

Course Note for CPSC532D STAT Learn Theory

Lec ERM, uniform convergence bound.
empirical risk minimization

$$\text{ERM}(S) \in \underset{h \in \mathcal{H}}{\text{argmin}} L_S(h) \\ = \underset{h \in \mathcal{H}}{\text{argmin}} \frac{1}{m} \sum_{i=1}^m \ell(h, z_i)$$

let $\hat{h}_S = \text{ERM}(S)$
let $h^* = \underset{h \in \mathcal{H}}{\text{arg inf}} L_0(h)$

how ERM generalize? uniform convergence ≤ 0 as $\hat{h}_S \in \text{argmin} L_S(h)$ Hoeffding inequality

$$L_0(\hat{h}_S) - L_0(h^*) = \underbrace{L_0(\hat{h}_S) - L_S(\hat{h}_S)}_{\text{estimation error}} + \underbrace{L_S(\hat{h}_S) - L_S(h^*)}_{\text{how } \hat{h}_S \text{ overfits}} + \underbrace{L_S(h^*) - L_0(h^*)}_{\text{how } \hat{h}_S \text{ fits } S} + \underbrace{L_0(h^*) - L_0(h^*)}_{\text{how } h^* \text{ underfits}}$$

$$L_0(\hat{h}_S) - L_{\text{Bayes}} = \underbrace{L_0(\hat{h}_S) - \inf_{h \in \mathcal{H}} L_0(h)}_{\text{estimation error}} + \underbrace{\inf_{h \in \mathcal{H}} L_0(h) - L_{\text{Bayes}}}_{\text{approximation error}}$$

excess error

increase $|\mathcal{H}|$

increase m

estimation error

\uparrow

$\rightarrow 0$ (ideally)

approximation error

\downarrow

-

concentration inequalities

bound the deviation of independent random variables from their expectations

$$t_i = \ell(h, z_i) - \mathbb{E}_{z \sim D} \ell(h, z) \quad L_0(h) = \mathbb{E}_{S \sim D^m} L_S(h)$$

large law number: $m \rightarrow \infty \quad \frac{1}{m} \sum_{i=1}^m t_i \rightarrow 0$

central limit theorem: $\frac{1}{m} \sum_{i=1}^m t_i \xrightarrow{d} \mathcal{N}(0, \sigma^2/m) \quad \sigma^2 = \text{Var}(t_i)$

Hoeffding: Let x_1, \dots, x_m be independent with:

- $\mathbb{E} x_i = \mu$
- $\Pr(a \leq x_i \leq b) = 1$

let $\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$, then

$$\textcircled{1} \Pr(\bar{x} \leq \mu + (b-a) \sqrt{\frac{\log 1/\delta}{2m}}) \geq 1 - \delta$$

let $Y_i = -x_i$. Then $\Pr(\bar{Y} \leq -\mu + (b-a) \sqrt{\frac{\log 1/\delta}{2m}}) \geq 1 - \delta$

$$\textcircled{2} \Pr(\bar{x} \geq \mu - (b-a) \sqrt{\frac{\log 1/\delta}{2m}}) \geq 1 - \delta$$

$$\textcircled{3} \Pr(|\bar{x} - \mu| \leq (b-a) \sqrt{\frac{\log 2/\delta}{2m}}) \geq 1 - \delta$$

$$\Pr(\bar{x} > \mu + (b-a) \sqrt{\frac{\log 1/\delta}{2m}}) < \delta$$

$$\Pr(\bar{x} < \mu - (b-a) \sqrt{\frac{\log 1/\delta}{2m}}) < \delta$$

$$\Pr(|\bar{x} - \mu| > (b-a) \sqrt{\frac{\log 2/\delta}{2m}}) < \delta$$

$\left. \begin{array}{l} \Pr(A|B) \\ \Pr(A) + \Pr(B) \end{array} \right\}$

use Hoeffding to bound $L_S(h^*) - L_D(h^*)$
 $\Rightarrow \Pr(L_S(h^*) - L_D(h^*) \geq (b-a) \sqrt{\frac{\log(1/\delta)}{2m}}) \leq \delta$

we can't use Hoeffding on $L_D(\hat{h}_S) - L_S(\hat{h}_S)$
 $l(\hat{h}_S, z_i)$ are not independent of each other
 change z_i will change \hat{h}_S

we can't use Hoeffding directly on $L_D(\hat{h}_S) - L_S(\hat{h}_S)$
 because $l(\hat{h}_S, z_i)$ is not independent of each other. change z_i will change \hat{h}_S , thus affect all other $l(\hat{h}_S, z_j)$

For $L_D(\hat{h}_S) - L_S(\hat{h}_S)$, we could use uniform convergence:
 bound every $h \in \mathcal{H}$, the gap will not be too big. then it'll be small for \hat{h}_S
 $\sup_{h \in \mathcal{H}} L_D(h) - L_S(h) \leq \epsilon \Rightarrow \forall h \in \mathcal{H} L_D(h) - L_S(h) \leq \epsilon$

If $|\mathcal{H}| < \infty$: union bound

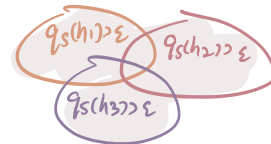
Let $g_S(h) = L_D(h) - L_S(h)$

$$\Pr(\exists h \in \mathcal{H}, g_S(h) > \epsilon) \leq \sum_{h \in \mathcal{H}} \Pr(g_S(h) > \epsilon)$$

now $l(h, z)$ is independent

$$\Pr(\exists h \in \mathcal{H}, L_D(h) - L_S(h) > \epsilon) \leq \sum_{h \in \mathcal{H}} \Pr(L_D(h) - L_S(h) > \epsilon) \leq |\mathcal{H}| \exp\left(\frac{-2m\epsilon^2}{(b-a)^2}\right)$$

$$\Rightarrow \Pr(L_S(h^*) - L_D(h^*) + \sup_{h \in \mathcal{H}} L_D(h) - L_S(h) > 2\epsilon) \leq (1 + |\mathcal{H}|) \exp\left(\frac{-2m\epsilon^2}{(b-a)^2}\right)$$



$$P(A \cup B) \leq P(A) + P(B)$$

Not tight enough because of $P(A \cap B) > 0$

so the overall bound would be:

Assume $\begin{cases} l(h, z) \in [a, b] \forall z, h, \\ \mathcal{H} \text{ finite} \\ \hat{h}_S \text{ is an ERM} \end{cases}$

$$L_D(\hat{h}_S) - \inf_{h \in \mathcal{H}} L_D(h) \leq L_D(\hat{h}_S) - L_S(\hat{h}_S) + L_S(h^*) - L_D(h^*) \leq \sup_{h \in \mathcal{H}} [L_D(h) - L_S(h)] + L_S(h^*) - L_D(h^*)$$

by showing: $\begin{cases} \bullet \forall h \in \mathcal{H}, \Pr(L_D(h) - L_S(h) > \epsilon) \leq \frac{\delta}{|\mathcal{H}|+1} \\ \bullet \Pr(L_S(h^*) - L_D(h^*) > \epsilon) \leq \frac{\delta}{|\mathcal{H}|+1} \end{cases}$

$$\therefore \Pr(L_D(\hat{h}_S) - \min_{h \in \mathcal{H}} L_D(h) \leq 2\epsilon) \geq 1 - \delta \quad \epsilon = (b-a) \sqrt{\frac{1}{2m} \log \frac{|\mathcal{H}|+1}{\delta}}$$

$$\Pr(L_D(\hat{h}_S) - \min_{h \in \mathcal{H}} L_D(h) \leq \sqrt{\frac{2}{m} \log \frac{|\mathcal{H}|+1}{\delta}}) \geq 1 - \delta$$

L3. MARKOV Inequality to prove Hoeffding

x cannot be too big with high prob $\mathbb{P}(X > 0) = 1$

- Markov's inequality: $\leftarrow X$ nonnegative-valued random variable
 $\Pr(X \geq t) \leq \frac{\mathbb{E}X}{t}$ for all $t > 0$ is weak because only assume non-neg

Pf $\begin{cases} x \geq 0 \\ x \geq t \text{ when } x \geq t \end{cases} \therefore \frac{x \geq t \mathbb{1}(x \geq t)}{\text{always holds}} \therefore \mathbb{E}X \geq t \mathbb{E} \mathbb{1}(x \geq t)$
 t is scalar $= t \Pr(X \geq t)$

\downarrow Let $\delta = \frac{\mathbb{E}X}{t}$ then $\Pr(X \geq \frac{\mathbb{E}X}{\delta}) \leq \delta \Rightarrow \Pr(X \leq \frac{\mathbb{E}X}{\delta}) \geq 1 - \delta$

- $\leftarrow Y = (X - \mathbb{E}X)^2$ just random variable
 Chebyshev's inequality: for any X , $\Pr(|X - \mathbb{E}X| \geq \epsilon) \leq \frac{\text{Var}X}{\epsilon^2}$

Pf $\Pr(|X - \mathbb{E}X| \geq \epsilon) = \Pr((X - \mathbb{E}X)^2 \geq \epsilon^2) \leq \frac{1}{\epsilon^2} \mathbb{E}(X - \mathbb{E}X)^2 = \frac{\text{Var}X}{\epsilon^2}$

with probability at least $1 - \delta$, $|X - \mathbb{E}X| \leq \sqrt{\frac{\text{Var}X}{\delta}} = \frac{\text{Var}X}{\sqrt{\delta}}$

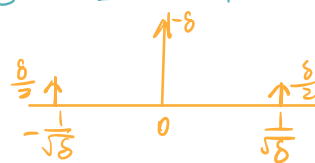
- Consider iid X_1, X_2, \dots, X_m , $\mathbb{E}X_i = \mu$, $\text{Var}(X_i) = \sigma^2$
 $\bar{X} = \frac{1}{m} \sum_{i=1}^m X_i$ $\mathbb{E}\bar{X} = \frac{1}{m} \sum_{i=1}^m \mathbb{E}X_i = \mu$ $\text{Var}(\bar{X}) = \frac{1}{m^2} \cdot m\sigma^2 = \frac{\sigma^2}{m}$
 \Rightarrow with $\Pr \geq 1 - \delta$, $|\bar{X} - \mu| \leq \frac{\sigma}{\sqrt{m\delta}}$

chebyshev's inequality may not be good enough!

$\frac{1}{\sqrt{m}}$ is fine, but the dependence on δ is the problem $\sqrt{\log \frac{1}{\delta}}$ better than $\frac{1}{\sqrt{\delta}}$

$\Pr(X = 0) = 1 - \delta$, $\Pr(X = \frac{1}{\sqrt{\delta}}) = \Pr(X = -\frac{1}{\sqrt{\delta}}) = \frac{1}{2} \delta$ $\mathbb{E}X = 0$ $\text{Var}X = 1$

$\Pr(|\bar{X}| \leq \frac{1}{\sqrt{m\delta}}) \geq 1 - \delta$



- Chernoff bounds: construct a non-negative random variable

$\lambda > 0$ arbitrary $Y = e^{\lambda(X - \mu)}$ $\Pr(Y \geq t) \leq \frac{\mathbb{E}Y}{t}$
 $\Pr(e^{\lambda(X - \mu)} \geq e^{\lambda\epsilon}) \leq e^{-\lambda\epsilon} \mathbb{E} e^{\lambda(X - \mu)}$

$= \Pr(\lambda(X - \mu) \geq \lambda\epsilon)$ centred moment-generating function $M_X(\lambda)$

$e^t = 1 + t + \frac{t^2}{2!} + \frac{t^3}{3!} + \dots$

$$M_X(\lambda) = \mathbb{E} e^{\lambda(X-\mu)} = 1 + \lambda \underbrace{\mathbb{E}[X-\mu]}_0 + \frac{\lambda^2}{2!} \mathbb{E}[(X-\mu)^2] + \frac{\lambda^3}{3!} \mathbb{E}[(X-\mu)^3] + \dots$$

k th derivative of $M_X(\lambda)$ at $\lambda=0$: $M_X^{(k)}(0) = \mathbb{E}[(X-\mu)^k]$

If $X \sim \mathcal{N}(\mu, \sigma^2)$, then $\mathbb{E} e^{\lambda(X-\mu)} = e^{\frac{1}{2}\lambda^2\sigma^2}$

for $X \sim \mathcal{N}(0, 1)$

$$\mathbb{E} e^{\lambda X} = \int \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} e^{\lambda x} dx = e^{-\frac{1}{2}\lambda^2} \int \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\lambda)^2} dx = e^{-\frac{1}{2}\lambda^2}$$

$\underbrace{\hspace{10em}}_{\mathcal{N}(\lambda, 1)}$

for $Y \sim \mathcal{N}(0, \sigma^2)$

$$\mathbb{E} e^{\lambda Y} = \mathbb{E} e^{\lambda \sigma X} = \mathbb{E} e^{(\lambda \sigma) X} = e^{-\frac{1}{2}\lambda^2\sigma^2}$$

$\underbrace{\hspace{2em}}_{Y \sim \mathcal{N}(0, \sigma^2)} \quad \underbrace{\hspace{2em}}_{\sigma X \sim \mathcal{N}(0, \sigma^2)} \quad \underbrace{\hspace{2em}}_{X \sim \mathcal{N}(0, 1)}$

$$\Pr(X \geq \mu + \varepsilon) \leq e^{-\lambda \varepsilon} e^{\frac{1}{2}\lambda^2\sigma^2}$$

we can optimize λ to get tighter bound

$$\sigma^2 \lambda = \varepsilon \rightarrow \lambda = \frac{\varepsilon}{\sigma^2} \rightarrow e^{\frac{1}{2}\frac{\varepsilon^2}{\sigma^2} - \frac{\varepsilon^2}{\sigma^2}} = e^{-\frac{\varepsilon^2}{2\sigma^2}}$$

$$\Pr(X \geq \mu + \varepsilon) \leq e^{-\frac{\varepsilon^2}{2\sigma^2}} \quad e^{-\frac{\varepsilon^2}{2\sigma^2}} = \delta$$

For $X \sim \mathcal{N}(\mu, \sigma^2)$:

with probability at least $1 - \delta$ $X \leq \mu + \sigma \sqrt{2 \log \frac{1}{\delta}}$

If X is st. $\mathbb{E} e^{\lambda(X-\mathbb{E}X)} \leq e^{\frac{1}{2}\lambda^2\sigma^2}$, then $\Pr(X \leq \mathbb{E}X + \sigma \sqrt{2 \log \frac{1}{\delta}}) \geq 1 - \delta$

σ SG(σ) sub-gaussian with σ (no heavier tails than Gaussian)

σ does not imply the variance of the sub-gaussian X may have relation but not equal!

• sub-Gaussian - $X \sim \text{SG}(\sigma)$

random variable X with mean $\mu = \mathbb{E}[X]$ is sub-Gaussian with σ

(i) if $\mathbb{E}[e^{\lambda(X-\mu)}]$ exists and for all $\lambda \in \mathbb{R}$, $\mathbb{E}[e^{\lambda(X-\mu)}] \leq e^{\frac{1}{2}\lambda^2\sigma^2}$

if $\sigma_1 < \sigma_2$ then $X \sim \text{SG}(\sigma_1) \Rightarrow X \sim \text{SG}(\sigma_2)$ $\text{SG}(\sigma_1) \subseteq \text{SG}(\sigma_2)$

(ii) If $X \in \text{SG}(\sigma)$, $aX \in \text{SG}(|a|\sigma)$ for $\forall a \in \mathbb{R}$

$$M_X(\lambda) \leq e^{\frac{1}{2}\lambda^2\sigma^2} \Rightarrow \mathbb{E} e^{\lambda(aX - a\mathbb{E}X)} = \mathbb{E} e^{a\lambda(X - \mathbb{E}X)} \leq e^{\frac{1}{2}(a\lambda)^2\sigma^2} = e^{\frac{1}{2}\lambda^2(|a|\sigma)^2}$$

(iii) If $X_1 \in \text{SG}(\sigma_1)$, $X_2 \in \text{SG}(\sigma_2)$ are independent $X_1 + X_2 \in \text{SG}(\sqrt{\sigma_1^2 + \sigma_2^2})$

$$\mathbb{E} e^{\lambda(X_1 + X_2 - \mathbb{E}X_1 - \mathbb{E}X_2)} = \mathbb{E} e^{\lambda(X_1 - \mathbb{E}X_1)} \mathbb{E} e^{\lambda(X_2 - \mathbb{E}X_2)} \leq e^{\frac{1}{2}\lambda^2\sigma_1^2} \cdot e^{\frac{1}{2}\lambda^2\sigma_2^2}$$

we can use the above properties to construct a set of X

Hoeffding's Lemma :

a real-valued random variable bounded in $[a, b]$ is $SG(\frac{b-a}{2})$

i.e. $\Pr(a \leq X \leq b) = 1, X \in SG(\frac{b-a}{2})$ see pf in notes 3.

- Hoeffding: If X_1, X_2, \dots, X_m iid with $\mathbb{E}X_i = \mu, \Pr(a \leq X_i \leq b) = 1$
 then $\Pr(\frac{1}{m} \sum_{i=1}^m X_i > \mu + (b-a) \sqrt{\frac{1}{2m} \log \frac{1}{\delta}}) \leq \delta$
 $\frac{1}{m} \sum_{i=1}^m X_i \in SG(\frac{1}{m} \sqrt{m (\frac{b-a}{2})^2}) = SG(\frac{b-a}{2\sqrt{m}})$
 with prob at least $1-\delta$ $\frac{1}{m} \sum_{i=1}^m X_i \leq \mu + \frac{b-a}{2\sqrt{m}} \sqrt{2 \log \frac{1}{\delta}} = \mu + (b-a) \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$

Lec 4. PAC Learning, infinite \mathcal{H}

learning algorithm from sample to \exists

Defn: an algorithm A is agnostically PAC learns \mathcal{H} with loss l

if there is a function $m(\epsilon, \delta)$ sample complexity function

st. for any D , for any $\epsilon, \delta \in (0, 1)$

if $S \sim P^m$, with $m \geq m(\epsilon, \delta)$

then $\Pr(L_0(A(S)) \leq \inf_{h \in \mathcal{H}} L_0(h) + \epsilon) \geq 1 - \delta$

solve for m will give sample complexity

efficient = polynomial run time

Def \mathcal{H} is agnostically PAC learnable if $\exists A$ that agnostically PAC learn \mathcal{H}

agnostically PAC is the worst case, \mathcal{H} has nothing to do with data distribution.

an $m(\epsilon, \delta)$ should work for any D

PAC learnable does not show ^{how} learn quickly in terms of m .

