

Mathematical Tools for machine learning

Definition : (The realizability assumption) if $\exists f^* \text{ s.t. } L_{\text{exp}}(f^*) = 0$.

$$L_{\text{exp}}(f^*) = P(f^* \neq f)$$

$$L_S(f^*) = \text{empirical risk} \Rightarrow L_S(f^*) = L_S(f_{\text{EM}}^*)$$

A formal learning model (chapter 3)

- Waive the realizability assumption: Given a hypothesis class H , we do not assume anymore that $\exists h^* \in H : L_D(h^*) = P_{(x,y)}(h^*(x) \neq y) = 0$.

This leads to agnostic PAC learning:

Definition: A hypothesis class H is called agnostic PAC learnable if there exist a function $m_H : (\varepsilon, \delta)^2 \rightarrow \mathbb{N}$ and a learning algorithm st. $\forall (\varepsilon, \delta) \in (\varepsilon, \delta)^2$, \forall distribution D on $X \times Y$, if the algorithm is run on $m \geq m_H(\varepsilon, \delta)$ iid examples generated by D , this algorithm returns a predictor h_S , st. $\overline{\mathbb{P}_{S \sim D^m} (L_D(h_S) \leq \min_{h \in H} L_D(h) + \varepsilon)} \geq 1 - \delta$

\uparrow
training set

Remark: • $m_H(\varepsilon, \delta)$ is viewed as the smallest sample size for which the learning algorithm satisfied.

• In case the realizability assumption, $\min_{h \in H} L_D(h) = 0$.
 \Rightarrow agnostic PAC = PAC.

- Extending classification to other types of learning

- Multiclass classification: consider the problem of classifying some document into 4 categories

Domainset = $\mathcal{X} = \mathbb{N}^+$ with $p \gg$ some integer.

$\hookrightarrow x \in \mathcal{X}$ is a vector storing counts for some specific key words.

label set : $\mathcal{Y} = \{1, 2, \dots, k\}$ where $k \geq 2$ some integer.

Training set: $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$

A prediction rule h_s is the output of some learning algorithm for a new document with associate feature $x \in \mathcal{X}$. $h_s(x) = y \in \{1, -1\}$ the prediction class

- Regression: find a linear relationship between $y \in \mathbb{R}$ and covariates $x \in \mathbb{R}^p$ for some $p > 0$.

A linear model assumes that $E[y|x] = \beta^\top x$, β : unknown. the true risk (error) is

$$L_{\text{sq}}(h) = E_{(x,y) \sim P} [(h(x) - y)^2] : \text{expected squared error.}$$

The measure of error can be generalized to any loss function.

- Generalized loss function

Consider a hypothesis class H and some domain \mathcal{Z} (so far $= \mathcal{X} \times \mathcal{Y}$). Also consider a function $l: H \times \mathcal{Z} \rightarrow [0, \infty)$

such a function l is called a loss-function, for a given $h \in H$, the true error (risk) is now defined,

$$\text{as } L_{\text{sq}}(h) = E_{z \sim D} [l(h, z)].$$

the training error (empirical risk) is defined as

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m l(h, z_i) \quad \text{where } S = \{z_1, \dots, z_m\}$$

Typical loss functions are:

$$\begin{aligned} * \text{ 0-1 loss } l_{0-1}(h, z) &= l_{0-1}(h, (x, y)) = \begin{cases} 1 & \text{if } h(x) \neq y \\ 0 & \text{o.w.} \end{cases} \\ &= \mathbb{1}_{\{h(x) \neq y\}} \end{aligned}$$

$$* \text{ squared loss } l_{\text{sq}}(h, z) = (h(x) - y)^2$$

Definition (more generalized Agnostic PAC):

H is agnostic PAC to the domain \mathcal{Z} and loss function if (the same as before), where $L_P(h) = E_{z \sim D} [l(h, z)]$.

Learning via Uniform Convergence (Chapter 4)

Introduction: We know that under realizability assumption, finite classes

is PAC learnable.

Question: How can we show that finite classes are agnostic PAC learnable? and what assumption on the loss function?

Answer: Uniform convergence.

