$$f(tx + (1-t)y) \leq t f(x) + (1-t) f(y)$$

$\forall t \in [0,1] \quad \& \quad x, y \in \text{domain}(f).$



Lemma If f is differentiable, then f is convex if & only if domain (f) is convex &

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

Proof Exercise, use definition of convexity.

Also,

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

Corollary If f is convex & differentiable then $\nabla f(x) = 0$ implies x is a global minimum of f (local minimum implies global minimum).

Def (Strong convexity) A function f is α -strongly convex if its domain is convex & if x, y

$$f(y) \geq f(x) + \langle y-x, \nabla f(x) \rangle + \frac{\alpha}{2} \|y-x\|^2.$$

E.g. $f(x) = \|x\|^2$

Lemma (Some properties of convex functions)

① If f is twice differentiable, then f is convex if & only if $\text{domain}(f)$ is convex

$$\& \nabla^2 f(x) \geq 0 \quad \forall x \in \text{domain}(f).$$

(Recall $A \geq 0 \Leftrightarrow x^\top A x \geq 0 \quad \forall x$).

② If $f_i(x)$ is convex $\forall i \in [n]$ then

$$g(x) = \sum_{i=1}^n w_i f_i(x) \text{ is convex, where } w_i \geq 0.$$

③ If $f_i(x)$ is convex $\forall i \in [n]$ then

$$g(x) = \max_{i \in [n]} f_i(x) \text{ is convex.}$$

Convex Learning Problems

An optimization problem

$$\min_{x \in A} f(x)$$

is called a convex optimization problem if

① $f(x)$ is convex ② A is convex.

Convex optimization problems can be solved efficiently.

Recall the ERM problem of finding the ERM wrt some hypothesis \mathcal{H} on some training set $S = (z_1, \dots, z_n)$ where $z_i = (x_i, y_i)$.

$$\text{ERM}_{\mathcal{H}}(S) = \underset{h \in \mathcal{H}}{\operatorname{arg\,min}} \sum_{i=1}^n l(h, z_i)$$

Let $\mathcal{H} \subseteq \mathbb{R}^d$ parameterized by $w \in \mathcal{W}$

$$\text{ERM}_{\mathcal{H}}(S) = \underset{w \in \mathcal{W}}{\operatorname{arg\,min}} \frac{1}{n} \sum_{i=1}^n l(w, z_i)$$

Lemma If l is a convex loss (in terms of w), & \mathcal{H} is convex, then $\text{ERM}_{\mathcal{H}}(S)$ is a convex optimization problem.

Proof Average of convex functions is convex..

Example (linear regression with squared loss)

Let $\mathcal{H} = \{x \mapsto \langle w, x \rangle, w \in \mathbb{R}^d\}$

$$l(h, (x, y)) = (h(x) - y)^2$$

$$l(w, (x, y)) = (\langle w, x \rangle - y)^2$$

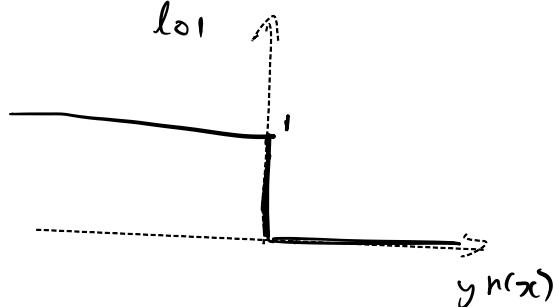
$$\mathcal{H} = \mathbb{R}^d \text{ (convex)}$$

Exercise $l(w, (x, y))$ is convex in terms of w .

$\therefore \text{ERM}_h(\gamma)$ is a convex problem.

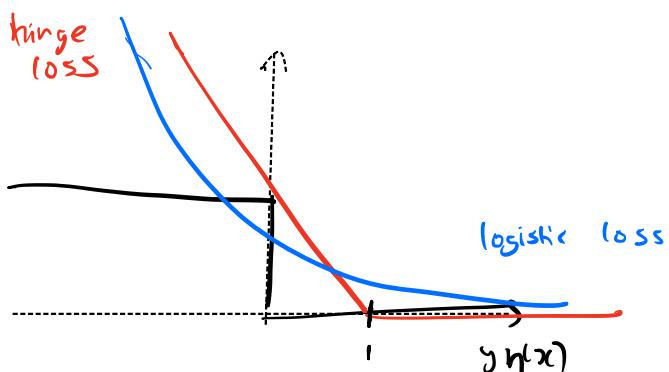
What about the 0/1 loss?

$$l_{01}(h, (x, y)) = \mathbb{1}(h(x) \neq y) = \mathbb{1}(y h(x) \leq 0)$$



Non-convex!

A common technique to handle a non-convex loss is to instead consider a convex surrogate.



hinge loss: $l(h, (x, y)) = \max(1 - y h(x), 0)$

logistic loss: $l(h, (x, y)) = \log(1 + e^{-y h(x)})$