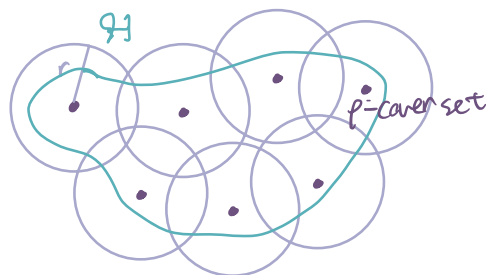
- Agnostically PAC learning: work for any distribution ERM infinite \mathcal{H}
- PAC learning is weaker: work only for realizable distribution

Def (PAC-learn) A PAC-learns \mathcal{H} if
 $\exists m = (0, 1)^2 \rightarrow \mathbb{N}$ s.t.
 for any (ϵ, δ) , for any realizable \mathcal{D} , if $m \geq m(\epsilon, \delta)$
 $\Pr_{S \sim \mathcal{D}^m} (L_0(A(S)) \leq \epsilon) \geq 1 - \delta$

Def (realizable). for a non-negative loss $\ell(h, z) \geq 0$.
 \mathcal{D} is realizable by \mathcal{H} if there exists an $h^* \in \mathcal{H}$ s.t. $L_0(h^*) = 0$

We already know finite case, what about infinite \mathcal{H} ?

idea of covering number



still: uniform convergence

$$\begin{aligned} & \sup_{h \in \mathcal{H}} L_0(h) - L_S(h) \\ &= \sup_{h \in \mathcal{H}} L_0(h) - L_0(h_j) + L_0(h_j) - L_S(h_j) + L_S(h_j) - L_S(h) \\ &\leq \sup_{h \in \mathcal{H}} \underbrace{[L_0(h) - L_0(h_j)]}_{\text{Lipschitz}} + \max_{\substack{h \in \mathcal{T} \\ j}} \underbrace{[L_0(h_j) - L_S(h_j)]}_{\text{union bound on size of } \mathcal{T}} + \sup_{h \in \mathcal{H}} [L_S(h_j) - L_S(h)] \end{aligned}$$

as p decrease, the set to cover points become tighter

\mathcal{T} is a p -cover for \mathcal{H} of size $N(\mathcal{H}, p)$

specific case: Logistic regression

$$\mathcal{Z} = \mathcal{X} \times \mathcal{Y} \quad \mathcal{X} = \mathbb{R}^d \quad \mathcal{Y} = \{-1, 1\}$$

$$\mathcal{H} = \{x \mapsto w \cdot x = W \in \mathbb{R}^d, \|w\| \leq B\} \quad \text{equivalent to } L_2 \text{ reg}$$

$$\ell_{\log}(h, (x, y)) = \ell_y(h(x)) = \log(1 + \exp(-y h(x)))$$

use absolute value to make analysis easier, but results looser

$$|L_0(\underline{w}) - L_0(\underline{v})| = \left| \mathbb{E}_{(x, y) \sim \mathcal{D}} \ell_{\log}(h_w, (x, y)) - \ell_{\log}(h_v, (x, y)) \right|$$

weight w weight v

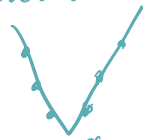
$$\leq \mathbb{E} |l_y(hw(x)) - l_y(hv(x))| \quad l_y \text{ is } \underline{G}\text{-Lipschitz}$$

$$\leq \underline{G} \mathbb{E} |hw(x) - hv(x)|$$

$G=1$ for l_y

input close \rightarrow output close
smoothness

G -Lipschitz: l_y is G -Lips if

$$|l_y(\hat{y}_1) - l_y(\hat{y}_2)| \leq G |\hat{y}_1 - \hat{y}_2| \quad \forall \hat{y}_1, \hat{y}_2$$


Euclidean distance (can be other distance if particular specified)

\Rightarrow if f' exist everywhere, and $\sup |f'(x)| \leq G$, f is G -Lips

$$|f(x) - f(x')| = \left| \int_x^{x'} f'(t) dt \right| \leq \int_x^{x'} |f'(t)| dt \leq \int_x^{x'} G dt = G |x' - x|$$

if f is differentiable the Lips constant is the largest |derivative| (upper bound of |derivative|)

for the LR: $|hw(x) - hv(x)| = |w \cdot x - v \cdot x| = |\langle w-v, x \rangle| \leq \|w-v\| \cdot \|x\|$

$$l_y(\hat{y}) = \log(1 + \exp(-y\hat{y})) \quad l_y'(\hat{y}) = \frac{-y \exp(-y\hat{y})}{1 + \exp(-y\hat{y})} = \frac{-y}{\exp(y\hat{y}) + 1}$$

$$|l_y'(\hat{y})| \leq 1 \quad l_y \text{ is } 1\text{-Lipschitz}$$

$$ax_1 + bx_2 \leq \sqrt{a^2 + b^2} \sqrt{x_1^2 + x_2^2}$$

$$a^2 x_1^2 + b^2 x_2^2 \leq (ax_1 + bx_2)^2$$

$$|hw(x) - hv(x)| = |w \cdot x - v \cdot x| = |\langle w-v, x \rangle| \leq \|w-v\| \cdot \|x\|$$

Then $\mathbb{E} |L_D(hw) - L_D(hv)| \leq \underbrace{\mathbb{E} \|x\|}_{\text{constant of distribution}} \cdot \|w-v\|$ distance from h to h_j at most p

$\leq Cp$

Assume $\|x\| \leq C$ on \mathcal{D} .

① $|L_S(hw) - L_S(hv)| \leq \frac{1}{m} \sum_{i=1}^m \|x_i\| \|w-v\| \leq Cp$

② $|L_D(h_j) - L_S(h_j)| \leq \underbrace{(b-a)}_{\text{range}} \sqrt{\frac{1}{2m} \log \frac{N(B, \rho)}{\delta}}$ l_y is not bounded, but we assume $\|w\| \leq B$ and $\|x\| \leq C$ so we can use Hoeffding's

$$\log(1 + \exp(BC)) \leq BC + 1 \quad (\text{just simpler})$$

So $\sup_{h \in \mathcal{H}} L_D(h) - L_S(h) \leq 2Cp + (BC+1) \sqrt{\frac{1}{2m} \log \frac{N(B, \rho)}{\delta}}$

what is the covering number?

$$N(B, \rho) \leq \left(\frac{3B}{\rho}\right)^d \quad (\text{proof in notes } \Psi)$$

$$\sup_{h \in \mathcal{H}} L_D(h) - L_S(h) \leq 2Cp + (BC+1) \sqrt{\frac{1}{2m} \left(\log \frac{1}{\delta} + d \log \frac{3B}{\rho} \right)}$$

how to choose ρ to minimize the right term?

Let $\rho = \frac{B}{\sqrt{m}}$
(roughly optimal)

$$\leq \frac{2CB}{\sqrt{m}} + (BC+1) \sqrt{\frac{\log \frac{1}{\delta} + d \log(9m)}{2m}}$$

Big Op notation means

dominate by $O_p\left(\sqrt{\frac{\log m}{m}}\right)$ some random variable is stochastically bounded. (Hoeffding)

$$\Pr(\sup_{h \in \mathcal{H}} L_D(h) - L_S(h) > \frac{2BC}{\sqrt{m}} + (BC+1) \sqrt{\frac{\log 1/\delta + \frac{d}{2} \log 9m}{2m}}) < \delta$$

- the sample complexity depends on the input distribution $\mathbb{E}\|X\| \leq c$
 X not allowed in agnostically PAC learning
 so this bound doesn't show Linear Regression is agnostically PAC learnable ^{ERM+}
- we can never achieve 0 loss on any \mathcal{D} , so it's not realizable

L3 Rademacher complexity distribution-specific complexity

- covering number approach: depend on dimension d .
 need bounded norm. B . scale sensitive
- Rademacher complexity: do not depend on dimension d .

here we are going to bound on average instead of high probability.

$$\mathbb{E}_{S \sim \mathcal{D}^m} \sup_{h \in \mathcal{H}} L_D(h) - L_S(h) = \mathbb{E}_{S \sim \mathcal{D}^m} \left(\sup_{h \in \mathcal{H}} \underbrace{\mathbb{E}_{S' \sim \mathcal{D}^m} [L_{S'}(h) - L_S(h)]}_{L_D(h)} \right)$$

how much can I overfit to sample S

$$\textcircled{1} \sup_x \mathbb{E} f_y(x) \leq \mathbb{E}_x \sup_y f_y(x)$$

Pf: for any y , we have $f_y(x) \leq \sup_{y'} f_{y'}(x)$

$$\text{then } \mathbb{E}_x f_y(x) \leq \mathbb{E}_x \sup_{y'} f_{y'}(x)$$

$$\text{and then } \sup_x \mathbb{E} f_y(x) \leq (\sup_x) \mathbb{E} \sup_{y'} f_{y'}(x)$$

$$\leq \mathbb{E}_{S, S' \sim \mathcal{D}^m} \sup_{h \in \mathcal{H}} [L_{S'}(h) - L_S(h)]$$

$$= \mathbb{E}_{S, S'} \sup_h \frac{1}{m} \sum_{i=1}^m (l(h, z_i) - l(h, z'_i))$$

$$\textcircled{2} \mathbb{E}_{S, S', \vec{\sigma}} \mathbb{E}_{U, U'} \left[\sup_h \frac{1}{m} \sum_{i=1}^m \sigma_i (l(h, U_i) - l(h, U'_i)) \mid S, S', \vec{\sigma} \right]$$

not doing anything (deterministic U, U' ($S, S', \vec{\sigma}$)) condition

(switch order of expec $P(S, S') \cdot P(\vec{\sigma}) \cdot P(U, U' \mid S, S', \vec{\sigma})$)

$$= \mathbb{E}_{U, U', \vec{\sigma}} \mathbb{E}_{S, S'} \left[\sup_h \frac{1}{m} \sum_{i=1}^m \sigma_i (l(h, U_i) - l(h, U'_i)) \mid U, U', \vec{\sigma} \right]$$

deterministic we can drop S

$\textcircled{3}$ let $\sigma_i \in \{-1, 1\}$ for $i \in [m]$,

$$(u_i, u'_i) = \begin{cases} (z_i, z'_i) & \text{if } \sigma_i = 1 \\ (z'_i, z_i) & \text{if } \sigma_i = -1 \end{cases}$$

for any $\vec{\sigma} = (\sigma_1, \dots, \sigma_m)$

$$\text{then } l(h, z_i) - l(h, z'_i) = \sigma_i (l(h, z_i) - l(h, z'_i))$$

$$\sigma_i \sim \text{Unif}; P(\sigma_i = 1) = P(\sigma_i = -1) = \frac{1}{2}$$

$U, U' | S, S', \vec{\sigma}$ is determined
 but $U, U' \sim D^m$, $U, U', \vec{\sigma}$ are independent
 $S, S' | U, U', \vec{\sigma}$ is determined

$$\begin{aligned}
 &= \mathbb{E}_{U, U', \vec{\sigma}} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i (\ell(h, u_i) - \ell(h, u_i')) \right] \\
 &\quad \text{rename } U \rightarrow S \\
 &= \mathbb{E}_{S, S' \sim D^m} \mathbb{E}_{\vec{\sigma} \sim \text{Rad}^m} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i (\ell(h, z_i) - \ell(h, z_i')) \right] \\
 &\quad \text{sup } [\sigma_i \ell(h, z_i) + (-\sigma_i) \ell(h, z_i')]
 \end{aligned}$$

$$\begin{aligned}
 &\leq \mathbb{E}_{S, \vec{\sigma}} \mathbb{E}_h \sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i \ell(h, z_i) + \mathbb{E}_{S', \vec{\sigma}} \mathbb{E}_h \sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m (-\sigma_i) \ell(h, z_i') \\
 &\quad \leftarrow \text{same} \rightarrow \\
 &= 2 \mathbb{E}_{S, \vec{\sigma}} \sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i \ell(h, z_i) \\
 &\stackrel{②}{=} 2 \mathbb{E}_S \text{Rad}(\ell \circ \mathcal{H} | S)
 \end{aligned}$$

$$\mathbb{E}_{S \sim D^m} \sup_{h \in \mathcal{H}} L_D(h) - L_S(h) \leq 2 \mathbb{E}_S \text{Rad}(\ell \circ \mathcal{H} | S)$$

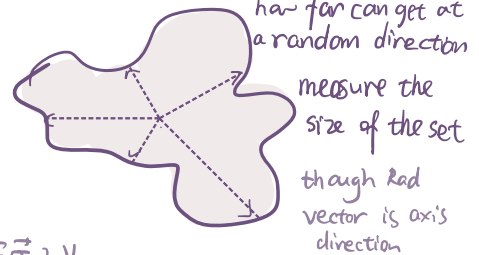
$$\mathbb{E}_{S \sim D^m} \sup_{h \in \mathcal{H}} L_S(h) - L_D(h) \leq 2 \mathbb{E}_S \text{Rad}(\ell \circ \mathcal{H} | S)$$

③ def $\ell \circ \mathcal{H} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$

$\mathcal{F}|_S = \{(f(z_1), f(z_2), \dots, f(z_m)) : f \in \mathcal{F}\} \subseteq \mathbb{R}^m$ per-sample risk

$$\text{Rad}(V) = \mathbb{E}_{\vec{\sigma} \sim \text{Rad}^m} \sup_{v \in V} \frac{1}{m} \sum_{i=1}^m \sigma_i v_i$$

the biggest product one can get given a Rademacher vector



(i) suppose one-element hypothesis :

$$\text{Rad}(\{v\}) = \mathbb{E}_{\vec{\sigma}} \sup_{v \in \{v\}} \frac{\vec{\sigma} \cdot v}{m} = \mathbb{E}_{\vec{\sigma}} \frac{\vec{\sigma} \cdot v}{m} = \frac{(\mathbb{E} \vec{\sigma}) \cdot v}{m} = 0$$

(ii) suppose \mathcal{H} is bigger :

$$\text{Rad}(\{-1, 1\}^m) = \mathbb{E}_{\vec{\sigma}} \sup_{v \in \{-1, 1\}^m} \frac{\vec{\sigma} \cdot v}{m} = 1$$

$$\text{Rad}(\{-1, 1\}^m) = 1$$

$$\text{Rad}(\{v : \|v\| \leq \sqrt{m}\}) = 1$$

\mathcal{H} bigger, but Rademacher is equal

because Rad vector only looks at ± 1 Gaussian complexity would look at any direction

property :

$$\begin{aligned}
 \text{(i) } \text{Rad}(cV) &= \text{Rad}(\{cv : v \in V\}) = \frac{1}{m} \mathbb{E}_{\vec{\sigma}} \sup_{v \in V} \sigma \cdot (cv) = \frac{|c|}{m} \mathbb{E}_{\vec{\sigma}} \sup_{v \in V} \text{sign}(c) \vec{\sigma} \cdot v \\
 c \in \mathbb{R} &= |c| \text{Rad}(V)
 \end{aligned}$$

$$\begin{aligned} \text{(ii) Rad}(V+W) &= \text{Rad}(\{v+w : v \in V, w \in W\}) = \frac{1}{m} \mathbb{E} \sup_{v,w} \sigma(v+w) = \frac{1}{m} \mathbb{E} \sqrt{\sup_v \sigma(v)^2 + \sup_w \sigma(w)^2} \\ &= \text{Rad}(V) + \text{Rad}(W) \end{aligned}$$

if w is constant vector $\Rightarrow \text{Rad}(w) = 0$

How to compute $\text{Rad}(\ell \circ \mathcal{H}|_S)$ practically?

• Talagrand's contraction lemma: l -Lipschitz function

Let $\phi: \mathbb{R}^m \rightarrow \mathbb{R}^m$ be given by $\phi(t) = (\phi_1(t), \dots, \phi_m(t))$, where each ϕ_i is l -Lip then $\text{Rad}(\phi \circ V) = \text{Rad}(\{\phi(v) : v \in V\}) \leq l \text{Rad}(V)$

Then for typical supervised learning losses

$$\begin{aligned} (\ell \circ \mathcal{H})|_S &= \{(\ell(h, z_1), \dots, \ell(h, z_m)) : h \in \mathcal{H}\} \\ &= \{(l_{y_1}(h(x_1)), \dots, l_{y_m}(h(x_m))) : h \in \mathcal{H}\} \\ &= (l_{S_y} \circ \mathcal{H})|_{S_x} \end{aligned}$$

might depend on S_y

if l_{y_i} (the loss func of a prediction for y_i) are all l -Lipschitz. then Talagrand's lemma gives:

$$\text{Rad}(\ell \circ \mathcal{H}|_S) \leq l \text{Rad}(\mathcal{H}|_{S_x})$$

$S_x = (x_1, \dots, x_m)$

only holds for real-valued hypothesis (H don't need to be function) (can be param of Gaussian...)

for linear regression, l_x is l -Lipschitz $\mathcal{H}_B = \{x \mapsto w \cdot x : \|w\| \leq B\}$

$$\text{Rad}(l_y \circ \mathcal{H}|_S) = \text{Rad}(\{l_{y_1}(h(x_1)), l_{y_2}(h(x_2)), \dots, l_{y_m}(h(x_m))\})$$

$$\leq \text{Rad}(\mathcal{H}|_S)$$

$$\text{Rad}(\mathcal{H}_B) = \mathbb{E} \sup_{\sigma} \frac{1}{m} \sum_{i=1}^m \sigma_i \langle w, x_i \rangle$$

$$= \mathbb{E} \sup_{\sigma} \frac{1}{m} \langle w, \sum_{i=1}^m \sigma_i x_i \rangle$$

$$\leq \frac{1}{m} \mathbb{E} B \|\sum_{i=1}^m \sigma_i x_i\|$$

$$\leq \frac{1}{m} B \sqrt{\mathbb{E} \|\sum_{i=1}^m \sigma_i x_i\|^2}$$

to get rid of σ

$$\begin{aligned} \mathbb{E} \|\sum_{i=1}^m \sigma_i x_i\|^2 &= \langle \sum_{i=1}^m \sigma_i x_i, \sum_{j=1}^m \sigma_j x_j \rangle \\ &= \sum_{i,j} \sigma_i \sigma_j \langle x_i, x_j \rangle \end{aligned}$$

$$= \frac{B}{m} \sqrt{\sum_i \mathbb{E} \sigma_i^2 \|x_i\|^2 + \sum_{i \neq j} \mathbb{E} \sigma_i \sigma_j \langle x_i, x_j \rangle}$$

$$= \frac{B}{m} \sqrt{\sum_i \|x_i\|^2}$$

independent $= \mathbb{E} \sigma_i \mathbb{E} \sigma_j = 0$

$$\mathbb{E}_S \text{Rad}(\mathcal{H}|_S) \leq \frac{B}{\sqrt{m}} \mathbb{E}_{S_x} \sqrt{\sum_i \|x_i\|^2} \leq \frac{B}{\sqrt{m}} \sqrt{\mathbb{E}_{S_x} \|x\|^2}$$

If $\Pr(\|x\| \leq C) = 1$, $\text{Rad}(f_B(x)) \leq \frac{BC}{\sigma_m}$

It is fine in some case to only look the average

Proof for Talagrand's Lemma

① use Lemma 2 to prove Talagrand's contraction Lemma

Lemma 2 if $\phi: \mathbb{R} \rightarrow \mathbb{R}$ is l -Lipschitz, $\text{Rad}(\{\phi(v_1), v_2, \dots, v_m\}) \leq \text{Rad}(V)$

Pf: $\text{Rad}(\{\frac{1}{l}\phi_1(v_1), v_2, \dots, v_m\}) \leq \text{Rad}(V)$ ← l -Lipschitz

$\text{Rad}(\{v_2, \frac{1}{l}\phi_1(v_1), \dots, v_m\})$ Rotate doesn't change Rad