Definition (ε -representative sample)

A training set S is called ε -representative w.r.t. a domain Z , a hypothesis class H , a loss function l and a distribution D if there exists $|L_S(h) - L_D(h)| \leq \varepsilon$

Lemma: Assume that a training set S is $(\frac{\varepsilon}{2})$ -representative w.r.t. Z, H, l and D . Then for any output h_S of $\text{ERM}_H(S)$, that is $h_S \in \arg \min_{h \in H} L_S(h)$ we have that

$$L_D(h_S) \leq \min_{h \in H} L_D(h) + \varepsilon.$$

Proof: Since S is $(\frac{\varepsilon}{2})$ -representative, we have

$$\begin{aligned} L_D(h_S) &\leq L_S(h_S) + \frac{\varepsilon}{2} \\ &\leq L_S(h) + \frac{\varepsilon}{2} \quad \forall h \in H \text{ since } h \in \arg \min_{h \in H} L_S(h) \\ &\leq L_D(h) + \frac{\varepsilon}{2} + \frac{\varepsilon}{2} \\ &= L_D(h) + \varepsilon. \\ (\Rightarrow) L_D(h_S) &\leq \min_{h \in H} L_D(h) + \varepsilon. \end{aligned}$$

□

Definition (uniform convergence)

We say that a hypothesis class H has the uniform convergence property (w.r.t. Z and l) if there exists a function $m_H^{uc}: (0, 1)^2 \rightarrow N$ s.t. $\forall (\varepsilon, \delta) \in (0, 1)^2 \quad \forall$ distribution D if S is a training set with $m \geq m_H^{uc}(\varepsilon, \delta)$ example iid $\sim D$, then

$$\begin{aligned} \Pr_{S \sim D^m} (S \text{ is } \varepsilon\text{-representative}) &\geq 1 - \delta \\ \Leftrightarrow (\forall h \in H) |L_S(h) - L_D(h)| &\leq \varepsilon \end{aligned}$$

smallest one.

Corollary If a hypothesis class H has the uniform convergence property. (w.r.t. Z and l) with a function m_H^{uc} , then this function class is agnostic PAC learnable with sample complexity $m_H(\varepsilon, \delta) \leq m_H^{uc}(\frac{\varepsilon}{2}, \delta)$.

Furthermore, the ERM predictors is an agnostic PAC learner for H .

$$\text{proof: } \Pr_{\substack{\delta \\ \text{ERM}}} \left(L(h) \leq \min_{\mathcal{H}} L + \varepsilon \right) \geq \Pr \left(L(h_1) - L(h_2) + \dots + L(h_k) - L(h^*) \leq \frac{\varepsilon}{k} \right) = 1 - \delta.$$

Probability to get $\forall i$
 δ -representable $\Delta m_i > m_H^{uc}(\frac{\varepsilon}{k}, \delta)$

This means H is agnostic PAC learnable with the ERM paradigm with sample complexity $m_H(\varepsilon, \delta) \leq m_H^{uc}(\frac{\varepsilon}{k}, \delta)$. \square

- Finite classes are agnostic PAC learnable.

Lemma (Hoeffding's inequality): let $\theta_1, \dots, \theta_m$ be iid random variables s.t. $E[\theta_i] = \mu$, $i = \{1, \dots, m\}$ and $P(a \leq \theta_i \leq b) = 1$, then,

$$P \left(\left| \frac{1}{m} \sum_{i=1}^m (\theta_i - \mu) \right| \geq \varepsilon \right) \leq 2 \exp \left(- \frac{2m\varepsilon^2}{(b-a)^2} \right)$$

Auxiliary lemma: Let X be a random variable s.t. $E[X] = 0$ and $P(X \in [a, b]) = 1$. Then, $\forall \lambda > 0$ we have $E[e^{\lambda X}] \leq e^{\frac{\lambda^2(b-a)^2}{8}}$