\vdots

$\text{Rad}(\{v_2, v_3, \dots, \frac{1}{l}\phi_1(v_1)\})$

the same for ϕ_1, ϕ_2, \dots

$\text{Rad}(\{\frac{1}{l}\phi_1(v_1), \frac{1}{l}\phi_2(v_2), \dots, v_m\}) \leq \text{Rad}(V)$

$\text{Rad}(\{\frac{1}{l}\phi_2(v_2), \frac{1}{l}\phi_1(v_1), \dots, v_m\}) \leq \text{Rad}(V)$

$\therefore \text{Rad}(\frac{1}{l}\phi \circ V) \leq \text{Rad}(V)$

$\therefore \text{Rad}(\phi \circ V) \leq l \text{Rad}(V)$

② prove this Lemma 2

$m \cdot \text{Rad}(\{\phi(v_1), v_2, \dots, v_m\})$

$= \mathbb{E} \sup_{\vec{\sigma} \sim \text{Rad}^m} \sigma_1 \phi(v_1) + \underbrace{\sigma_2 v_2 + \dots + \sigma_m v_m}_{\vec{\sigma}_2 \cdot v_2}$

$= \frac{1}{2} \mathbb{E} \sup_{\vec{\sigma}_2 \sim \text{Rad}^m} [\phi(v_1) + \vec{\sigma}_2 \cdot v_2] + \frac{1}{2} \mathbb{E} \sup_{\vec{\sigma}_2 \sim \text{Rad}^m} [-\phi(v_1) + \vec{\sigma}_2 \cdot v_1]$ two independent parts

$= \frac{1}{2} \mathbb{E} \sup_{\vec{\sigma}_2} [\phi(v_1) - \phi(v_1') + \vec{\sigma}_2 \cdot (v_2 + v_2')]$ switch v, v' won't change

$\leq \frac{1}{2} \mathbb{E} \sup_{\vec{\sigma}_2} [l|v_1 - v_1'| + \vec{\sigma}_2 \cdot (v_2 + v_2')]$ 1-Lipschitz

$= \frac{1}{2} \mathbb{E} \sup_{\vec{\sigma}_2} [l|v_1 - v_1'| + \vec{\sigma}_2 \cdot (v_2 + v_2')]$

$= \frac{1}{2} \mathbb{E} \sup_{\vec{\sigma}_2} [l|v_1 - v_1'| + \vec{\sigma}_2 \cdot v_2] + \sup_{v \in V} [v_1 + \vec{\sigma}_2 \cdot v_2] = \mathbb{E} \sup_{\vec{\sigma}} \vec{\sigma} \cdot v = m \text{Rad}(V)$

Lee McDiarmid's inequality

use Randermacher to get a high probability bound
 \Rightarrow McDiarmid's inequality (concentration inequality)

McDiarmid's inequality: Let X_1, \dots, X_m be independent, let $f(x_1, \dots, x_m)$ be a real-valued function satisfying **banded differences**
 $\forall i \in [m] \sup_{x_1, \dots, x_m, x_i'} |f(x_1, \dots, x_m) - f(x_1, \dots, x_{i-1}, x_i', x_{i+1}, \dots, x_m)| \leq C_i$
change any variable at any one index doesn't change result too much

Then with probability at least $1 - \delta$,

$$\begin{cases} f(x_1, \dots, x_m) \leq \mathbb{E}f(x_1, \dots, x_m) + \sqrt{\frac{1}{2} (\sum_{i=1}^m C_i^2) \log \frac{1}{\delta}} \\ f(x_1, \dots, x_m) \geq \mathbb{E}f(x_1, \dots, x_m) - \sqrt{\frac{1}{2} (\sum_{i=1}^m C_i^2) \log \frac{1}{\delta}} \end{cases}$$

(special case: $f(x) = \frac{1}{m} \sum_i x_i$, each $x_i \in [a, b]$
 then $\forall i \in [m] \frac{1}{m} \sup_{x_1, \dots, x_m, x_i'} |x_1 + \dots + x_m - (x_1 + \dots + x_{i-1} + x_i' + x_{i+1} + \dots + x_m)| = \frac{1}{m} \sup_{x_i, x_i'} |x_i - x_i'| = \frac{1}{m} (b-a)$

so $C_i = \frac{b-a}{m}$, $\frac{1}{m} \sum_i x_i - \mathbb{E} x_i \leq \sqrt{\frac{m}{2} (\frac{b-a}{m})^2 \log \frac{1}{\delta}} = (b-a) \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$

(McDiarmid's inequality implies Hoeffding)

in case of loss $\ell(h, z) \in [a, b]$ for all h, z , then with prob $\geq 1 - \delta$

$$\sup_{h \in \mathcal{H}} L_0(h) - L_S(h) \leq \mathbb{E} \sup_{h \in \mathcal{H}} L_0(h) - L_S(h) + (b-a) \sqrt{\frac{1}{2m} \log \frac{1}{\delta}} \quad (1)$$

if \hat{h}_S is an ERM, with prob $\geq 1 - \delta$ that

$$L_0(\hat{h}_S) - L_0(h^*) \leq \mathbb{E} \sup_{h \in \mathcal{H}} [L_0(h) - L_S(h)] + (b-a) \sqrt{\frac{2}{m} \log \frac{2}{\delta}} \quad (2)$$

Pf. $L_0(\hat{h}_S) - L_0(h^*) \leq \underbrace{L_0(\hat{h}_S) - L_S(\hat{h}_S)}_{(1)} + \underbrace{L_S(h^*) - L_0(h^*)}_{(2)}$

prove (1) = Let $S^{(i)} = (z_1, \dots, z_{i-1}, z_i', z_{i+1}, \dots, z_m)$, we have

$$L_0(h) - L_S(h) = L_0(h) - L_{S^{(i)}}(h) + L_{S^{(i)}}(h) - L_S(h)$$

$$\sup_h L_0(h) - L_S(h) \leq \sup_h L_0(h) - L_{S^{(i)}}(h) + \sup_h L_{S^{(i)}}(h) - L_S(h)$$

$$|\sup_h [L_0(h) - L_S(h)] - \sup_h [L_0(h) - L_{S^{(i)}}(h)]| \leq |\sup_h [L_S^{(i)}(h) - L_S(h)]| \quad \underline{s, s^{(i)} \text{ changable}}$$

$$|\sup_h [L_0(h) - L_S(h)] - \sup_h [L_0(h) - L_{S^{(i)}}(h)]| \leq \sup_h |L_{S^{(i)}}(h) - L_S(h)| = \frac{b-a}{m}$$

apply McDiarmid's inequality \rightarrow

with prob at least of $1-\delta$:

$$\sup_h L_D(h) - L_S(h) \leq \underbrace{\mathbb{E} \sup_h L_D(h) - L_S(h)}_{\mathbb{E} \text{Rad}(L \circ H)} + (b-a) \sqrt{\frac{1}{2m} \log \frac{1}{\delta}}$$

$\neq \sup_h \mathbb{E}(L_D(h) - L_S(h))$ \mathbb{E} and max does not commute!

Pf for McDiarmid **very general commonly used.**

$X_{1:k-1} = (X_1, \dots, X_{k-1})$. Fix some $k \in [m]$ freeze some arbitrary values $X_{1:k-1}$

$\mathbb{E}_{X_{k+1:m}} f(X_{1:k-1}, X_k, X_{k+1:m})$ is random depend on X_k

note that $\sup_{X_k} f(X_{1:m}) - \inf_{X_k} f(X_{1:m}) \leq C_k$ by assumption

$$\Rightarrow \mathbb{E}_{X_{k+1:m}} \sup_{X_k} f(X_{1:k-1}, X_k, X_{k+1:m}) + \sup_{X_k} \mathbb{E}_{X_{k+1:m}} f(X_{1:k-1}, X_k, X_{k+1:m}) \leq C_k \quad \sup \mathbb{E} \leq \mathbb{E} \sup$$

$$\sup_{X_k} \mathbb{E}_{X_{k+1:m}} f(X_{1:k-1}, X_k, X_{k+1:m}) + \sup_{X_k} \mathbb{E}_{X_{k+1:m}} (-f(X_{1:k-1}, X_k, X_{k+1:m})) \leq C_k$$

$$\sup_{X_k} \mathbb{E}_{X_{k+1:m}} f(X_{1:k-1}, X_k, X_{k+1:m}) - \inf_{X_k} \mathbb{E}_{X_{k+1:m}} f(X_{1:k-1}, X_k, X_{k+1:m}) \leq C_k$$

by Hoeffding's Lemma $\mathbb{E}_{X_{k+1:m}}$ is SG($C_k/2$) $\mathbb{E}[e^{\lambda(X-\mu)}] \leq e^{\frac{1}{2}\lambda^2\sigma^2}$

$$\text{then } \mathbb{E}_{X_k} \exp(\lambda \mathbb{E}_{X_{k+1:m}} f(X_{1:k-1}, X_k, X_{k+1:m})) \leq \exp(\lambda \mathbb{E}_{X_k} \mathbb{E}_{X_{k+1:m}} f(X_{1:k-1}, X_k, X_{k+1:m}) + \frac{1}{2}\lambda^2(\frac{C_k}{2})^2)$$

holds for any $X_{1:k-1}$, then average no $X_{1:k-1}$

$$\mathbb{E}_{X_{1:k}} \exp(\lambda \mathbb{E}_{X_{k+1:m}} f(X_{1:m})) \leq \mathbb{E}_{X_{1:k-1}} \exp(\lambda \mathbb{E}_{X_{k+1:m}} f(X_{1:m}) + \frac{1}{8}\lambda^2 C_k^2)$$

take log

$$\log \mathbb{E}_{X_{1:k}} \exp(\lambda \mathbb{E}_{X_{k+1:m}} f(X_{1:m})) \leq \log \mathbb{E}_{X_{1:k-1}} \exp(\lambda \mathbb{E}_{X_{k+1:m}} f(X_{1:m})) + \frac{1}{8}\lambda^2 C_k^2$$

sum over $k=1, \dots, m$ $a_k \leq a_1 + \sum_{k=1}^m \frac{1}{8}\lambda^2 C_k^2$ a_{k-1}

$$\log \mathbb{E}_{X_{1:m}} \exp(\lambda f(X_{1:m})) \leq \log \exp(\lambda \mathbb{E}_{X_{1:m}} f(X_{1:m})) + \sum_{k=1}^m \frac{1}{8}\lambda^2 C_k^2$$

$$\mathbb{E}_{X_{1:m}} \exp(\lambda f(X_{1:m})) \leq \exp(\lambda \mathbb{E}_{X_{1:m}} f(X_{1:m})) \cdot \exp(\sum_{k=1}^m \frac{1}{8}\lambda^2 C_k^2)$$

$$\mathbb{E}_{X_{1:m}} \exp(\lambda f(X_{1:m}) - \lambda \mathbb{E}_{X_{1:m}} f(X_{1:m})) \leq \exp(\frac{1}{2}\lambda^2 \sum_{k=1}^m (\frac{C_k}{2})^2) \quad f(X_{1:m}) \in \text{SG}(\frac{1}{2}\sqrt{\sum_{k=1}^m C_k^2})$$

with chernoff bound for sub-gaussians

with prob at least $1-\delta$ $f(X_{1:m}) \leq \mathbb{E} f(X_{1:m}) + \sqrt{\sum_{i=1}^m C_i^2} \cdot \sqrt{2 \log \frac{1}{\delta}}$

Rad is scale-sensitive $\hat{=}$ depend on the choice of function

Lec 6 Growth Function and VC Dimension

How to bound Rad for binary classifier?

For binary $\mathcal{H}|_{S_x} = \{h(x_1), h(x_2), \dots, h(x_m)\} = h \in \mathcal{H} \subseteq \{0, 1\}^m$

At most 2^m possible vectors even if \mathcal{H} is infinite

• $L_{0,1}(h, (x, y)) = \mathbb{1}(h(x) \neq y)$

$\mathcal{H}|_{S_x}$ finite

• If $|V| < \infty, \|v\| \leq B$ for all $v \in V$

Then $\text{Rad}(V) \leq \frac{B}{m} \sqrt{2 \log |V|}$

special case: if $V = \mathcal{H}_{\pm 1}|_{S_x}, \|v\| = \sqrt{1 + \dots + 1} = \sqrt{m}$
 $\text{Rad}(\mathcal{H}_{\pm 1}|_{S_x}) \leq \sqrt{\frac{2}{m} \log |\mathcal{H}_{\pm 1}|_{S_x}|}$

Pf $\text{Rad}(V) = \mathbb{E}_{\sigma} \sup_{v \in V} \sum_{i=1}^m \frac{\sigma_i v_i}{m}$

$\sigma_i \in SG(1)$

$\frac{\sum \sigma_i v_i}{m} \in SG(\frac{\|v\|}{m}) \Rightarrow \sum_i \frac{\sigma_i v_i}{m} \in SG(\sqrt{\sum_i \frac{v_i^2}{m^2}})$ independent
 $= SG(\frac{1}{m} \|v\|) \leq SG(\frac{B}{m})$

• $T_i \in SG(\sigma), \mathbb{E} T_i = 0, T_i$ can be dependent

then $\mathbb{E} \max_{i \in [m]} T_i \leq \sigma \sqrt{2 \log m}$

Pf. see in A2 & C2.4 Jensen's Inequality: $\exp(\mathbb{E} Y) \leq \mathbb{E} \exp(Y)$

when taking $\sup_{v \in V} \sum_{i=1}^m \frac{\sigma_i v_i}{m}$ can be dependent, so $\mathbb{E} \sup_{v \in V} \sum_{i=1}^m \frac{\sigma_i v_i}{m} \leq \frac{B}{m} \sqrt{2 \log |V|}$

$\text{Rad}(\mathcal{H}|_{S_x})$ depends on particular distribution $|H|_{S_x}$

we can use $|H|_{S_x}| < |H|$ but it's very loose!

use growth function to drop dependence on particular S_x

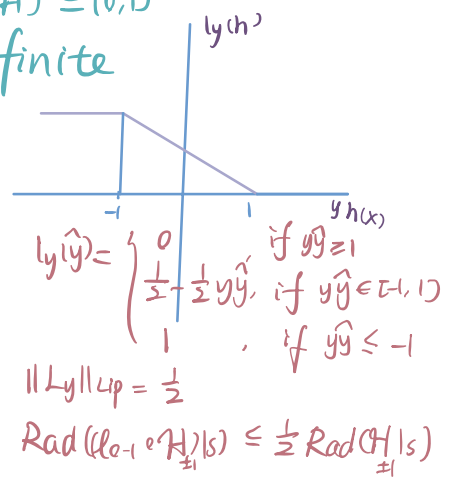
• The growth function of \mathcal{H} is

$\Pi_{\mathcal{H}}(m) = \sup_{x_1, \dots, x_m \in X} |\mathcal{H}|_{(x_1, \dots, x_m)}|, \quad |\mathcal{H}|_{S_x}| \leq \Pi_{\mathcal{H}}(m)$

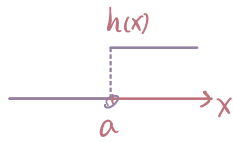
• (shatter) = if \mathcal{H} shatter a set S , means \mathcal{H} can assign any possible label to the set i.e. $|\mathcal{H}|_{S_x}| = 2^m$

if the theory could explain any outcome then \mathcal{H} is too general

VC is defined on the size of set you can shatter.

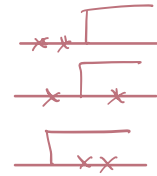


• (VC dimension) $VCdim(\mathcal{H}) = \max \{ m \geq 0 : T_{\mathcal{H}}(m) = 2^m \}$ worst case there is A set that can be shattered
 example = threshold $h_a(x) = \mathbb{1}(x \geq a)$ $VCdim = 1$



can't shatter two points

0	0
0	1
1	0
1	1



• $\mathcal{H} = \{ x \mapsto \text{sgn}(w \cdot x) : w \in \mathbb{R}^d \}$ want to show $VCdim(\mathcal{H}) = d$

$$\text{sgn}(t) = \begin{cases} 1 & \text{if } t \geq 0 \\ -1 & \text{if } t < 0 \end{cases}$$

① \mathcal{H} can shatter $\{e_1, e_2, \dots, e_d\}$ to get (y_1, \dots, y_d)

use $w = (y_1, \dots, y_d)$, $w \cdot e_j = y_j$

(assume) ② Let x_1, \dots, x_{d+1} be shatterable by \mathcal{H} (can't be linear dependent)
 $\therefore \exists \alpha \in \mathbb{R}^{d+1} \setminus \{0\}$, s.t. $\sum_{i=1}^{d+1} \alpha_i x_i = 0$

Let $I_+ = \{i \in [d+1] : \alpha_i > 0\}$, $I_0 = \{i \in [d+1] : \alpha_i = 0\}$, $I_- = \{i \in [d+1] : \alpha_i < 0\}$

$$\exists w \text{ s.t. } h_w(x_j) = \begin{cases} 1 & \text{if } j \in I_+ \\ -1 & \text{if } j \notin I_+ \end{cases}$$

$$0 = w \cdot 0 = w \sum_{i=1}^{d+1} \alpha_i x_i = w \sum_{i \in I_+} \alpha_i x_i + w \sum_{i \in I_-} \alpha_i x_i + w \sum_{i \in I_0} \alpha_i x_i$$

$$= \sum_{i \in I_+} \alpha_i w x_i + \sum_{i \in I_-} \alpha_i w x_i \geq 0 \quad \text{取等 iff } I_- = \emptyset$$

$\{w x_i = 0 \text{ for } \forall i \in I_+\}$

we can also find some \tilde{w} that $\tilde{w} \cdot x_i < 0$ for $\forall i \in I_+$

$$0 = \tilde{w} \cdot 0 = \tilde{w} \cdot \sum_{i \in I_+} \alpha_i x_i = \sum_{i \in I_+} \alpha_i \tilde{w} x_i < 0 \Rightarrow \text{contradiction!}$$

• $VCdim(\{x \mapsto w \cdot x + b : w \in \mathbb{R}^d, b \in \mathbb{R}\}) = d+1$

• $P_{\mathcal{H}}(m) = O(m^{VCdim(\mathcal{H})})$

• Binary classifier with $VCdim(\mathcal{H}) = d$, zero-one loss for any $m > d$:

$$\mathbb{E} \sup_{h \in \mathcal{H}} L_0(h) - L_S(h) \leq \mathbb{E} \sqrt{\frac{2}{m} \log |\mathcal{H}|_{S^X}} \leq \sqrt{\frac{2}{m} \log P_{\mathcal{H}}(m)} \leq \sqrt{\frac{2d}{m} (\log m + \log d)}$$

Corollary, If $m > d := VCdim(\mathcal{H})$, then $P_{\mathcal{H}}(m) \leq \left(\frac{em}{d}\right)^d$ $\log P_{\mathcal{H}}(m) \leq d(\log \frac{m}{d} + 1)$

Sauer-Shelah Lemma. ^① If $d = VCdim(\mathcal{H}) < \infty$, then $|\mathcal{H}(S)| \leq \sum_{i=0}^d \binom{m}{i}$

(we do not prove the lemma)

use S-S Lemma to prove corollary, want to show $\sum_{i=0}^d \binom{m}{i} \leq (\frac{em}{d})^d$ for $m \geq d$

$$\sum_{i=0}^d \binom{m}{i} \leq \sum_{i=0}^d \binom{m}{i} \left(\frac{m}{d}\right)^{d-i}$$

pf for corollary

$$\leq \sum_{i=0}^m \binom{m}{i} \left(\frac{m}{d}\right)^{d-i} \text{ add non-negative}$$

$$= \left(\frac{m}{d}\right)^d \sum_{i=0}^m \binom{m}{i} \left(\frac{d}{m}\right)^i$$

$$= \left(\frac{m}{d}\right)^d \left(1 + \frac{d}{m}\right)^m \quad \left(1 + \frac{x}{n}\right)^n \leq e^x$$

$$\leq \left(\frac{m}{d}\right)^d e^d$$

Lemma ^②: For all finite $S_x \subseteq X$, $|\mathcal{H}|_{S_x} \leq |\{T \subseteq S_x : T \text{ is shattered by } \mathcal{H}\}|$

pf for ① use ② $|\mathcal{H}|_{S_x}$ is upper bounded by the number of subsets of S of size at most d .

pf for ② inductive

These are equivalent for binary class, 0-1 loss:

$\Pr(\sup_{h \in \mathcal{H}} L_D(h) - L_S(h) \leq \epsilon) \geq 1 - \delta$ for $m \geq m(\epsilon, \delta)$

ERM agnostically PAC-learn \mathcal{H}

\mathcal{H} is ag PAC-learnable

ERM PAC-learns \mathcal{H}

\mathcal{H} is PAC-learnable

$VCdim(\mathcal{H}) < \infty$

implies uniform convergence.

Theorem. \mathcal{H} is a binary classifier on X , $d = VCdim(\mathcal{H})$

Assume $m \leq \frac{d}{2}$. then $\inf_A \sup_{\substack{D \\ \text{realizable}}} \Pr_{S \sim D^m} (L_D(A(S)) \geq \frac{1}{8}) \geq \frac{1}{7}$

for a best-case learning algorithm ^{by \mathcal{H}} there is a worst-case distribution

$$\mathbb{E}_{f \sim \text{unif}(\mathcal{X} \rightarrow \mathcal{Y})} \mathbb{E}_{S \sim D^m(f)} L_D(f, \hat{h}_S) = \mathbb{E}_{f \sim S_x} \mathbb{E}_{x \sim D_x} \mathbb{1}(\hat{h}_S(x) \neq f(x))$$

not in the training set

$$\mathbb{E}_{f \sim S_x, x} \mathbb{1}(\hat{h}_S(x) \neq f(x)) = \mathbb{E}_{f \sim S_x} [\Pr(x \notin S_x) \mathbb{E}_{x \sim D_x} [\mathbb{1}(\hat{h}_S(x) \neq f(x)) | x \notin S_x]] + [\Pr(x \in S_x) \mathbb{E}_{x \sim D_x} [\mathbb{1}(\hat{h}_S(x) \neq f(x)) | x \in S_x]]$$

$$P(x \notin S_x) = \frac{|\bar{X} \setminus S_x|}{|\bar{X}|} \geq \frac{m}{2m} = \frac{1}{2} \quad \checkmark \quad S_x = ((x_1, f(x_1)), \dots, (x_m, f(x_m)))$$

$$\geq \frac{1}{2} \mathbb{E}_{f, S_x, x} [\mathbb{1}(\hat{h}_{S_x}(x) \neq h(x) | x \notin S_x)]$$

$$= \frac{1}{2}$$

average over $f \geq \frac{1}{2}$, so $\exists g \in \mathcal{F}$. $\mathbb{E}_{S \sim D^m} L_D(g) \geq \frac{1}{2}$

Want to show $\Pr_{S \sim D^m} (L_{D(g)}(\hat{h}_S) < \frac{1}{8}) \leq \frac{6}{7}$

$$\text{Markov: } \Pr(1 - L_D(g) \geq \frac{7}{8}) \leq \frac{\mathbb{E}(1 - L_D(g))}{\frac{7}{8}} \leq \frac{1 - \frac{1}{2}}{\frac{7}{8}} = \frac{6}{7}$$

with no prior, all algorithms perform the same on average.

Let \mathcal{H} be a set of binary classifier over X . For any $m \geq 1$

$$\inf_{\substack{A \text{ realizable} \\ \text{by } \mathcal{H}}} \sup_{S \sim D^m} \Pr(L_D(A(S)) > \frac{VC \dim(\mathcal{H}) - 1}{32m}) \geq \frac{1}{100}$$

Let $d = VC \dim(\mathcal{H})$

① $d=1$, holds almost trivially

② $d \geq 2m$, this is a corollary of the above theorem


③ $2 \leq d < 2m$

VC dimension doesn't say anything about distribution

If $\mathcal{H} = \{x \rightarrow w \cdot x : w \in \mathbb{R}^d, \|w\| \leq B\}$ $VC(\mathcal{H}) = d$ (scale doesn't matter)

What do we do for approximation error?

lec structural Risk Minimization *not uniform over all H_k*

- $\mathcal{H} = \mathcal{H}_1 \cup \mathcal{H}_2 \cup \dots$  weight: how much do I like H_1, H_2, \dots
 $\sum_{k \geq 1} w_k \leq 1$ preference on each hypothesis's $m \rightarrow \infty$
 $w_k > 0$
- $\forall k \forall D. \Pr_{S \sim D^m} (\sup_{h \in \mathcal{H}_k} L_D(h) - L_S(h) \leq \epsilon_k(m, \delta)) \geq 1 - \delta$ with $\epsilon_k(m, \delta) \rightarrow 0$

$$\therefore \Pr_{S \sim D^m} (\forall h \in \mathcal{H}. L_D(h) \leq L_S(h) + \epsilon_{k_h}(m, \delta w_{k_h})) \geq 1 - \delta$$

$k_h := \operatorname{argmin}_{k \geq 1} \epsilon_k(m, \delta w_k)$ index that minimize ϵ_k corresponding to h
 $h \in \mathcal{H}_{k_h}$ h lives in some of the hypo classes at least one H_k

- ① Def: $SRM_{\mathcal{H}}(S) \in \operatorname{argmin}_{h \in \mathcal{H}} [L_S(h) + \epsilon_{k_h}(m, \delta w_{k_h})]$ commit to a δ beforehand
higher w_k for simpler hypothesis, tighter for complex hypo
- SRM algorithm

We can implement this minimization by a finite number of calls to an "ERM oracle", as long as our loss is lower-bounded by $a \leq \ell(h, z)$ (typically $a = 0$):

```
function SRMH(S)
  best ← ∞
  for k = 1, 2, ... do
    hk ← ERMHk(S)
    cand.loss ← LS(hk) + εk(m, wkδ)
    if cand < best then
      ĥ ← hk
      best ← cand
    if mink > k a + εk(m, wkδ) > best then
      break
  return ĥ
```

wired because depend on δ

Note that if we "decompose" as $\mathcal{H}_1 = \mathcal{H}$, then SRM becomes just $ERM_{\mathcal{H}}$.

- general bound of SRM *non-uniform learnability because m function depend on h*
- $\hat{h}_S := SRM_{\mathcal{H}}(S) \in \operatorname{argmin}_{h \in \mathcal{H}} L_S(h) + \epsilon_{k_h}(m, \delta w_{k_h})$
- $L_D(\hat{h}_S) \leq L_S(\hat{h}_S) + \epsilon_{k_{\hat{h}_S}}(m, \delta w_{k_{\hat{h}_S}})$ (def. prob $\geq 1 - \delta$)
- $\leq L_S(h^*) + \epsilon_{k_{h^*}}(m, \delta w_{k_{h^*}})$ *some h can be learned with less data*
- $= L_D(h^*) + L_S(h^*) - L_D(h^*) + \epsilon_{k_{h^*}}(m, \delta w_{k_{h^*}})$
- $\leq L_D(h^*) + (b-a) \sqrt{\frac{1}{Sm} \log \frac{1}{\delta}} + \epsilon_{k_{h^*}}(m, \delta w_{k_{h^*}})$
- $\Pr(L_D(\hat{h}_S) - L_D(h^*) \leq \epsilon_{k_{h^*}}(m, \delta w_{k_{h^*}}) + (b-a) \sqrt{\frac{1}{Sm} \log \frac{1}{\delta}}) \geq 1 - \delta$
- $h^* \in \mathcal{H}$ is any fixed hypothesis

• Specific bound use Rademacher

$$R_k = \mathbb{E}_{S \sim D^m} \text{Rad}(\ell \circ \mathcal{H}_k | S)$$

$$L_D(h) \leq L_S(h) + 2R_{k_h} + (b-a) \sqrt{\frac{1}{2m} \log \frac{1}{w_{k_h} \delta}} \quad \text{but depend on } \delta$$

$$w_k = \frac{6}{a^2 k^2} \log \frac{1}{w_{k_h} \delta} = \log(k_h^2 \cdot \frac{a^2}{6} \cdot \frac{1}{\delta}) = 2 \log k_h + \log \frac{a^2}{6} + \log \frac{1}{\delta}$$

$$< 2 \log k_h + \frac{1}{2} + \log \frac{1}{\delta}$$

$$= 2 \log k_h + \log \frac{\sqrt{e}}{\delta}$$

$$\sqrt{\log \frac{1}{w_k \delta}} < \sqrt{2 \log k_h} + \sqrt{\log \frac{\sqrt{e}}{\delta}}$$

$$\Pr(\forall h \in \mathcal{H}. L_D(h) \leq L_S(h) + 2R_{k_h} + (b-a) \sqrt{\frac{1}{m} \log k_h} + (b-a) \sqrt{\frac{1}{2m} \log \frac{\sqrt{e}}{\delta}}) \geq 1 - \delta$$

② SRM $\hat{h}_S \in \arg \min L_S(h) + 2R_{k_h} + (b-a) \sqrt{\frac{1}{m} \log k_h}$ does not commit to a δ
 does not depend on δ , better in practice

$$L_D(\hat{h}_S) \leq L_S(\hat{h}_S) + 2R_{k_{\hat{h}_S}} + (b-a) \sqrt{\frac{1}{m} \log k_{\hat{h}_S}} + (b-a) \sqrt{\frac{1}{2m} \log \frac{\sqrt{e}}{\delta}} \quad \left(\begin{array}{l} \text{hold for} \\ \forall h \\ \text{SRM } \hat{h}_S \\ \text{with } p=1-\delta \end{array} \right)$$

then

$$L_D(\hat{h}_S) \leq L_S(\hat{h}_S) + 2R_{k_{\hat{h}_S}} + (b-a) \sqrt{\frac{1}{m} \log k_{\hat{h}_S}} + (b-a) \sqrt{\frac{1}{2m} \log \frac{3}{\delta}} \quad \text{w.p. } > 3$$

(holds with probability $1 - \frac{\sqrt{e}}{3} \delta$)

$$\leq L_S(h^*) + 2R_{k_{h^*}} + (b-a) \sqrt{\frac{1}{m} \log k_{h^*}} + (b-a) \sqrt{\frac{1}{2m} \log \frac{3}{\delta}} \quad (1 - \frac{\sqrt{e}}{3} \delta)$$

Hoeffding $L_S(h^*) \leq L_D(h^*) + (b-a) \sqrt{\frac{1}{2m} \log \frac{3}{\delta}} \quad (1 - \frac{\delta}{3})$

$$\leq L_D(h^*) + 2R_{k_{h^*}} + (b-a) \sqrt{\frac{1}{m} \log k_{h^*}} + (b-a) \sqrt{\frac{3}{2m} \log \frac{3}{\delta}} \quad (1 - \frac{\sqrt{e}+1}{3} \delta \geq 1 - \delta)$$

Comparison: if we know the correct \mathcal{H}_{k_h} from the start and run ERM

$$\text{ERM: } L_D(\text{ERM}_{\mathcal{H}_{k_{h^*}}}) \leq L_D(h^*) + 2R_{k_{h^*}} + (b-a) \sqrt{\frac{2}{m} \log \frac{1}{\delta}}$$

• varying bounded losses

e.g. logistic regression $\mathcal{H}_k = \{x \mapsto w \cdot x \mid \|w\| \leq B_k\}$ $\Pr(\|x\| \leq C) = 1$
 $\hat{h}(x) = \hat{w}^* \cdot x$

if choose $B_k = 2^k$ then $2^{k-1} < \|w\| \leq 2^k \Rightarrow k_h < \log_2(2 \|w\|) \Rightarrow \sqrt{\log k_h} < \sqrt{\log \log_2(2 \|w\|)}$
 $(b-a)k_h < (2C \|w\| + 1)$, thus (not commit to a δ)

$$\hat{h}_S = \hat{h}_{\hat{w}_S}; \quad \hat{w}_S \in \arg \min_{w \in \mathcal{H}} L_S(hw) + \frac{4C \|w\|}{\sqrt{m}} + (2C \|w\| + 1) \sqrt{\frac{1}{m} \log \log_2(2 \|w\|)}$$

$$L_D(\hat{h}_S) \leq L_D(h^*) + \frac{4C \|w^*\|}{\sqrt{m}} + (2C \|w^*\| + 1) \sqrt{\frac{1}{m} \log \log_2(2 \|w\|)} + (2C \|\hat{w}_S\| + 2) \sqrt{\frac{1}{2m} \log \frac{3}{\delta}}$$

\approx

• relationship to regularization

$$\vec{w}_S \in \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} L_D(h^*) + \frac{\lambda}{\sqrt{m}} \|w^*\| \quad \text{SRM} \sim \text{regularization}$$

• non-uniform learnability =

we have shown a bound

$$\Pr(L_D(\text{SRM}_{\mathcal{H}, S}(S)) \leq L_D(h^*) + \epsilon_{\text{KVC}}(m, w_{\text{KVC}} \delta) + (b-a) \sqrt{\frac{1}{2m} \log \frac{1}{\delta}}) \geq 1 - 2\delta$$

the bound does not show PAC learnability because ϵ depends on h^* !!

Def (competes) : $\mathcal{A}(S)(\epsilon, \delta)$ - competes with $h \in \mathcal{H}$ on \mathcal{D}
with m samples, $\Pr_{S \sim \mathcal{D}^m}(L_D(\mathcal{A}(S)) \leq L_D(h) + \epsilon) \geq 1 - \delta$

Def (non-uniform learning) : \mathcal{A} nonuniformly learns \mathcal{H} if
 \exists finite $m(\epsilon, \delta, h)$, s.t. $\forall \epsilon, \delta \in (0, 1)$, $\forall h \in \mathcal{H}$, $\forall \mathcal{D}$,
 $\forall m \geq m(\epsilon, \delta, h)$, $\mathcal{A}(S)(\epsilon, \delta)$ - competes with h on \mathcal{D}

this is looser than PAC because m depend on h

$$\text{PAC: } \exists m \geq m(\epsilon, \delta) \quad \Pr(L_D(\mathcal{A}(S)) \leq \inf_h (L_D(h) + \epsilon)) \geq 1 - \delta \quad \text{for } \forall \epsilon, \delta \in (0, 1)$$

If $\mathcal{H} = \mathcal{H}_1 \cup \mathcal{H}_2 \cup \dots$ with $\text{VCdim}(\mathcal{H}_k) < \infty$

then SRM non-uniformly learns \mathcal{H} (in 0-1 binary classification)

Pf Let $\mathcal{H}_k = \{h \in \mathcal{H} : m(\frac{1}{8}, \frac{1}{7}, h) \leq k\}$

$\therefore \mathcal{H} = \mathcal{H}_1 \cup \mathcal{H}_2 \cup \dots$

consider \mathcal{D} realizable by \mathcal{H}_k

then $\Pr_{S \sim \mathcal{D}^k}(L_D(\mathcal{A}(S)) \leq \frac{1}{8}) \geq \frac{1}{7}$ using $h^* \in \mathcal{H}_k$ with $L_D(h^*) = 0$

Assume $\mathcal{H} = \{h_1, h_2, \dots\}$ ^{countable} - divide singleton classes

Use $\mathcal{H}_k = \{h_k\}$

$$\text{Then } \sum_k (m, w_k \delta) \leq (b-a) \sqrt{\frac{1}{2m} \log \frac{1}{w_k \delta}}$$

$$\leq (b-a) \sqrt{\frac{1}{2m} \log \frac{1}{w_k}} + (b-a) \sqrt{\frac{1}{2m} \log \frac{1}{\delta}}$$

how to assign weights?
how to decide the order?

Minimum Description Length ^{prefix-free}

choose a weight according to a binary language for determining \mathcal{H}

Kraft's inequality :=

If $S \subseteq \{0, 1\}^*$ (a set of binary strings) is prefix-free

(if '00' is a valid s, we don't have anything else that starts with '00')
(or read the data before '0' in (language))

then $\sum_{s \in S} 2^{-|s|} \leq 1$ (probability distribution)



Then, we can choose a representation for \mathcal{H} and assign $w_h = 2^{-|h|}$

$$\text{MDL}_S \in \underset{h \in \mathcal{H}}{\text{argmin}} L_S(h) + (b-a) \sqrt{\frac{1}{2m} \log \frac{1}{\delta 2^{-|h|}}}$$

$$\sqrt{\frac{1}{2m} (\log(\frac{1}{\delta}) + |h| \log 2)}$$

or $\text{MDL}_S \in \underset{h \in \mathcal{H}}{\text{argmin}} L_S(h) + \sqrt{\frac{\log 2}{2} \frac{|h|}{m}}$

$$\therefore L_0(\text{MDL}(S)) \leq L_0(h^*) + (b-a) \left(\sqrt{\frac{\log 2}{2} \cdot \frac{|h^*|}{m}} + \sqrt{\frac{2}{m} \log \frac{2}{\delta}} \right)$$

Occam's razor: if there are multiple explanations of the data ($L_S(h_1) = L_S(h_2)$) prefer the simpler one (shortest explanation)

If we choose $|h|$ to be the length of the shortest possible implementation of h in some programming language, is known Kolmogorov Complexity
MAP inference with a Kolmogorov prior

Margins Theory

We do ERM with 0-1 loss

VC theory

$$\mathcal{H} = \{x \mapsto \text{sgn}(w \cdot x) : x \in \mathbb{R}^d\} \quad \text{VCdim}(\mathcal{H}) = d \quad \text{known from previous lec}$$

$$\sup_{h \in \mathcal{H}} L_D(h) - L_S(h) \leq \sqrt{\frac{2d}{m} (\log m + 1 - \log d)} + \sqrt{\frac{1}{2m} \log \frac{1}{\delta}} \quad (\text{prob} \geq 1 - \delta)$$

$$L_D(\hat{h}_S) - \inf_{h \in \mathcal{H}} L_D(h) \leq \sqrt{\frac{2d}{m} (\log m + 1 - \log d)} + \sqrt{\frac{2}{m} \log \frac{2}{\delta}} \quad (\text{prob} \geq 1 - \delta)$$

two problems:

① ERM with 0-1 loss is NP hard if \mathcal{D} is not realizable by \mathcal{H}

② if d is really big, the bound doesn't tell anything until $\frac{m}{\log m} > 2d$

(like in kernel methods, d is sometimes infinite!)

we like a better bound when in high dimension (big d) (hope it is not d -dependent)