Proof: Put $f(x) = e^{\lambda x}$, $x \in \mathbb{R}$ for $\lambda \in (0, \infty)$. f is convex on \mathbb{R} and hence $f(\alpha a + (1-\alpha)b) \leq \alpha f(a) + (1-\alpha)f(b)$. for $X \in [a, b]$ choose $a = \frac{bx}{b-a}$, $\alpha a + (1-\alpha)b = X \Rightarrow f(X) \leq \frac{bx}{b-a} f(a) + \frac{x-a}{b-a} f(b)$

$$\Leftrightarrow e^{\lambda X} \leq \frac{bx}{b-a} e^{\lambda a} + \frac{x-a}{b-a} e^{\lambda b}$$

$$\begin{aligned} \Rightarrow E[e^{\lambda X}] &\leq \frac{b - E[X]}{b-a} e^{\lambda a} + \frac{E[X] - a}{b-a} e^{\lambda b} \\ &= \frac{b}{b-a} e^{\lambda a} - \frac{a}{b-a} e^{\lambda b} \end{aligned}$$

$$h = \lambda(b-a) \text{ and } p = -\frac{a}{b-a}$$

Also, define the function

$$L(h) = -hp + \log(1-p+pe^h) \text{ with } h \in [0, \infty)$$

$$\begin{aligned} \text{we have } \frac{b}{b-a} e^{\lambda a} - \frac{a}{b-a} e^{\lambda b} &= e^{\lambda a} (1-p+pe^{\lambda(b-a)}) \\ &= \frac{\lambda a}{b-a} (b-a) e^{-hp} (1-p+pe^h) \\ &\quad - hp \\ &= e^{L(h)} \end{aligned}$$

In the following, we show that $L(h) \leq \frac{h^2}{8}$

$$L(0) = 0, L'(h) = -p + \frac{pe^h}{1-p+pe^h} \Rightarrow L'(0) = 0.$$

$$L'(h) = +P \cdot \frac{e^h (1-P+Pe^h)}{(1-P+Pe^h)^2} - Pe^h$$

$$= \frac{P(1-P)e^h}{(1-P+Pe^h)^2}$$

$$L''(h) - \frac{1}{4} = \frac{P(1-P)e^h}{(1-P+Pe^h)^2} - \frac{1}{4} = \frac{4P(1-P)e^h - (1-P+Pe^h)^2}{4(1-P+Pe^h)^2} = -\frac{(1-P+Pe^h)^2}{-11} \leq 0$$

$\forall h \in (0, +\infty)$, $L''(h) \leq \frac{1}{4}$. Using Taylor expansion of L up to the 2nd order we can write:

$$L(h) = L(0) + hL'(0) + \frac{h^2}{2}L''(h^*) \text{ where } h^* \in [0, h]$$

$$L(0) = L'(0) = 0$$

$$\Rightarrow L(h) = \frac{h^2}{2}L''(h^*) \leq \frac{h^2}{8} \quad \square$$

Now, we prove that the Hoeffding's inequality. Recall $\theta_1, \dots, \theta_m$, iid $\in [a, b]$
 $X_i = \theta_i - \mu$ where $\mu = \mathbb{E}[\theta_i]$ for $i \in \{1, \dots, m\}$.

Note that X_1, \dots, X_m are iid and $X_i \in [a-\mu, b-\mu]$. Using monotonicity of $x \mapsto e^{x\lambda}$ for $\lambda \in (0, +\infty)$ and the Markov's inequality, we know that $P(X > \varepsilon) = P(e^{X\lambda} > e^{\lambda\varepsilon}) \leq E[e^{X\lambda}] e^{-\lambda\varepsilon}$

where $\bar{X} = \frac{1}{m} \sum_{i=1}^m X_i$ (the empirical mean of the X 's)

$$E[e^{X\lambda}] = \prod_{i=1}^m E[e^{+\frac{\lambda}{m}X_i}] = (E[e^{\frac{\lambda}{m}X_1}])^m.$$