Rademacher

$$\text{If } \mathcal{H}_B = \{x \mapsto w \cdot x : \|w\| \leq B\} \text{ then } \mathbb{E} \text{Rad}(\mathcal{H}_B | \mathcal{S}) \leq \frac{B \sqrt{\mathbb{E} \|x\|^2}}{\sqrt{m}}$$

but \uparrow is continuous, (logistic loss)

$$\text{Rad}(\text{logistic} \circ \mathcal{H}_B | \mathcal{S}) \leq 1 \cdot \text{Rad}(\mathcal{H}_B | \mathcal{S})$$

$$\therefore L_D^{\text{logistic}}(\text{argmin}_{h \in \mathcal{H}_B} L_D^{\text{logistic}}(h)) \leq \frac{2B \sqrt{\mathbb{E} \|x\|^2}}{\sqrt{m}} + (B \sqrt{\mathbb{E} \|x\|^2} + 1) \sqrt{\frac{2}{m} \log \frac{2}{\delta}}$$

what if we want to bound on accuracy

$$\text{Rad}(l_{0-1} \circ \text{sgn} \circ \mathcal{H}_B) \quad \left\{ \begin{array}{l} l_{0-1} \text{ not Lipschitz} \\ \text{VCdim}(\mathcal{H}_B) = d \end{array} \right.$$

we can use a surrogate loss

$$\forall h, z \quad l_{\text{surr}}(h, z) \geq l_{0-1}(h, z)$$

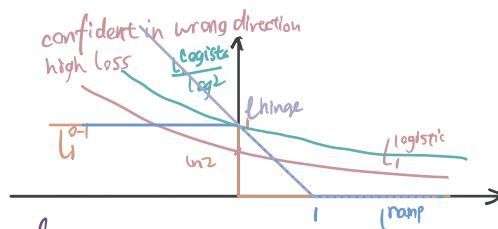
$$L_D^{\text{surr}}(h) \geq L_D^{0-1}(h)$$

$$L_D^{0-1}(\text{sgn} \circ h) \leq L_D^{\text{surr}}(h) \leq L_S^{\text{surr}}(h) + 2 \text{Rad}(l_{\text{surr}} \circ \mathcal{H} | \mathcal{S}) + (b-a)_{\text{surr}} \sqrt{\frac{1}{2m} \log \frac{1}{\delta}}$$

① use logistic loss

$$\text{let } l_{\text{surr}}(h, z) = \frac{1}{\log 2} l_{\text{logistic}}(h, z)$$

but logistic loss gives a loose bounds (logistic \sim l_{0-1} not close enough)



if $y = 1$

0-1 loss $1_{h \leq 0}$

logistic loss (LR) $\log(1 + \exp(-h))$

hinge loss (SVM) $[1-h]_+$ ($s=0$)

ramp loss (S) $[1-h]_+ - [s-h]_+$

② use ramp loss $L_{\text{ramp}}(h, (x, y)) = L_y(h(x)) = \begin{cases} 1 & y h(x) \leq 0 \\ 1 - y h(x) & 0 \leq y h(x) \leq 1 \\ 0 & 1 \leq y h(x) \end{cases}$

1-Lipschitz, and bounded in $[0, 1]$

$$\forall h \quad L_D^{0-1}(\text{sgn} \circ h) \leq L_D^{\text{ramp}}(h) \leq L_S^{\text{ramp}}(h) + 2 \frac{BC}{\sqrt{m}} + \sqrt{\frac{1}{2m} \log \frac{1}{\delta}} \quad (B \geq 1 - \delta)$$

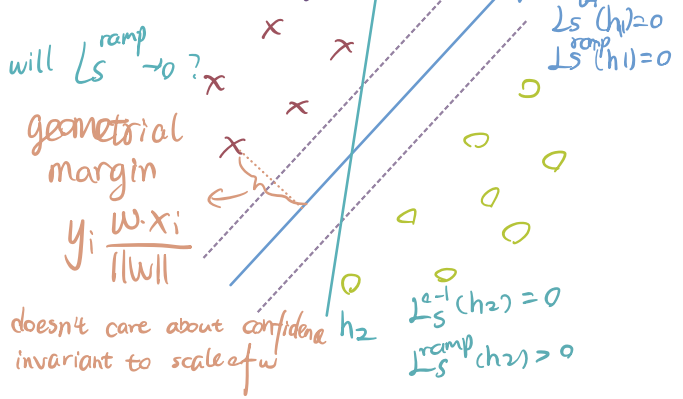
better than logistic loss

so ... what about L_S^{ramp} term?

$$\text{margin} = \min_i w \cdot x_i y_i$$

least confident sample is sensitive to scale of w

the larger the norm of w the smaller the margin \times



Assume $\exists h^* \in \mathcal{H}$, s.t. $L_D^{\text{ramp}}(h^*) = 0$. And $\mathbb{E} \|x\|^2 \leq C^2$

D is linearly separable with margin

$$\therefore L_S^{\text{ramp}}(h^*) = 0 \quad w^* = \arg \min \|w\| \text{ s.t. } L_D^{\text{ramp}}(h^*) = 0$$

$$\text{For ERM: } \hat{w} = \arg \min \|w\| \text{ s.t. } L_S^{\text{ramp}}(hw) = 0 \Rightarrow \|\hat{w}\| \leq \|w^*\|$$

$$\therefore L_D^{0-1}(\hat{h}_S) \leq \frac{2C \|w^*\|}{\sqrt{m}} + \sqrt{\frac{1}{2m} \log \frac{1}{\delta}} = \frac{2C}{\rho \sqrt{m}} + \sqrt{\frac{1}{2m} \log \frac{1}{\delta}}$$

$$\min_i \frac{w^* \cdot x_i}{\|w^*\|} y_i \geq \frac{1}{\|w^*\|} \quad \text{let } \rho = \frac{1}{\|w^*\|} \text{ (scale } w_p = \rho w^*)$$

$$\text{then } \min_i \frac{w^* \cdot x_i y_i}{\|w^*\|} = \min_i \rho w^* \cdot x_i y_i$$

this bound depends on $\|w^*\|$, which we don't really know

so use a SRM-like argument to get bound only on $\|\hat{w}\|$

$$L_D^{0-1}(\hat{h}_S) \leq L_S^{\text{ramp}}(hw) + \frac{1}{\sqrt{m}} \left(\sqrt{\frac{1}{2} \log \frac{1}{\delta}} + \begin{cases} 4Cr & \text{if } \|w\| \leq r \\ 4C\|w\| + \sqrt{\log \log \frac{2\|w\|}{r}} & \text{if } \|w\| > r \end{cases} \right)$$

Pf Define $B_k = r 2^k$ $\delta_k = \frac{6\delta}{2^k k^2}$ for all $k \geq 1$ $\sum_{k=1}^{\infty} \delta_k = \delta$

$$\forall k, \Pr_{S \sim D^m} (\forall h \in \mathcal{H}_{B_k}, L_D^{0-1}(\text{sgn} \circ h) \leq L_D^{\text{ramp}}(h) \leq L_S^{\text{ramp}}(h) + \underbrace{2\mathbb{E} \text{Rad}(\mathcal{H}_{B_k} |_{S_k})}_{\leq \frac{2B_k C}{\sqrt{m}}} + \underbrace{\sqrt{\frac{1}{m} \log \frac{1}{\delta_k}}}_{\leq \frac{1}{\sqrt{m}} (\sqrt{\log \frac{3}{\delta}} + \sqrt{2 \log k})} \geq 1 - \delta_k$$

To get a union bound, look a particular w and ask which k

If $\|w\| \leq 2r$, $h \in \mathcal{H}_B \Rightarrow k=1 \Rightarrow \log k = 0$, $B_k = 2r$

$$\forall w: \|w\| \leq 2r, L_D^{0-1}(\text{sgn} \circ h) \leq L_S^{\text{ramp}}(h) + \frac{4Cr}{\sqrt{m}} + \sqrt{\frac{1}{2m} \log \frac{1}{\delta}}$$

$$\text{If } \|w\| > 2r, L_D^{0-1}(\text{sgn} \circ h) \leq L_S^{\text{ramp}}(h) + \frac{4C\|w\|}{\sqrt{m}} + \sqrt{\frac{1}{2m} \log \frac{1}{\delta}} + \sqrt{\frac{1}{m} \log \log \frac{2\|w\|}{r}}$$

\uparrow If $\|w\| > 2r$, $\|w\| \leq r2^k \Rightarrow k = \lceil \log_2 \frac{\|w\|}{r} \rceil \Rightarrow h \in H_{B_{kw}}$
 $B_{kw} = r \cdot 2^{\lceil \log_2 \frac{\|w\|}{r} \rceil} < r \cdot 2^{\log_2 \frac{\|w\|}{r} + 1} = 2\|w\|$
 $\log k < \log(\log_2 \frac{\|w\|}{r} + 1) = \log_2 \log_2 \frac{2\|w\|}{r}$ can choose $r = \|w\|$

this bound doesn't imply non-uniform learning because its dependence on data distribution

Lec SVM

D is separable with a margin if $\exists h^* \text{ s.t. } L_D(h^*) = 0$

Expand L_S^{ramp}

$\hat{h} = h_w$; $\hat{w} \in \text{argmin}_w \|w\|^2 \text{ s.t. } \forall i \in [m], y_i w \cdot x_i \geq 1$ (Hard SVM)

convex quadratic program

$\hat{w} = \text{argmin}_w \|w\|^2 \text{ s.t. } \forall i \in [m], y_i w \cdot x_i \geq 1$
 $= \text{argmax}_w \frac{1}{\|w\|} [\min_i y_i w \cdot x_i]$, s.t. $\forall i \in [m], y_i w \cdot x_i \geq 1$
at least 1, also no bigger than 1 because could scale down w

$\geq \text{argmax}_w \min_j \frac{y_j w \cdot x_j}{\|w\|} \text{ s.t. } \forall i \in [m], y_i w \cdot x_i > 0$
maximize the worst case geometric margin
 \uparrow invariant to scale of w \rightarrow so we could relax it to positive by scaling $L_S^{-1}(hw) = 0$

Duality:

$\min_w \frac{1}{2} \|w\|^2 \text{ s.t. } \forall i \in [m], y_i w \cdot x_i \geq 1$
 $= \min_w \max_{\alpha_i \geq 0} \frac{1}{2} \|w\|^2 + \sum_{i=1}^m \alpha_i (1 - y_i w \cdot x_i)$
— dual variable

$= \max_{\alpha_i \geq 0} \min_w \frac{1}{2} \|w\|^2 + \sum_{i=1}^m \alpha_i (1 - y_i w \cdot x_i)$

$\nabla_w \frac{1}{2} \|w\|^2 + \sum_{i=1}^m \alpha_i (1 - y_i w \cdot x_i)$

$= w - \sum_{i=1}^m \alpha_i y_i x_i = 0$

$\Rightarrow w = \sum_{i=1}^m \alpha_i y_i x_i \quad \|w\|^2 = \sum_{i=1}^m \sum_{j=1}^m \alpha_i y_i x_i \cdot x_j y_j \alpha_j = \alpha^T \text{diag}(y) X X^T \text{diag}(y) \alpha$
 $X \in \mathbb{R}^{m \times d}$
 $w \cdot x_j = \sum_{i=1}^m \alpha_i y_i x_i \cdot x_j$

$= \max_{\alpha \geq 0} \frac{1}{2} \alpha - \frac{1}{2} \alpha^T \text{diag}(y) X X^T \text{diag}(y) \alpha$

$X X^T \in \mathbb{R}^{m \times m}$: Gram matrix $(X X^T)_{ij} = x_i \cdot x_j$

$\min_x \max_y f(x, y) \geq \max_y \min_x f(x, y)$
(weak duality)

$\min_x \max_y f(x, y) = \max_y \min_x f(x, y)$
(strong duality)

\rightarrow we have strong duality here by Slater's condition

which problem to solve (min or max?) depends on m, d which is smaller in \mathbb{R}^d , a lot of α_i will be zero

what if the data is not linearly separable?

Soft SVM:

recall the SRM-like bound, two problems:

1) the "r" is annoying

we choose a small enough "r" that $r \leq \|w\|$ for all reasonable w but not so small that $\sqrt{\log \log \frac{2\|w\|}{r}}$ is relevant to anything.

Ignore $\sqrt{\log \log \frac{2\|w\|}{r}}$ term to pick

$$\min_w L_S^{\text{ramp}}(hw) + \frac{4C}{\sqrt{m}} \|w\| \quad \text{still NP hard}$$

2) the problem is still NP hard, not a practical algorithm.

we can take $l_{\text{hinge}} = l_{\text{ramp}} \circ l_{\circ-1}$, $l_{\text{hinge}}(h, (x, y)) = \begin{cases} 1 - yhx & \text{if } yhx \leq 1 \\ 0 & \text{if } yhx > 1 \end{cases}$

① 1-Lipschitz ② not bounded ③ convex

get more loss for a more-confident wrong answer

Algorithm =

$$\text{hinge version} \approx \min_w L_S^{\text{hinge}}(hw) + \frac{4C}{\sqrt{m}} \|w\|$$

$$(a) \approx \min_w L_S^{\text{hinge}}(hw) + \lambda \|w\|^2 \quad \begin{matrix} \uparrow \\ \text{square norm is much easier} \\ \text{to optimize} \end{matrix}$$

C does not necessarily provide tight bound, λ should scale with $\frac{1}{\sqrt{m}}$, usually

$$(b) = \frac{1}{2\lambda} \min_{0 \leq \alpha_i \leq \frac{1}{2\lambda m}} 1^T \alpha - \frac{1}{2} \alpha^T \text{diag}(y) X X^T \text{diag}(y) \alpha$$

not equal in value but in min operation (X X^T)_{ij} = X_i \cdot X_j

$$(b) \text{ predict with } h(x) = \sum_{i=1}^m \alpha_i y_i X_i \cdot X$$

hinge ERM + Bounded weights

$$H_B = \{x \rightarrow w \cdot x : w \in \mathbb{R}^d, \|w\| \leq B\}$$

if $\hat{h}_B = \arg \min_{h \in H_B} L_S^{\text{hinge}}(hw)$ since $l^{\text{hinge}} \geq l^{\text{ramp}}$

$$L_D^{\alpha-1}(\text{sgn} \circ \hat{h}_B) \leq L_S^{\text{hinge}}(\hat{h}_B) + \frac{2BC}{\sqrt{m}} + \sqrt{\frac{1}{2m} \log \frac{1}{\delta}} \quad (Pr \geq 1-\delta)$$

l^{hinge} is unbounded, but $h(x)$ is $\sup_{h \in H} |h(x)| \leq H \sup_{h \in H} \|w\| \|x\| = HBC$

$$L_D^{\alpha-1}(\text{sgn} \circ \hat{h}_B) \leq \inf_{h \in H_B} L_D^{\text{hinge}}(h) + \frac{2BC}{\sqrt{m}} + (2+BC) \sqrt{\frac{1}{2m} \log \frac{1}{\delta}} \quad (Pr \geq 1-\delta)$$

Lec. kernel (real-valued)

$$h(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3 \quad c \in \mathbb{R}^+$$

$$= \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix} \cdot \underbrace{\begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \end{bmatrix}}_{\phi(x)} = \begin{bmatrix} \sqrt{c^3} w_0 \\ \sqrt{3c^2} w_1 \\ \sqrt{3c} w_2 \\ w_3 \end{bmatrix} \cdot \underbrace{\begin{bmatrix} \sqrt{c^3} \\ \sqrt{3c^2} x \\ \sqrt{3c} x^2 \\ x^3 \end{bmatrix}}_{\phi_c(x)} \quad \text{scale them with constant } c$$

$$\phi(x) \cdot \phi(x') = 1 + xx' + (xx')^2 + (xx')^3$$

$$\underbrace{\phi_c(x) \cdot \phi_c(x')}_{k(x, x')} = c^3 + 3c^2(xx') + 3c(xx')^2 + (xx')^3 = (xx' + c)^3 \quad \leftarrow \text{easier to compute}$$

kernel function

\mathcal{F} is a vector space of functions

$$f, f' \in \mathcal{F} \quad f + f' \in \mathcal{F}, \quad a f \in \mathcal{F} \quad \forall a \in \mathbb{R}, \quad a(bf) = (ab)f \dots$$

$$\langle f, f' \rangle_{\mathcal{F}} = \int f(x) f'(x) dx$$

so $\mathcal{F}_c = \{x \mapsto w \cdot \phi_c(x) : w \in \mathbb{R}^d\}$

$\mathcal{F}' = \{x \mapsto w \cdot \phi(x) : w \in \mathbb{R}^d\}$

Def: f has weights w , f' has weights w' , f and $f' \in \mathcal{F}$

$$\Rightarrow \left\{ \begin{array}{l} \textcircled{1} f + f' \text{ have weights } w + w' \\ \textcircled{2} \alpha f \text{ have weights } \alpha w \end{array} \right\} \text{ vector space}$$

$$\textcircled{3} \langle f, f' \rangle_{\mathcal{F}} = w \cdot w' \quad \left. \vphantom{\textcircled{3}} \right\} \text{ (real) Hilbert space}$$

$$\textcircled{4} \|f\|_{\mathcal{F}} = \sqrt{\langle f, f \rangle_{\mathcal{F}}} = \|w\|$$

$\textcircled{1}$ symmetric $\langle y, x \rangle = \langle x, y \rangle$

$\textcircled{2}$ linear $\langle a x_1 + b x_2, y \rangle = a \langle x_1, y \rangle + b \langle x_2, y \rangle$

$\textcircled{3}$ complete $\sum_{i=1}^{\infty} \|x_i\| < \infty$

property: $\forall x, \alpha f(x) + f'(x) = [\alpha f + f'](x)$

$$\langle 0, 0 \rangle = 0$$

$$\langle f, f \rangle > 0 \text{ if } f \neq 0$$

$$\langle \alpha f + g, h \rangle = \langle \alpha f, h \rangle + \langle g, h \rangle$$

$$\langle f, g \rangle = \langle g, f \rangle$$

although \mathcal{F}_c and \mathcal{F}' are the same set, $\|f\|_{\mathcal{F}_c} \neq \|f\|_{\mathcal{F}'}$ for any f

identity fun

if $\mathcal{F} = \{x \mapsto w \cdot \phi(x) : w \in \mathbb{R}^d\}$ $\phi(x) = x$ then $\|f\|_{\mathcal{F}} = \|w\|$

Soft SVM: $\min_w L_S^{\text{hinge}}(hw) + \lambda \|w\|^2$

$\Rightarrow \min_{h \in \mathcal{F}} L_S^{\text{hinge}}(hw) + \lambda \|h\|_{\mathcal{F}}^2$

what kind of function can be a kernel?

kernel is the inner product of feature maps over Hilbert space
if and only if

An example of kernel

let $\phi(x) = (\sqrt{c^3}, \sqrt{3c^2}x, \sqrt{3c}x^2, x^3) \in \mathbb{R}^4$, $\mathcal{F} = \{x \mapsto w \cdot \phi(x) : w \in \mathbb{R}^4\}$

consider $\phi(x)$ as a weight vector for an element in \mathcal{F}

$$\begin{aligned} \Rightarrow f(x) : x' \rightarrow \sqrt{c^3}\sqrt{c^3} + \sqrt{3c^2}x\sqrt{3c^2}x' + \sqrt{3c}x^2\sqrt{3c}x'^2 + x^3x'^3 \\ = (xx' + c)^3 = k(x, x') = \langle \phi(x), \phi(x') \rangle \end{aligned}$$

define a kernel

⌋ know the feature function $\phi(x) \in \mathcal{F}$, then $k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{F}}$
⌋ don't know $\phi(x)$, make sure the kernel matrix is positive semi-definite

Def. A function $k: X \times X \rightarrow \mathbb{R}$ is a positive definite kernel if and only if there exist some Hilbert space \mathcal{F} and feature map $\phi(x) : X \rightarrow \mathcal{F}$ such that $k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{F}}$

positive semi-definite matrix is equivalently characterized as

- ① For all $\alpha \in \mathbb{R}^m$, $\alpha^T K \alpha \geq 0$
- ② All eigenvalues are non-negative
- ③ $K = LL^T$ for some $L \in \mathbb{R}^{m \times m}$

Thm. A function $k: X \times X \rightarrow \mathbb{R}$ is a positive definite kernel if and only if for all $m \geq 1$ and $x_1, \dots, x_m \in X$, the kernel matrix K is positive semi-definite

$$K = \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_m) \\ \vdots & \ddots & \vdots \\ k(x_m, x_1) & \dots & k(x_m, x_m) \end{bmatrix} \in \mathbb{R}^{m \times m}$$

Pf. ① k is a kernel $\Rightarrow \alpha^T K \alpha \geq 0$

if $k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{F}}$, then

$$\alpha^T K \alpha = \sum_{i=1}^m \sum_{j=1}^m \alpha_i \langle \phi(x_i), \phi(x_j) \rangle \alpha_j = \left\| \sum_{i=1}^m \alpha_i \phi(x_i) \right\|_{\mathcal{F}}^2 \geq 0$$

② $\alpha^T K \alpha \geq 0 \Rightarrow k$ is a kernel $\begin{cases} \exists \text{ Hilbert space } \mathcal{H} \\ \exists \phi: X \rightarrow \mathcal{H} \end{cases}$ s.t. $k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$

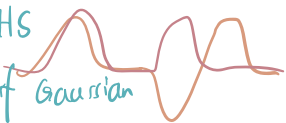
Reproducing kernels

define $\varphi(x) = \{x' \mapsto k(x, x')\}$ for all x

let \mathcal{F}_0 be set of all linear combinations of these functions $\sum_{i=1}^m \alpha_i \varphi(x_i)$
 for all $x_1, \dots, x_m \in X, \alpha_1, \dots, \alpha_m \in \mathbb{R}$. $\text{Range}(\varphi) = \{\varphi(x) : x \in X\} \subseteq \mathcal{F}_0$ linearity

example (1) $k(x, x') = x \cdot x', x \in \mathbb{R}^d, \varphi(x) = [x' \mapsto x \cdot x']$ is linear on x'
 the set of all $\varphi(x) \{ \varphi(x) : x \in X \} = \{ x' \mapsto w \cdot x' : w \in \mathbb{R}^d \}$
 so $\sum \alpha_i \varphi(x_i)$ also in \uparrow don't need to worry about linearity
 $\|\varphi(x)\|_{\mathcal{F}} = \sqrt{\langle \varphi(x), \varphi(x) \rangle_{\mathcal{F}}} = \sqrt{k(x, x)} = \sqrt{\|x\|^2} = \|x\|$

(2) $k(x, x') = (x \cdot x' + c)^3, \phi(x) = (\sqrt{c}, \sqrt{c^2}x, \sqrt{c^3}x^2, \sqrt{c^3}x^3) \in \mathbb{R}^4$ or $(1, x, x^2, x^3) \dots$
 $\varphi(x) = [x' \mapsto (xx' + c)^3] \in \mathcal{F}$ feature map $\nexists w$ s.t. $f = \varphi(w)$
 \uparrow it cannot represent functions like $f(x) = -1 + x^3$, but this could be represented through linear combination of $\varphi(x)$, so $f(x) \in \mathcal{F}$

(3) $k(x, x') = \exp(-\frac{1}{2\sigma^2} \|x - x'\|^2)$ functions in RKHS
 Gaussian kernel 
 linear combination of Gaussian

to make \mathcal{F}_0 Hilbert

linear symmetric define $\langle \sum_{i=1}^m \alpha_i \varphi(x_i), \sum_{j=1}^n \beta_j \varphi(x'_j) \rangle_{\mathcal{F}_0} = \sum_{i=1}^m \sum_{j=1}^n \alpha_i \beta_j k(x_i, x'_j)$

then we have $\langle \varphi(x), \varphi(x') \rangle_{\mathcal{F}_0} = k(x, x') =$

reproducing $\rightarrow \langle f, \varphi(x) \rangle_{\mathcal{F}_0} = f(x) \quad f \in \mathcal{F}_0$

property $\langle \sum_{i=1}^m \alpha_i \varphi(x_i), f \rangle_{\mathcal{F}_0} = \sum_{i=1}^m \alpha_i f(x_i)$

$$\|f\|_{\mathcal{F}} = |\langle f, \varphi(x) \rangle_{\mathcal{F}}| \leq \|f\|_{\mathcal{F}} \|\varphi(x)\|_{\mathcal{F}} = \sqrt{k(x, x)} \|f\|_{\mathcal{F}}$$

$$|f(x) - f(x')| \leq |\langle f, \varphi(x) - \varphi(x') \rangle_{\mathcal{F}}| \leq \|f\|_{\mathcal{F}} \sqrt{k(x, x) + k(x', x') - 2k(x, x')}$$

\Rightarrow complete add the limits of all Cauchy sequences and define the inner product to construct RKHS \mathcal{F} based on \mathcal{F}_0

not all $f \in \mathcal{F}$ can be written as $\sum_{i=1}^n \alpha_i \varphi(x_i)$, but can get close

Optimizing in the RKHS

argmin $L_S(f)$ for Gaussian kernel f is in infinite dimension

$f: \|f\|_{\mathcal{F}} \leq B$ space, $\|f\|_{\mathcal{F}}$ is a problem

argmin $L_S(f) + \lambda \|f\|_{\mathcal{F}}^2$ $f = \sum_{i=1}^m \alpha_i y_i \varphi(x_i)$ (soft SVM)

general form: $\operatorname{argmin}_{f \in \mathcal{F}} A(f(x_1), \dots, f(x_m)) + R(\|f\|_{\mathcal{F}})$

↑ non-decreasing, like 0-inf indicator func

the solution will be some linear combination of $\varphi(x_i)$

Theorem (Representer theorem). If \mathcal{F} is an RKHS with feature map φ , then for any function $A: \mathbb{R}^m \rightarrow \mathbb{R}$ and nondecreasing function $R: \mathbb{R} \rightarrow \mathbb{R} \cup \{\infty\}$

$$\operatorname{argmin}_{f \in \mathcal{F}} A(f(x_1), \dots, f(x_m)) + R(\|f\|)$$

contains a solution $f = \sum_{i=1}^m \alpha_i \varphi(x_i)$, if R is strictly increasing then every solution is ↑ this form.

so can solve $\operatorname{argmin}_{\alpha \in \mathbb{R}^m} A(\dots) + R(\sqrt{\alpha^T K \alpha})$

and then $\hat{f} = \sum_i \hat{\alpha}_i \varphi(x_i)$

↑ $f(x_j) = \sum_{i=1}^m \alpha_i k(x_i, x_j)$

Pf.

$A(f(x_1), \dots, f(x_m)) = A(\langle f, \varphi(x_1) \rangle_{\mathcal{F}}, \dots, \langle f, \varphi(x_m) \rangle_{\mathcal{F}})$

Let $\mathcal{F}_{\parallel} = \operatorname{span}(\{\varphi(x_i) : i \in [m]\}) = \{ \sum_i \alpha_i \varphi(x_i) : \alpha \in \mathbb{R} \}$ parallel to data

Let \mathcal{F}_{\perp} = orthogonal subspace of \mathcal{F}_{\parallel} in \mathcal{F}

$f_{\parallel} \in \mathcal{F}_{\parallel}, f_{\perp} \in \mathcal{F}_{\perp}$, then $\langle f_{\parallel}, f_{\perp} \rangle_{\mathcal{F}} = 0$

$\forall f \in \mathcal{F}, f = f_{\parallel} + f_{\perp} \quad f(x_i) = \langle f_{\parallel} + f_{\perp}, \varphi(x_i) \rangle_{\mathcal{F}} = \langle f_{\parallel}, \varphi(x_i) \rangle_{\mathcal{F}}$

$\therefore A(f(x_1), \dots, f(x_m)) = A(\underbrace{f_{\parallel}(x_1), \dots, f_{\parallel}(x_m)}_{\text{having nonzero } f_{\perp} \text{ doesn't change } A})$

$\|f\|_{\mathcal{F}}^2 = \|f_{\parallel} + f_{\perp}\|_{\mathcal{F}}^2 = \|f_{\parallel}\|_{\mathcal{F}}^2 + \|f_{\perp}\|_{\mathcal{F}}^2 + 2 \langle f_{\parallel}, f_{\perp} \rangle_{\mathcal{F}}$

\therefore always a solution in \mathcal{F}_{\parallel} having nonzero f_{\perp} does not help R

Eg kernel ridge regression

$\operatorname{argmin}_{f \in \mathcal{F}} L_s^{sq}(f) + \lambda \|f\|_{\mathcal{F}}^2$

$L_s^{sq}(f) = \frac{1}{m} \sum_{i=1}^m (f(x_i) - y_i)^2$

$\operatorname{argmin}_{\alpha \in \mathbb{R}^m} \frac{1}{m} \sum_{i=1}^m \left(\sum_{j=1}^m \alpha_j k(x_i, x_j) - y_i \right)^2 + \lambda \alpha^T K \alpha$

$= \operatorname{argmin}_{\alpha \in \mathbb{R}^m} \frac{1}{m} \frac{\|K\alpha - y\|^2}{\alpha^T K \alpha - 2y^T K \alpha + y^T y} + \lambda \alpha^T K \alpha$

$$= \operatorname{argmin}_{\alpha \in \mathbb{R}^m} \alpha^T K (\kappa \alpha + m \lambda I) \alpha - 2 y^T K \alpha$$

let gradient = 0 $2K(K + m\lambda I)\alpha = 2Ky$

$$\alpha = (K + m\lambda I)^{-1} y$$

what kind of things kernel can / can't do ?

Rademacher complexity

$f \in \mathcal{F}$, RKHS for kernel. $K(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{F}}$

$f: X \rightarrow \mathbb{R}$

$$\mathcal{H}_B = \{f \in \mathcal{F} : \|f\|_{\mathcal{F}} \leq B\}$$

$$|f(x)| = |\langle f, \varphi(x) \rangle_{\mathcal{F}}| \leq \|f\|_{\mathcal{F}} \|\varphi(x)\|_{\mathcal{F}}$$

$$m \cdot \operatorname{Rad}(\mathcal{H}_B | s_X) = \mathbb{E}_{\sigma} \sup_{\|f\|_{\mathcal{F}} \leq B} \sum_i \sigma_i f(x_i)$$

$$\|\varphi(x)\|_{\mathcal{F}}^2 = \langle \varphi(x), \varphi(x) \rangle_{\mathcal{F}} = K(x, x)$$

$$= \mathbb{E}_{\sigma} \sup_{\|f\|_{\mathcal{F}} \leq B} \langle f, \sum_i \sigma_i \varphi(x_i) \rangle_{\mathcal{F}}$$

$$\leq B \mathbb{E}_{\sigma} \left\| \sum_i \sigma_i \varphi(x_i) \right\|_{\mathcal{F}}$$

$$\leq B \sqrt{\mathbb{E}_{\sigma} \left\| \sum_i \sigma_i \varphi(x_i) \right\|_{\mathcal{F}}^2}$$

$$= B \sqrt{\sum_i \mathbb{E}_{\sigma} \|\varphi(x_i)\|_{\mathcal{F}}^2} = B \sqrt{\sum_i K(x_i, x_i)} = 1 \text{ for gaussian kernel}$$

$$\operatorname{Rad}(\mathcal{H}_B | s_X) \leq \frac{B}{\sqrt{m}} \sqrt{\frac{1}{m} \sum_i K(x_i, x_i)}$$

Lec Universal approximation

includes all limiting values of points

Euclidean space

X is a compact metric space e.g. a closed bounded subset of $\mathbb{R}^d, [0, 1]^d$

$C(X)$ is the space of continuous functions $X \rightarrow \mathbb{R}$

with $\|f\|_{\infty} = \sup_{x \in X} |f(x)|$

universal approximation

Def (universal kernel). A kernel $K: X \times X \rightarrow \mathbb{R}$ is universal if its

RKHS \mathcal{F} is dense in $C(X)$:

$$\forall g \in C(X), \forall \epsilon > 0, \exists f \in \mathcal{F} \text{ st. } \|g - f\|_{\infty} = \sup_{x \in X} |g(x) - f(x)| \leq \epsilon$$

every finite dimension kernel will approach 0 if $x \rightarrow \infty$
 linear kernel is not universal

Prop. For any universal RKHS, $VCdim(\mathcal{F}) = \infty$

Let $V, W \in X$ be compact, disjoint sets. Let K be universal

Then $\exists f \in \mathcal{F}_K$ s.t. $\forall x \in V, f(x) > 0$
 $\forall x \in W, f(x) < 0$ shattering

\Rightarrow which means $VCdim(\mathcal{F}) = \infty$

Pf. Let $dist_V(x) = \min_{v \in V} \|x - v\|$

$$\text{Let } g(x) = \frac{dist_W(x) - dist_V(x)}{dist_W(x) + dist_V(x)} \quad \begin{cases} \forall x \in V, g(x) = 1 & (dist_V(x) = 0) \\ \forall x \in W, g(x) = -1 & (dist_W(x) = 0) \end{cases}$$

then $\exists f \in \mathcal{F}_{st} \|f - g\|_\infty < 1$ separate compact set

Prop. For any universal RKHS, $Rad(\mathcal{F}|_{S_x}) = \infty$

for any σ_i , can find some value that are positively correlated.

so Rad at least positive. If multiply f by a (a really large)

then the whole thing is scaled by a . thus $Rad(\mathcal{F}|_{S_x})$ is infinite.

universal approximation of neural networks

$$f(x) = f^{(D)}(x) \quad D \text{ layers}$$

$$f^k(x) = \sigma_k(W_k f^{k-1}(x) + b_k) \quad f^0(x) = x$$

$$ReLU(z) = [\max(z, 0)]_+$$

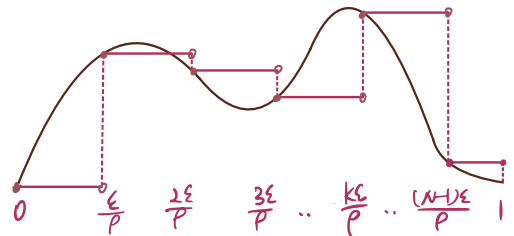
$$LTU(z) = [\mathbb{1}(z > 0)]_+$$

Thm. Let $g: [0, 1]^d \rightarrow \mathbb{R}$ be L -Lipschitz. For any $\epsilon > 0$, \exists a two-layer net f with one hidden layer with $N = \lceil \frac{L}{\epsilon} \rceil$ LTU units and a linear output unit

$$\text{s.t. } \|f - g\|_\infty \leq \epsilon$$

Pf. For $i \in \{0, \dots, N-1\}$, let $b_i = \frac{i\epsilon}{L}$
 $b_0 = 0, b_1 = \frac{\epsilon}{L}, \dots, b_{N-1} = (\lceil \frac{L}{\epsilon} \rceil - 1) \frac{\epsilon}{L} < 1$

$$f(x) = \begin{cases} g(0) & 0 \leq x < b_1 \\ g(b_1) & b_1 \leq x < b_2 \\ \vdots & \vdots \\ g(b_{N-1}) & b_{N-1} \leq x \leq 1 \end{cases} \quad \text{piece-wise constant approximation}$$



$$f(x) = \sum_{k=0}^{N-1} a_k \mathbb{1}(x \geq b_k) \quad a_0 = g(0), a_1 = g(b_1) - a_0, \dots, \sum_{i=0}^k a_i = g(b_k), a_i = g(b_i) - g(b_{i-1})$$

For any input x , let $k = \max \{i : b_i \leq x\}$, then

$$|g(x) - f(x)| \leq |g(x) - g(b_k)| + \underbrace{|g(b_k) - f(b_k)|}_0 + \underbrace{|f(b_k) - f(x)|}_0$$

$$\leq p|x - b_k| \leq p \frac{\epsilon}{p} = \epsilon$$

another way to show universal

Thm. (Stone-Weierstrass) if \mathcal{F} of functions $X \rightarrow \mathbb{R}$ satisfies:

1. $\mathcal{F} \subseteq C(X)$ continuous

2. $\forall x \in X, \exists f \in \mathcal{F}$ with $f(x) \neq 0$

3. \mathcal{F} is an algebra: $\forall f, g \in \mathcal{F}, \forall \alpha \in \mathbb{R} \left\{ \begin{array}{l} \alpha f + g \in \mathcal{F} \\ fg = (x \mapsto f(x)g(x)) \in \mathcal{F} \end{array} \right.$

4. \mathcal{F} separates points: $\forall x, x' \in X$ with $x \neq x', \exists f \in \mathcal{F}$ with $f(x) \neq f(x')$

Then \mathcal{F} is dense in $C(X)$

$$\mathcal{F}_{\text{exp}} = \left\{ x \mapsto \sum_{i=1}^N \alpha_i \exp(w_i \cdot x) : N \geq 1, w_1, \dots, w_N \in \mathbb{R}^d, \alpha_1, \dots, \alpha_N \in \mathbb{R} \right\}$$

\mathcal{F}_{exp} is dense in $C(X)$

\checkmark continuous $\checkmark \exists f(x) \neq 0$ ($x \mapsto 1$) $\in \mathcal{F}$ $\checkmark \alpha f + g \in \mathcal{F}$

$\checkmark fg = \sum_{i=1}^N \sum_{j=1}^M \alpha_i \beta_j \exp((w_i + v_j) \cdot x) \in \mathcal{F}$

x_1 vs x_2 : use $f(x) = \exp((x_1 - x_2) \cdot x)$

$$\frac{f(x_1)}{f(x_2)} = \frac{\exp(\|x_1\|^2 - x_1 \cdot x_2)}{\exp(x_1 \cdot x_2 - \|x_2\|^2)} = \exp(\|x_1\|^2 + \|x_2\|^2 - 2x_1 \cdot x_2) = \exp(\|x_1 - x_2\|^2) \neq 1$$

MLP approximation: $\mathcal{F}_\sigma = \left\{ x \mapsto \sum_{i=1}^m \alpha_i \sigma(w_i \cdot x) : m \geq 1; w_1, \dots, w_m \in \mathbb{R}^d, \alpha_1, \dots, \alpha_m \in \mathbb{R} \right\}$

$\sigma: \mathbb{R} \rightarrow \mathbb{R}$ is continuous with $\lim_{z \rightarrow -\infty} \sigma(z) = 0, \lim_{z \rightarrow +\infty} \sigma(z) = 1$

⊙ use $f_0 \in \mathcal{F}_{\text{exp}}$ to approximate g s.t. $\|f_0 - g\|_\infty \leq \frac{\epsilon}{2}$

⊕ use linear combination of σ s.t. $\exp(z) \approx \sum_j c_j \sigma(t_j z)$

then replace each $\exp(w_i \cdot x)$ in f_0 by $\sum_j c_j \sigma(t_j w_i \cdot x)$ to find

$f \in \mathcal{F}_\sigma$ s.t. $\|f - f_0\| \leq \frac{\epsilon}{2}$

this also holds if σ is anything but polynomial

because linear combination of polynomials of degree d is also degree d .

can't approximate any function (needs to be arbitrarily high to do this!)