By the auxiliary result, we have

$$E[e^{\frac{\lambda}{m}X_1}] \leq e^{\frac{(\lambda(b-a))^2}{8m}}$$

$$\text{Hence } P(\bar{X} > \varepsilon) \leq e^{\frac{\lambda^2(b-a)^2}{8m}} - \lambda\varepsilon \quad \forall \lambda \in (0, +\infty)$$

$$\Leftrightarrow P(\bar{X} > \varepsilon) \leq \inf_{\lambda \in (0, +\infty)} e^{\frac{\lambda^2(b-a)^2}{8m}} - \lambda\varepsilon \Leftrightarrow \lambda^* = \frac{4\varepsilon m}{(b-a)^2} > 0.$$

$$\leq e^{-\frac{2m\varepsilon^2}{(b-a)^2}}$$

The same argument can be used to $-X_i$ to show that $P(-\bar{X} > \varepsilon) \leq e^{-\frac{2m\varepsilon^2}{(b-a)^2}}$

$$\Rightarrow P(|\bar{X}| > \varepsilon) \leq 2 \cdot e^{-\frac{2m\varepsilon^2}{(b-a)^2}}.$$

The proof is complete by noting that $\bar{X} = \frac{1}{m} \sum_{i=1}^m \theta_i - \mu$. \square

$$|\bar{X}|(w) = \frac{-X}{X}$$

Now we go back to the original problem: showing that any finite hypothesis class is Agnostic PAC learnable. In the following, we assume that the loss function, $\ell: \mathcal{H} \times \mathcal{Z} \rightarrow [0, 1]$, satisfies $\ell(h, z) \in \{0, 1\} \forall h, z \in \mathcal{H} \times \mathcal{Z}$.

If we fix $h \in \mathcal{H}$, $\ell = \ell(h, \cdot)$, $S = \text{training set } \{z_1, \dots, z_m\}$ are iid, write $\ell_i = \ell(h, z_i)$, $i=1 \dots m$

ℓ_1, \dots, ℓ_m are also iid with true mean $\mu = \mathbb{E}_{z \in \mathcal{Z}} [\ell(h, z)]$ (true risk associated to h)
 $L_S(h) \stackrel{\text{def}}{=} \frac{1}{m} \sum_{i=1}^m \ell(h, z_i) = \frac{1}{m} \sum_{i=1}^m \ell_i$ (training error)

Note that $\ell_1, \dots, \ell_m \in [0, 1]$ by assumption, using the Hoeffding's inequality $\mathbb{P}^m(S: |L_S(h) - \mu| > \varepsilon) \leq 2e^{-2m\varepsilon^2}$ (loss function bounded)

$$\Rightarrow \mathbb{P}^m(S: \exists h \in \mathcal{H} \text{ s.t. } |L_S(h) - \mu| > \varepsilon) \leq |\mathcal{H}| \times 2e^{-2m\varepsilon^2}$$

$$\text{we want } |\mathcal{H}| \times 2e^{-2m\varepsilon^2} \leq \frac{1}{\delta}. \text{ If it is true, then } \mathbb{P}^m(S: \forall h \in \mathcal{H} \mid |L_S(h) - \mu| \leq \varepsilon) \geq 1 - \delta$$

$$\leq m \geq \frac{-\log(\frac{2|\mathcal{H}|}{\delta})}{\varepsilon^2}$$

Corollary Let \mathcal{H} be a finite hypothesis class, \mathcal{Z} a domain and $\ell: \mathcal{H} \times \mathcal{Z} \rightarrow [0, 1]$ be a loss function. Then, \mathcal{H} enjoys the uniform convergence property with $m_H^{uc}(\varepsilon, \delta) \leq \frac{-\log(\frac{2|\mathcal{H}|}{\delta})}{\varepsilon^2}$.

furthermore, \mathcal{H} is Agnostic PAC learnable with the ERM paradigm.

The sample complexity satisfies

$$m_H(\varepsilon, \delta) \leq m_H^{uc}(\varepsilon, \delta) \leq \frac{2\log(\frac{2|\mathcal{H}|}{\delta})}{\varepsilon^2}$$

Question What if the class is infinite?