Circuit complexity

Two-layer network, threshold activations,
can represent all $g: \{\pm 1\}^d \rightarrow \pm 1$

- might take exponential width to do it
- if g can be computed in time T , \exists a network of size $O(T^2)$ that implements g
there exists doesn't mean we can find it

Is ERM enough?

- ERM is NP-hard to compute for NN
 $f(x) = \text{ReLU}(w \cdot x)$ with square loss is NP-hard
- Uniform convergence ERM bounds might not be good enough

VC dimension param counting bound

P : params D : depth ReLU (piecewise-linear) net

$$\text{VCdim} = O(PD \log P), \Omega(PD \log \frac{1}{D})$$

big O : less than equal
big Ω : greater than equal
big Θ : equal

in practice, more params improve generalization,

but worse bounds

DNN has lots of params!

norm-based bounds

Rademacher complexity, covering number ... tend to be vacuous

- ERM might not generalize well

Lec Stability, Regularization, Convex problems

Regularized loss minimization:

$$\operatorname{argmin}_{h \in \mathcal{H}} L_S(h) + \lambda R(h)$$

eg

kernel ridge regression

$$\operatorname{argmin}_{h \in \mathcal{F}} L_S^{sq}(h) + \frac{1}{2} \lambda \|h\|_{\mathcal{F}}^2$$

soft SVM

$$\operatorname{argmin}_{h \in \mathcal{F}} L_S^{\text{hinge}}(h) + \frac{1}{2} \lambda \|h\|_{\mathcal{F}}^2$$

Lagrange dual

$$\operatorname{argmin}_{h: R(h) \leq \frac{1}{2} B^2} L_S(h)$$

$$\operatorname{argmin}_{h: \|h\|_{\mathcal{F}} \leq B} L_S^{sq}(h)$$

$$\operatorname{argmin}_{h: \|h\|_{\mathcal{F}} \leq B} L_S^{\text{hinge}}(h)$$

vice versa

equivalent in the sense that for any λ there is some B s.t. they are the same
easier to find some λ

motivation: the failure of boundingERM with uniform convergence

we want a "stable" algorithm; If $S \approx S'$, $A(S) \approx A(S')$

the algorithm is not very sensitive to the small change in samples
implies a small \mathcal{H} !

in RLM., $R(h)$ could be seen as a stabilizer

$$S = (z_1, \dots, z_m) \quad S_{\text{change } z_i}^{(i \leftarrow z')} = (z_1, \dots, z_{i-1}, z', z_{i+1}, \dots, z_m)$$

Theorem $\mathbb{E}_{S \sim \mathcal{D}^m, A} [L_D(A(S)) - L_S(A(S))] = \mathbb{E}_{S \sim \mathcal{D}^m} \underbrace{[l(A(S^{i \leftarrow z'}), z_i)]}_{\text{generalization error}} - \underbrace{[l(A(S), z_i)]}_{\text{training error}}$

A could be random but independent of data
eg. SGD

$S \sim \mathcal{D}^m, z_i \sim \mathcal{D}$
 $i \sim \text{Uniform}(\{1, \dots, m\}), A$

Def (stable) A is $\epsilon(m)$ -on-average-replace-one stable if for all \mathcal{D} .

$$\mathbb{E}_{S \sim \mathcal{D}^m, z \sim \mathcal{D}, i \sim \text{unif}(\{1, \dots, m\}), A} [l(A(S^{i \leftarrow z}), z_i) - l(A(S), z_i)] \leq \epsilon(m)$$

average-case generalization gap

Def (uniformly stable) A is $\beta(m)$ -uniformly stable if for all $m \geq 1$

stronger notion $\sup_{S \in \mathcal{D}^m} \sup_{z, z' \in \mathcal{D}} \left| \mathbb{E}_A \ell(A(S^{i \leftarrow z'}), z) - \mathbb{E}_A \ell(A(S), z) \right| \leq \beta(m)$

Thm. (Uniform stability) $\ell \in [a, b]$, A is $\beta(m)$ -uniformly stable, with Prob $\geq 1 - \delta$ over $S \sim \mathcal{D}^m$ $\mathbb{E}_A [L_D(A(S)) - L_S(A(S))] \leq \beta(m) + (2m\beta(m) + b - a) \sqrt{\frac{1}{2m} \log \frac{1}{\delta}}$

if A is $\beta(m)$ -uniformly stable, it's average-replace-one stable

let $f(S) = \mathbb{E}_A [L_D(A(S)) - L_S(A(S))] \leq \mathbb{E}_{z \sim \mathcal{P}, i \sim \text{unif}[m], A} [\ell(A(S^{i \leftarrow z'}), z) - \ell(A(S), z)]$

If A is $\beta(m)$ -uniformly stable. $\mathbb{E}_S f(S) \leq \beta(m)$

Then we can use McDiarmid's inequality if $f(S)$ has bounded difference $|f(S) - f(S^{i \leftarrow z'})| = \left| \mathbb{E}_A [L_D(\hat{h}) - L_S(\hat{h}) - L_D(\hat{h}^i) + L_S(\hat{h}^i)] \right| \leq \left| \mathbb{E}_A [L_D(\hat{h}) - L_D(\hat{h}^i)] \right| + \left| \mathbb{E}_A [L_S(\hat{h}) - L_S(\hat{h}^i)] \right|$

Assume $\ell \in [a, b]$, let $\hat{h} = A(S)$, $\hat{h}^i = A(S^{i \leftarrow z'})$

$$\begin{aligned} \textcircled{1} \quad \left| \mathbb{E}_A L_D(\hat{h}^i) - \mathbb{E}_A L_D(\hat{h}) \right| &= \left| \mathbb{E}_{A, z} [\ell(\hat{h}^i, z) - \ell(\hat{h}, z)] \right| \\ &\leq \mathbb{E}_{z \sim \mathcal{D}} \left| \mathbb{E}_A \ell(\hat{h}^i, z) - \mathbb{E}_A \ell(\hat{h}, z) \right| \\ &\leq \beta(m) \end{aligned}$$

$$\begin{aligned} \textcircled{2} \quad \left| \mathbb{E}_A L_S(\hat{h}) - \mathbb{E}_A L_{S^{i \leftarrow z'}}(\hat{h}^i) \right| &\leq \frac{1}{m} \sum_{j \neq i} \left| \mathbb{E}_A \ell(\hat{h}, z_j) - \mathbb{E}_A \ell(\hat{h}^i, z_j) \right| + \frac{1}{m} \left| \mathbb{E}_A \ell(\hat{h}, z_i) - \mathbb{E}_A \ell(\hat{h}^i, z') \right| \\ &\stackrel{\frac{m-1}{m} \leq 1}{\leq} \beta(m) + \frac{b-a}{m} \quad \ell \text{ is bounded} \end{aligned}$$

Combine $\textcircled{1} + \textcircled{2}$

f has bounded diffs with rate $C_i = 2\beta(m) + \frac{b-a}{m}$

\Rightarrow with prob at least $1 - \delta$ over $S \sim \mathcal{D}^m$

$$f(S) \leq \mathbb{E} f(S) + C \sqrt{\frac{m}{2} \log \frac{1}{\delta}} \leq \beta(m) + (2m\beta(m) + b - a) \sqrt{\frac{1}{2m} \log \frac{1}{\delta}}$$

if $\beta(m)$ decays with $\frac{1}{m}$, then $f(S)$ decays $\frac{1}{\sqrt{m}}$

Convex Function

Def (convex set) set $C \subseteq X$ is convex if for all $x_0, x_1 \in C$ and $\alpha \in [0, 1]$

$$x_\alpha = (1 - \alpha)x_0 + \alpha x_1 \in C$$

instead define a restricted domain, we define $f(x) = \infty$ for x out of domain. $\text{dom } f = \{x \in X : f(x) < \infty\}$

Def (convex function). A function $f : X \rightarrow \mathbb{R} \cup \{\infty\}$ is called

- convex iff $\forall x_0, x_1 \in X, \alpha \in (0, 1) \quad f((1-\alpha)x_0 + \alpha x_1) \leq (1-\alpha)f(x_0) + \alpha f(x_1)$
- m -strongly convex, $\dots \quad f((1-\alpha)x_0 + \alpha x_1) \leq (1-\alpha)f(x_0) + \alpha f(x_1) - \frac{1}{2}m\alpha(1-\alpha)\|x_1 - x_0\|^2$
- strictly convex, $\dots \quad f((1-\alpha)x_0 + \alpha x_1) < (1-\alpha)f(x_0) + \alpha f(x_1)$

If f is differentiable, f is

- convex, iff $\forall x, x', \quad f(x') \geq f(x) + \langle \nabla f(x), x' - x \rangle$
- m -strongly convex, $\quad f(x') \geq f(x) + \langle \nabla f(x), x' - x \rangle + \frac{1}{2}m\|x' - x\|^2$

If f is continuously differentiable, f is

- convex, iff $\forall x, x', \quad \langle \nabla f(x) - \nabla f(x'), x - x' \rangle \geq 0$
- m -strongly convex, $\quad \langle \nabla f(x) - \nabla f(x'), x - x' \rangle \geq m\|x - x'\|^2$

If f is continuously twice-differentiable, f is

- convex, iff $\forall x, x', \quad \nabla^2 f(x) \succeq 0$ all eigs $\geq m$
- m -strongly convex, $\quad \nabla^2 f(x) \succeq mI$

$$\begin{aligned} A \succeq 0 &\Leftrightarrow A \text{ ps.d} \\ A \succeq B &\Leftrightarrow A - B \succeq 0 \end{aligned}$$

RLM is often uniformly stable under some conditions

$$f_s(h) = \underbrace{L_s(h)}_{\text{convex}} + \underbrace{\lambda R(h)}_{\text{strongly convex (set to 1-strongly convex here)}}$$

$$\begin{aligned} \underbrace{f_s(h) - f_s(g)}_{\text{}} &= \underbrace{L_s(h) - L_s(g)}_{\text{}} + \underbrace{\lambda R(h) - \lambda R(g)}_{\text{}} \quad s' = s^{i \leftarrow z'} \\ &= \underbrace{f_s'(h) - f_s'(g)}_{\text{}} + \frac{1}{m}(\underbrace{l(h, z') - l(g, z')}_{\text{}}) + \frac{1}{m}(\underbrace{l(g, z') - l(h, z')}_{\text{}}) \end{aligned}$$

let \hat{h}_i minimize $f_{s^{i \leftarrow z'}}$ and \hat{h} minimize f_s , then $f_s'(\hat{h}_i) - f_s'(\hat{h}) \leq 0$

$$\hat{h}_i = A(s^{i \leftarrow z'}) \in \arg \min_h f_{s^{i \leftarrow z'}}(h) \quad \hat{h} = A(s)$$

$$\circ f_S(\hat{h}^i) - f_S(\hat{h}) \leq \frac{1}{m} (l(\hat{h}^i, z_i) - l(\hat{h}, z_i)) + \frac{1}{m} (l(\hat{h}, z'_i) - l(\hat{h}^i, z'_i)) + \underbrace{f_S(\hat{h}^i) - f_S(\hat{h})}_{\leq 0}$$

$\circ \nabla f_S(\hat{h}) = 0$ as f_S is convex \circ

$$\circ f_S(\hat{h}^i) - f_S(\hat{h}) \geq \langle \nabla f_S(\hat{h}), \hat{h}^i - \hat{h} \rangle + \frac{1}{2} \lambda \|\hat{h}^i - \hat{h}\|^2$$

$\circ \text{+} \circ$ gives

$$\frac{1}{2} \lambda \|\hat{h}^i - \hat{h}\|^2 \leq \frac{1}{m} (l(\hat{h}^i, z_i) - l(\hat{h}, z_i)) + \frac{1}{m} (l(\hat{h}, z'_i) - l(\hat{h}^i, z'_i))$$

Assume $\forall z, h \mapsto l(h, z)$ is ρ -Lipschitz, then $\frac{1}{2} \lambda \|\hat{h}^i - \hat{h}\|^2 \leq \frac{2\rho}{m} \|\hat{h} - \hat{h}^i\|$

$$\Rightarrow \|\hat{h} - \hat{h}^i\| \leq \frac{4\rho}{m\lambda}$$

use Lips twice

$$\Rightarrow |l(\hat{h}, z) - l(\hat{h}^i, z)| \leq \frac{4\rho^2}{m\lambda}$$

convex + Lipschitz RLM is $\frac{4\rho^2}{m\lambda}$ -uniformly stable

Assume R is non-negative

$$\mathbb{E}_S L_S(A(S)) \leq \mathbb{E}_S L_S(A(S)) + \frac{4\rho^2}{\lambda m} \leq L_D(h^*) + \lambda R(h^*) + \frac{4\rho^2}{\lambda m}$$

$$L_S(A(S)) \leq L_S(A(S)) + \lambda R(A(S)) \leq L_S(h^*) + \lambda R(h^*) \quad \forall h^* \in \mathcal{H}$$

$$\mathbb{E}_S L_S(A(S)) \leq L_D(h^*) + \lambda R(h^*)$$

Assume $R(h^*) \leq \frac{1}{2} B^2$

how to tradeoff λ to get good model?

$$\mathbb{E}_S L_D(A(S)) \leq \inf_{h: R(h) \leq \frac{1}{2} B^2} L_D(h) + \frac{1}{2} \lambda B^2 + \frac{4\rho^2}{\lambda m}$$

$$ax + \frac{b}{x} \geq 2\sqrt{ab} \quad (x > 0)$$

minimize the bound by opt λ

$$\text{if } \lambda = \sqrt{\frac{2 \times 4\rho^2}{m B^2}} = \frac{\rho}{B} \sqrt{\frac{8}{m}} \Leftrightarrow B = \frac{\rho}{\lambda} \sqrt{\frac{8}{m}}$$

$$\mathbb{E}_S L_D(A(S)) \leq \inf_{h: R(h) \leq \frac{1}{2} B^2} L_D(h) + B\rho \sqrt{\frac{8}{m}}$$

PAC learning

How to choose λ in practice?

$$\mathbb{E}_S L_D(A(S)) \leq \inf_{h: R(h) \leq \frac{1}{2} \left(\frac{\rho}{\lambda} \sqrt{\frac{8}{m}}\right)^2} L_D(h) + \frac{8\rho^2}{\lambda m}$$

constant λ can converge

but B shrinks, can't compare

every thing in \mathcal{H} .

λ can be $\frac{1}{\sqrt{m}}$, or we let $\left\{ \begin{array}{l} \frac{1}{\lambda \sqrt{m}} \rightarrow \infty \\ \frac{1}{\lambda m} \rightarrow 0 \end{array} \right.$

$$\Rightarrow \lambda \propto \frac{1}{m^\gamma} \quad \gamma \in \left(\frac{1}{2}, 1\right)$$

Lec Gradient Descent

GD tries to find $\min_w f(w)$. start at w_1 , then update

$$w_{t+1} = w_t - \eta_t \nabla f(w_t) \quad f: W \rightarrow \mathbb{R}, W \subseteq \mathbb{R}^d$$

repeat T steps then return

last iterate = w_T

average iterate = $\frac{1}{T} \sum_{i=1}^T w_i$

best iterate : $w_{\hat{t}} = \hat{t} \in \operatorname{argmin}_{t \in \{1, \dots, T\}} f(w_t)$, $w \in \{w_1, \dots, w_T\}$

Assume η_t independent of data

what if w is optimized in constrained space?

Projected Gradient Descent

define $\operatorname{proj}_W(w) \in \operatorname{argmin}_{w' \in W} \|w - w'\|$ find the closest point in the set

Thm. If W is closed and convex, $\forall v \in W, \|\operatorname{proj}_W(w) - v\| \leq \|w - v\|$



what if the function isn't differentiable?

subgradient descent

subgradient: slope of planes that lower bound f

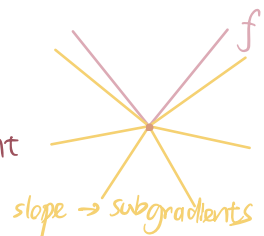
convex function \Leftrightarrow has subgradient over every point

$$\text{convexity: } f(w') \geq f(w) + \langle g, w' - w \rangle$$

Def (subgradient). g is a subgradient of f at w if

$$\forall w', f(w') \geq f(w) + \langle g, w' - w \rangle$$

The subdifferential of f at w , $\partial f(w)$, is the set of all subgradients of f at w



Prop 1. If f is convex and differentiable at w $\partial f(w) = \{\nabla f(w)\}$

Prop 2. f is convex iff $\forall w, \partial f(w)$ is non-empty

Prop 3. If f is ρ -Lipschitz, $\forall w, \forall g \in \partial f(w), \|g\| \leq \rho$

Pf. $\forall w$ in the interior of domain of f , let $g \in \partial f(w)$

$$\rho \varepsilon \geq f(w + \varepsilon \frac{g}{\|g\|}) - f(w) \geq \langle g, \varepsilon \frac{g}{\|g\|} \rangle = \frac{\varepsilon \langle g, g \rangle}{\|g\|} = \varepsilon \|g\|$$

Prop 4. Let $f(w) = \max_{i \in X} f_i(w)$, $|X| < \infty$, each f_i is convex.