A possible answer "Discretization trick"

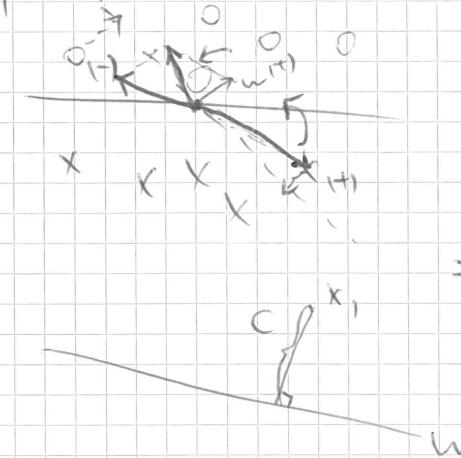
Example $X = \mathbb{R}$ and $y = \{+1\}$, $\mathcal{H} = \{h_\theta, \theta \in \mathbb{R}\}$, where $h_\theta(x) = \text{sign}(x - \theta)$, $x \in \mathbb{R}$, $|\mathcal{H}| = +\infty$

Since any real number can be represented by 64 bits, we can "approximate" \mathcal{H} by $\widehat{\mathcal{H}}$ where $|\widehat{\mathcal{H}}| = 2^{64}$. Since we are in the framework of binary classification, the loss function

$$\ell(h, z) = \begin{cases} 1 & h_\theta(x) \neq y \\ 0 & h_\theta(x) = y \end{cases} \in \{0, 1\}$$

Bayes estimator which is not consistent. for Jeffreys' Prior.

Remark:



Take bias into matrix \Rightarrow all have no bias passing through zero.

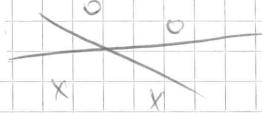
$$w^{(t+1)} = w^{(t)} + \gamma_i x_i \quad \text{if } y_i > 0 \quad \forall i.$$

$$\Rightarrow \bar{w} = \sum \alpha_i x_i$$

$\left\langle \frac{w}{\|w\|}, x_i \right\rangle = \text{distance to hyperplane } w \text{ from } x_i \geq \frac{1}{B}$

$\Rightarrow B$ large ~~more difficult to classify~~

\rightarrow more set B small



$\Rightarrow R$ large difficult to classify.

Chapter 5

1 proof of No-Free-Lunch theorem:

Theorem: Let A be any learning algorithm for the task of binary classification with respect to some domain X and the 0-1 loss function. Let $m < \frac{|X|}{2}$ represent the training set size \exists distribution D on $X \times \{0,1\}$ such that

1. $\exists f: X \rightarrow \{0,1\}$ such that $L_D(f) = 0$

2. $P_{S \sim D^m} (L_D(A(S)) \geq \frac{1}{2}) \geq \frac{1}{2}$

proof: Consider a subset $C \subset X$ such that $|C| = m$.

There are $T = 2^{|C|} = 2^m$ possible functions $C \rightarrow \{0,1\}$

Let us denote all these functions by f_1, \dots, f_T .

Fix $i \in \{1, \dots, T\}$. Consider D_i the distribution defined on $C \times \{0,1\}$ by

$$D_i(\{x,y\}) = \begin{cases} \frac{1}{T} & \text{if } y = f_i(x) \\ 0 & \text{otherwise} \end{cases} \quad x \in C.$$

D_i draws x from C uniformly (with the same probability $\frac{1}{m}$), and conditionally on x , assigns $f_i(x)$. As y with probability $\frac{1}{T}$

$$\text{The sample complexity } m_{\mathcal{H}}(\varepsilon, \delta) \leq \frac{\frac{2d}{\delta} \ln\left(\frac{2}{\delta}\right) + 1}{\varepsilon^2} \quad \varepsilon = 0.1 \text{ and } \delta = 0.01$$

$$m_{\mathcal{H}} \leq 13860.$$

Chapter 9 Linear Predictors

1 Introduction consider the following class $L_d = \{ h_{w,b} : w \in \mathbb{R}^d, b \in \mathbb{R} \}$

where $h_{w,b}(x) = \langle w, x \rangle + b$ for $x = (x_1, \dots, x_d)^T \in \mathbb{R}^d$

$$= w^T x + b$$



$$w = (w_1, \dots, w_d)^T = \sum_{i=1}^d w_i x_i + b$$

w: vector of weights . b: bias

Given a function $\phi: \mathbb{R} \rightarrow Y$ we can consider also the class

$$\phi \circ L_d = \{ \phi \circ h_{w,b} : w \in \mathbb{R}^d, b \in \mathbb{R} \}$$

Example: $X = \mathbb{R}^d, Y = \{-1, 1\}$ and $\phi(t) = \text{sign}(t)$

The corresponding class (used for binary classification) is

$$\text{sign} \circ L_d = \{ \text{sign}(\langle w, x \rangle + b) : w \in \mathbb{R}^d, b \in \mathbb{R} \}$$

Remark: The bias b can be incorporated in the vector of weights

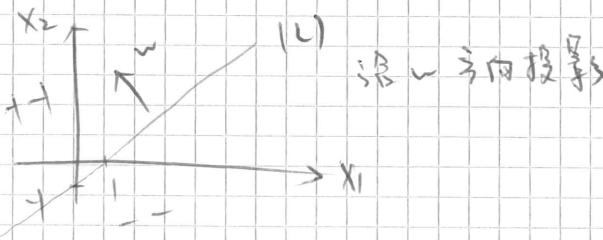
$$\langle w, x \rangle + b = (b, w_1, \dots, w_d) \cdot \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix}$$

2. Half-spaces

Consider again the binary classification problem with $X = \mathbb{R}^d, Y = \{-1, 1\}$

$$HS_d = \{ \text{sign}(\langle w, x \rangle + b) : w \in \mathbb{R}^d, b \in \mathbb{R} \}$$

Example: $d=2, w = (-1, 1)^T$ and $b=1$, $\langle w, x \rangle + b = -x_1 + x_2 + 1$



Remark: It can be shown that the class HS_d is VC of dimension $d+1$. Hence this class is Agnostic PAC learnable using ERM paradigm with sample complexity of order $\frac{d + \ln(1/\delta)}{\varepsilon^2}$

Sol 1:

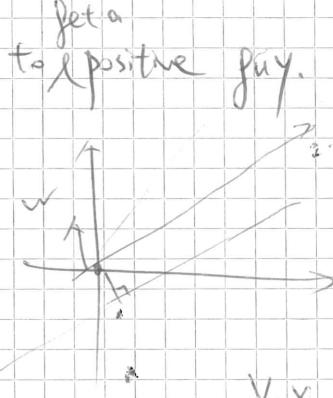
w and w^* both a perfect classifier. i.e. $y_i \langle w^T, x_i \rangle > 0$

Sol 2:

Rendomly take one update $\leq 0 \rightarrow$ add $\|x_i\|^2$ to positive guy.

In general, assume $y_i \langle w^T, x_i \rangle \geq 1$

It can stop in finite steps but B not known.



logistic regression

Interpretation: $P[Y=1 | x]$ hopes to close to 0 or 1.

$$\text{clearly } L_{D_i}(f_i) \stackrel{\text{def}}{=} P(x_i, y_i) \sim f_i(x_i) \neq y_i \\ = 0$$

We are going to show for any learning algorithm A which receives a training set S of m examples and outputs a classifier $A(S) : C \rightarrow \{0, 1\}$ it holds that $\max_{1 \leq i \leq T} \mathbb{E}_{S \sim D_i} [L_{D_i}(A(S))] \geq \frac{1}{4} (*)$.

构造的分布抽到它

There exists one distribution, the algorithm fails.

If $(*)$ is true then by Exercise 5-1 the result follows.

let us prove $(*)$

for $i \in \{1, \dots, T\}$, consider a training set S of m iid samples

$\sim D_i$. Such a training set S looks like

$$S = \{(x_1, f_i(x_1)), \dots, (x_m, f_i(x_m))\}$$

(x_1, \dots, x_m) is a random drawn with $|C|^m = (2m)^m$ possible outcomes which are all occurring with the same probability $\frac{1}{|C|^m}$. Put $k = |C|^m$

Also, let us list all possible training sets of sample $\sim D_i$ as

$$S_1, \dots, S_k.$$

$$\max_{1 \leq i \leq T} \mathbb{E}_{S \sim D_i} [L_{D_i}(A(S))] \geq \frac{1}{T} \sum_{i=1}^T \mathbb{E}_{S \sim D_i} [L_{D_i}(A(S))] \quad \text{finite.}$$