$\forall w$, if $j \in \operatorname{argmax}_i f_i(w)$, then $\partial f_j(w) \subseteq \partial f(w)$

Pf. $g \in \partial f_j(w), f(w') \geq f_j(w') \geq f_j(w) + \langle g, w' - w \rangle$
 $= f(w) + \langle g, w' - w \rangle$

Stochastic (projected) (sub)gradient descent

Lemma. Let v_1, \dots, v_T be an arbitrary sequence, $w_{t+1} = \operatorname{proj}_W(w_t - \eta v_t)$

Then $\sum_{t=1}^T \langle w_t - w^*, v_t \rangle \leq \frac{1}{2\eta} \|w_1 - w^*\|^2 + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2$

Pf. $\langle w_t - w^*, v_t \rangle = \frac{1}{\eta} \langle w_t - w^*, \eta v_t \rangle$ $2\langle a, b \rangle = \|a\|^2 + \|b\|^2 - \|a-b\|^2$

$$= \frac{1}{2\eta} (-\|w_t - \eta v_t - w^*\|_{\operatorname{proj}}^2 + \|w_t - w^*\|^2 + \eta^2 \|v_t\|^2)$$

$$\leq \frac{1}{2\eta} (-\|w_{t+1} - w^*\|_{d_{t+1}}^2 + \|w_t - w^*\|_{d_t}^2 + \eta^2 \|v_t\|^2)$$

$$\sum_{t=1}^T \langle w_t - w^*, v_t \rangle \leq \frac{1}{2\eta} (-\underbrace{d_{T+1}}_{< 0} + d_1) + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2$$

$$\leq \frac{1}{2\eta} \|w_1 - w^*\|^2 + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2$$

projected GD bound

To bound $f(\bar{w}) - f(w^*)$ any $w^* \in W$

$w^* \in W$; f is convex and ρ -Lipschitz

$$f(\bar{w}) - f(w^*) = f\left(\frac{1}{T} \sum_{t=1}^T w_t\right) - f(w^*)$$

$$\leq \frac{1}{T} \sum_{t=1}^T (f(w_t) - f(w^*))$$

$$g_t \in \partial f(w_t), f(w_t) - f(w^*) \leq \langle g_t, w_t - w^* \rangle$$

$$\leq \frac{1}{T} \sum_{t=1}^T \langle g_t, w_t - w^* \rangle$$

use Lemma \nearrow

if $\|w_1 - w^*\| \leq B$ $\leq \frac{1}{2\eta T} \|w_1 - w^*\|^2 + \frac{\eta}{2T} \sum_{t=1}^T \|g_t\|^2$

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

Jensen's inequality:

$$\text{convex } f \rightarrow f\left(\frac{1}{n} \sum_{i=1}^n x_i\right) \leq \frac{1}{n} \sum_{i=1}^n f(x_i)$$

If $\eta = \sqrt{\frac{B^2}{2T} \cdot \frac{2}{\rho^2}} = \frac{B}{\rho\sqrt{T}}$, $f(\bar{w}) - f(w^*) \leq \frac{B\rho}{\sqrt{T}}$ tradeoff learning rate η

SGD bound get an average direction of gradients

$f(w) = \frac{1}{m} \sum_{i=1}^m f_i(w) + \lambda R(w)$ $f_i(w) = l(hw, z_i)$

$f(w) = \mathbb{E}_{z \sim D} l(hw, z) = \mathcal{L}(hw)$ but \hat{g}_t may be dependent as $\hat{g}_1 \rightarrow w_2 \rightarrow \hat{g}_2$

Assume $\hat{g}_t | w_t$ are independent of each other use data only one pass

$w_{t+1} = \text{proj}_W(w_t - \eta \hat{g}_t)$ for $\mathbb{E}[\hat{g}_t | w_t] \in \partial f(w_t)$

$\mathbb{E}_z \nabla f_t(w) = \mathbb{E}_z \nabla l(hw, z_t) = \nabla \mathbb{E}_{z \sim D} l(hw, z)$ unbiased gradient estimator

$\hat{g}_{t:z} = (\hat{g}_t, \hat{g}_{t+1}, \dots, \hat{g}_T)$

$\mathbb{E}[f(\bar{w}) - f(w^*)] \leq \mathbb{E} \frac{1}{T} \sum_{t=1}^T f(w_t) - f(w^*) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{\hat{g}_{1:t}} f(w_t) - f(w^*)$

$\hat{g}_{1:T}$ randomness in selected data

$\mathbb{E}_{\hat{g}_{1:t-1}} [f(w_t) - f(w^*)] \leq \mathbb{E}_{\hat{g}_{1:t-1}} [\langle w_t - w^*, \mathbb{E}_{\hat{g}_t} [\hat{g}_t | \hat{g}_{1:t-1}] \rangle]$ $g_t = \mathbb{E}[\hat{g}_t | w_t] = \mathbb{E}[\hat{g}_t | \hat{g}_{t-1}]$

$= \mathbb{E}_{\hat{g}_{1:t-1}} [\mathbb{E}_{\hat{g}_t} \langle w_t - w^*, \hat{g}_t \rangle | \hat{g}_{1:t-1}]$
 $= \mathbb{E}_{\hat{g}_{1:t}} [\langle w_t - w^*, \hat{g}_t \rangle]$

$\mathbb{E}_{\hat{g}_{1:T}} [f(\bar{w}) - f(w^*)] \leq \mathbb{E}_{\hat{g}_{1:T}} [\frac{1}{T} \sum_{t=1}^T \langle w_t - w^*, \hat{g}_t \rangle]$
 $\leq \frac{1}{2\eta T} \|w_1 - w^*\| + \frac{\eta}{2T} \sum_{t=1}^T \mathbb{E} \|\hat{g}_t\|^2$

Lec Non-convex Optimization

2 layer NN:

$hw(x) = w_2 \cdot \sigma(w_1 x)$ $\sigma(t) = t$ linear model

$h-w(x) = -w_2 \sigma(-w_1 x) = hw(x)$

$L_S(hw) = \frac{1}{m} \sum_{i=1}^m (hw(x_i) - y_i)^2$

If convex, $l(h_0, z) \leq l(h_w, z)$ as $l(h_w, z) = l(h-w, z)$ doesn't hold!
 So deep linear model is not convex

Def (β -smooth). A function f is β -smooth if f is differentiable everywhere, and gradient ∇f is β -Lipschitz

prop. If f is twice-differentiable, f β -smooth iff $\nabla^2 f(x) \preceq \beta I$ for $\forall x$

prop. f is β -smooth. Then for any x and y in its domain

$$|f(y) - f(x) - \langle \nabla f(x), y-x \rangle| \leq \frac{1}{2} \beta \|y-x\|^2$$

pf. Let $x_\alpha = (1-\alpha)x_0 + \alpha x_1$, $g(\alpha) = f(x_\alpha)$, $g'(\alpha) = \langle \nabla f(x_\alpha), x_1 - x_0 \rangle$

$$\begin{aligned} f(x_1) - f(x_0) &= g(1) - g(0) = \int_0^1 g'(\alpha) d\alpha = \int_0^1 \langle \nabla f(x_\alpha), x_1 - x_0 \rangle d\alpha \\ &= \int_0^1 \langle \nabla f(x_\alpha) - \nabla f(x_0) + \nabla f(x_0), x_1 - x_0 \rangle d\alpha \\ &= \langle \nabla f(x_0), x_1 - x_0 \rangle + \int_0^1 \langle \nabla f(x_\alpha) - \nabla f(x_0), x_1 - x_0 \rangle d\alpha \end{aligned}$$

$$|f(x_1) - f(x_0) - \langle \nabla f(x_0), x_1 - x_0 \rangle| \leq \int_0^1 |\langle \nabla f(x_\alpha) - \nabla f(x_0), x_1 - x_0 \rangle| d\alpha$$

$$\leq \int_0^1 \|\nabla f(x_\alpha) - \nabla f(x_0)\| \|x_1 - x_0\| d\alpha$$

$$\leq \int_0^1 \beta \|x_\alpha - x_0\| \|x_1 - x_0\| d\alpha$$

$$= \int_0^1 \alpha \beta \|x_1 - x_0\|^2 d\alpha$$

$$= \frac{1}{2} \beta \|x_1 - x_0\|^2$$

$$\|x_\alpha - x_0\| = \|\alpha(x_1 - x_0)\|$$

Descent Lemma: If ∇f is β -Lipschitz (or f β -smooth), $x_{t+1} = x_t - \eta \nabla f(x_t)$

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

$$= f(x_t) - \eta \langle \nabla f(x_t), \nabla f(x_t) \rangle + \frac{1}{2} \beta \eta^2 \|\nabla f(x_t)\|^2$$

$$= f(x_t) - \eta (1 - \frac{1}{2} \beta \eta) \|\nabla f(x_t)\|^2$$

need $1 - \frac{1}{2} \beta \eta > 0$ to minimize f , $\eta < \frac{2}{\beta}$

do GD on DNN will get to stationary point. ERM like

positive homogenous:

$$\text{ReLU } \sigma \quad W_3 \sigma(W_2 \sigma(W_1 x))$$

$\times \frac{1}{2} \quad \times 2$

Lec Neural Tangent Kernel

two layer NN

$$h(x; w) = \frac{1}{\sqrt{N}} \sum_{j=1}^N \underbrace{a_j}_{\text{fixed second layer}} \sigma(w_j \cdot x) \quad w^{t=1} \sim \mathcal{N}(0, Id)$$

$$\langle A, B \rangle_{\mathcal{F}} = \sum_{ij} A_{ij} B_{ij}$$

$$h(x; w) \approx h_{\tilde{w}}(x; w) = h(x; \tilde{w}) + \langle \nabla_w h(x; \tilde{w}), w - \tilde{w} \rangle_{\mathcal{F}} \quad \text{Taylor approximation}$$

$$\text{good approximation} \quad = \frac{1}{\sqrt{N}} \sum_{j=1}^N a_j [\sigma(\tilde{w}_j \cdot x) + \sigma'(\tilde{w}_j \cdot x) \cdot x \cdot (w_j - \tilde{w}_j)]$$

if w, \tilde{w} are near

$$= \frac{1}{\sqrt{N}} \sum_{j=1}^N a_j [\underbrace{\sigma(\tilde{w}_j \cdot x) - \sigma'(\tilde{w}_j \cdot x) \cdot x \cdot \tilde{w}_j}_{\text{doesn't depend on } w \text{ constant term}} + \underbrace{\sigma'(\tilde{w}_j \cdot x) \cdot x \cdot w_j}_{\text{linear in } w}]$$

For ReLU: $t > 0 \Rightarrow t - t = 0$ $t < 0 \Rightarrow 0 - 0 = 0$

$$h_{\tilde{w}}(x; W) = \langle \underbrace{\nabla_w h(x; \tilde{w})}_{\varphi_{\tilde{w}}(x)}, W \rangle : \{x \mapsto h_{\tilde{w}}(x; w) : w \in \mathbb{R}^{N \times d}\}$$

is RKHS

$$K_{\tilde{w}}(x, x') = \langle \nabla_w h(x; \tilde{w}), \nabla_w h(x'; \tilde{w}) \rangle$$

For general activations

$$h_{\tilde{w}}(x; w) - h(x; \tilde{w}) = \langle \nabla_w h(x; \tilde{w}), w - \tilde{w} \rangle_{\mathcal{F}} \quad \text{is in RKHS}$$

so how good is the linearization?

$$L_D(A(S)) - L^* = \underbrace{L_D(A(S)) - L_D(\text{ERM}_{h(S)})}_{\text{optimization error}} + \underbrace{L_D(\text{ERM}_{h(S)}) - \inf_{h^* \in \mathcal{H}} L_D(h^*)}_{\text{estimation error}} + \underbrace{\inf_{h^* \in \mathcal{H}} L_D(h^*) - L^*}_{\text{approximation error}}$$

$$|h(x; w) - h_{\tilde{w}}(x; w)| \leq \frac{1}{\sqrt{N}} \sum_{j=1}^N |a_j| |\sigma(w_j \cdot x) - \sigma(\tilde{w}_j \cdot x) - \sigma'(\tilde{w}_j \cdot x)(w_j \cdot x - \tilde{w}_j \cdot x)|$$

If σ is β -smooth assume σ β -smooth

$$|\sigma(t) - \sigma(\tilde{t}) - \sigma'(\tilde{t})(t - \tilde{t})| \leq \frac{1}{\sqrt{N}} \sum_{j=1}^N |a_j| \frac{1}{2} \beta \|w_j \cdot x - \tilde{w}_j \cdot x\|^2$$

$$\leq \frac{1}{2} \beta \|t - \tilde{t}\|^2$$

$$\leq \frac{\beta}{2\sqrt{N}} \left(\max_j |a_j| \right) \sum_j \|w_j - \tilde{w}_j\|^2 \|x\|^2$$

$$= \frac{\beta}{2\sqrt{N}} \left(\max_j |a_j| \right) \|x\|^2 \|W - \tilde{W}\|_{\mathcal{F}}^2 \quad \text{if } N \text{ is large enough}$$

$$\leq B \Rightarrow \|W - \tilde{W}\|_{\mathcal{F}} \leq \sqrt{\frac{2B}{\beta}} \frac{1}{\sqrt{\max_j |a_j|}} \frac{1}{\|x\|} N^{\frac{d}{2}} \quad \text{large } \|w - \tilde{w}\|_{\mathcal{F}}$$

gradient flow

$$\frac{dw_t}{dt} = -\eta \nabla L_S(x \mapsto h(x; w_t)) \quad \text{easier to analyze than GD}$$

square loss = $\frac{-2\eta}{m} \sum_{i=1}^m (h(x_i; w_t) - y_i) \nabla_w h(x_i; w_t)$

$$\begin{aligned} \frac{d}{dt} h(x_j; w_t) &= \langle \nabla h(x_j; w_t), \frac{dw_t}{dt} \rangle \\ &= \frac{-2\eta}{m} \sum_i (h(x_i; w_t) - y_i) \langle \nabla_w h(x_i; w_t), \nabla_w h(x_j; w_t) \rangle \end{aligned}$$

$$\frac{d}{dt} h|_{S_X}(t) = \frac{-2\eta}{m} (\text{kernel matrix mmm}) (h|_{S_X}(t) - S_Y)$$

suppose K_{w_t} is constant in t then this dynamics will be the same as kernel gradient descent for kernel regression.

$$\begin{aligned} L_S(h) &= \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i)^2 \\ &= \frac{1}{m} \sum_{i=1}^m (\langle h, \varphi(x_i) \rangle \langle \varphi(x_i), h \rangle - 2y_i \langle \varphi(x_i), h \rangle + y_i^2) \quad \text{define } [a \otimes b]b' = (b, b')a \\ &= \langle h, [\frac{1}{m} \sum_{i=1}^m \varphi(x_i) \otimes \varphi(x_i)] h \rangle - 2y_i \langle \frac{1}{m} \sum_{i=1}^m \varphi(x_i), h \rangle + \frac{1}{m} \sum_{i=1}^m y_i^2 \end{aligned}$$

$$\begin{aligned} \frac{dh_t}{dt} &= -\eta \nabla_n L_S(h) = \frac{-2\eta}{m} \sum_{i=1}^m [\varphi(x_i) \otimes \varphi(x_i)] h + \frac{2\eta}{m} \sum_{i=1}^m y_i \varphi(x_i) \\ &= \frac{-2\eta}{m} \sum_{i=1}^m (h(x_i) - y_i) \varphi(x_i) \end{aligned}$$

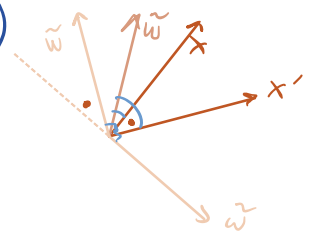
$$\begin{aligned} \langle \nabla_w h(x; \tilde{w}), \nabla_w h(x'; \tilde{w}) \rangle &= \frac{1}{N} \sum_{j=1}^N a_j^2 \langle x \sigma'(\tilde{w}_j \cdot x), x' \sigma'(\tilde{w}_j \cdot x') \rangle \\ &= x \cdot x' \frac{1}{N} \sum_{j=1}^N \sigma'(\tilde{w}_j \cdot x) \cdot \sigma'(\tilde{w}_j \cdot x') \end{aligned}$$

large number law; converge surely to, as $N \rightarrow \infty$

$$x \cdot x' \mathbb{E}_{\tilde{w} \sim \mathcal{N}(0, I)} [\sigma'(\tilde{w} \cdot x) \sigma'(\tilde{w} \cdot x')] = x \cdot x' \Pr\left(\frac{\tilde{w} \cdot x}{\|\tilde{w}\| \|x\|} > 0 \text{ and } \frac{\tilde{w} \cdot x'}{\|\tilde{w}\| \|x'\|} > 0\right) \quad \text{normalize}$$

$$\mathbb{E}[\mathbb{1}_{(a>0)}] = \Pr(a>0) = x \cdot x' \left(\frac{1}{2} - \frac{\arccos \frac{x \cdot x'}{\|x\| \|x'\|}}{2\pi}\right)$$

converge almost surely to some constant



empirical NTK

looking at the prediction on x , based on a single SGD update on x_t

$$h(x; w_{t+1}) - h(x; w_t) = \eta B(x) K_w(x, x_t) (y_t - h(x_t; w_t)) + O(\eta^2)$$

Lec Implicit Regularization

$f(w) = L_s^{\eta}(x \mapsto w \cdot x) = \frac{1}{m} \| \overset{m \times d}{X} w - \overset{m \times 1}{y} \|^2$
if $d > m$, have infinite many solutions

so what solution will we get when running GD? kernel GD flow is not actually used because η is not $\rightarrow 0$

$o f(w) = \frac{2}{m} X^T (Xw - y)$

$w_{t+1} = w_t - \frac{2\eta_0}{m} X^T X w_t + \frac{2\eta_0}{m} X^T y$ (set $\eta = \frac{2\eta_0}{m}$)

$= (I - \eta X^T X) w_t + \eta X^T y$

$= (I - \eta X^T X)^2 w_{t-1} + (I - \eta X^T X) \eta X^T y + \eta X^T y$

$= (I - \eta X^T X)^t w_1 + \sum_{k=0}^{t-1} (I - \eta X^T X)^k \eta X^T y$

$\stackrel{SVD}{=} (I - \eta V \Sigma^2 V^T)^t w_1 + \sum_{k=0}^{t-1} (I - \eta V \Sigma^2 V^T)^k \eta V \Sigma U^T y$

$\stackrel{\textcircled{4}}{=} (I - V V^T)^t w_1 + V (I - \eta \Sigma^2)^t V^T w_1 + \eta \sum_{k=0}^{t-1} ((I - V V^T) + V (I - \eta \Sigma^2)^k V^T) V \Sigma U^T y$

$= (I - V V^T)^t w_1 + V (I - \eta \Sigma^2)^t V^T w_1 + \eta V \sum_{k=0}^{t-1} (I - \eta \Sigma^2)^k \Sigma U^T y$

Assume $\eta < \frac{2}{\sigma_1^2}$ descending order, σ_i : largest singular value of X

$w_{\infty} = \underbrace{(I - V V^T)^t}_{\text{Neumann series}} w_1 + \underbrace{V \Sigma^{-1} U^T y}_{X^+ y} = (I - V V^T)^t w_1 + X^+ y$
pseudoinverse of X

$(I - V V^T)^{\infty} = I - 2V V^T + V V^T = I - V V^T$
there is a unique solution $X^+ y = (X^T X)^+ X^T y$

SVD singular value decomposition

If X is $m \times d$ of rank $r \leq \min(m, d)$, then Σ is $r \times r$ diagonal with $\Sigma_{ii} > 0$

U is $m \times r$, V is $d \times r$ with $U^T U = I_r = V^T V$

$X = \begin{matrix} U & \Sigma & V^T \\ m \times d & m \times r & r \times d \end{matrix}$

$\textcircled{1} X^T X = V \Sigma U^T U \Sigma V^T = V \Sigma^2 V^T$ (Σ^2)_{ii} is eigenvalue of $X^T X / X X^T$

$\textcircled{2} (V V^T)(V V^T) = V \underbrace{V^T V}_{I_r} V^T = V V^T$

$\textcircled{3} X(I - V V^T)y = U \Sigma V^T y - U \Sigma \overbrace{V^T V V^T}^I y = 0$

$\textcircled{4} (I - \eta X^T X)^k = (I - \eta V \Sigma^2 V^T)^k$
 $= (I - V V^T + \underbrace{V V^T - \eta V \Sigma^2 V^T}_{V(I - \eta \Sigma^2)V^T})^k$

$= (I - V V^T)^k + V (I - \eta \Sigma^2)^k V^T + \underbrace{V(I - \eta \Sigma^2)V^T}_0$