$\nexists j \times \text{fixed}$

Turn to a fixed j to estimate

$$= \frac{1}{k} \sum_{j=1}^k \frac{1}{T} \sum_{i=1}^T L_{D_i}(A(S_j))$$

station #1

$$\geq \min_{1 \leq i \leq k} \frac{1}{T} \sum_{i=1}^T L_{D_i}(A(S_j)) \quad A \text{ 值在所有 } x \text{ 组中} \\ \text{据 } D_i \text{ 分布的最小值.}$$

Now, fix $j \in \{1, \dots, k\}$. A training set S_j is of the form

$$S_j = \{(x_1^j, f_j(x_1^j)), \dots, (x_m^j, f_j(x_m^j))\} \quad (\text{# of } x_i^j \geq 1)$$

let $\{v_1, \dots, v_p\} = C \setminus \{x_1^j, \dots, x_m^j\}$ These are atmost m distinct values in $\{x_1^j, \dots, x_m^j\}$. Since $|C| = m + p \geq m$

let $h: C \rightarrow \{0, 1\}$

$$\begin{aligned} L_{S_j}(h) &\stackrel{\text{def}}{=} P_{(x,y) \sim S_j} (h(x) \neq y) \\ &= P_{x \sim \mathbb{P}_j^X} (h(x) \neq f_j(x)) \\ &\quad \text{marginal distribution of } x = \text{uniform on } C. \end{aligned}$$

$$\begin{aligned} &= \frac{1}{m} \sum_{c \in C} \mathbb{P}_{x \sim \mathbb{P}_j^X} \mathbb{1}_{\{h(c) \neq f_j(c)\}} \quad \uparrow \text{expectation of indicators.} \\ &\geq \frac{1}{m} \sum_{r=1}^p \mathbb{P}_{x \sim \mathbb{P}_j^X} \mathbb{1}_{\{h(v_r) \neq f_j(v_r)\}} \\ &\geq \frac{1}{p} \sum_{r=1}^p \mathbb{P}_{x \sim \mathbb{P}_j^X} \mathbb{1}_{\{h(v_r) \neq f_j(v_r)\}}. \end{aligned}$$

Replace h by $A(S_j)$, we get:

$$\begin{aligned} &\frac{1}{T} \sum_{i=1}^T \mathbb{1}_{\{A(S_j)(v_r) \neq f_i(v_r)\}} \\ &\leq \frac{1}{T} \sum_{i=1}^T L_{S_j}(A(S_j)) \\ &= \frac{1}{T} \sum_{r=1}^p \frac{1}{T} \sum_{i=1}^T \mathbb{1}_{\{A(S_j)(v_r) \neq f_i(v_r)\}} \end{aligned}$$

station #2 $\Rightarrow \min_{1 \leq r \leq p} \frac{1}{T} \sum_{i=1}^T \mathbb{1}_{\{A(S_j)(v_r) \neq f_i(v_r)\}}$

Fix $r \in \{1, \dots, p\}$. The classifier $\underline{f_i}$ can be constructed in the following way:

- label all elements in $C \setminus \{v_r\}$
- Assign either $\underline{f_i} \rightarrow 1$ or $\underline{f_i} \rightarrow 0$

Let us call $\hat{f}_i^{(0)}$ the first possible classifier and $\hat{f}_i^{(1)}$ the second one.

Note that $\forall c \in C \setminus \{v_r\}$ we have $\hat{f}_i^{(0)}(c) = \hat{f}_i^{(1)}(c) = f_i(c)$

Also $\{f_1, \dots, f_T\} = \{\hat{f}_1^{(0)}, \dots, \hat{f}_1^{(1)}, \hat{f}_2^{(0)}, \dots, \hat{f}_2^{(1)}, \dots, \hat{f}_T^{(0)}, \dots, \hat{f}_T^{(1)}\}$

$$\{(x_1^j, \hat{f}_i^{(0)}(x_1^j)), \dots, (x_m^j, \hat{f}_i^{(0)}(x_m^j))\} = \{(x_1^j, \hat{f}_i^{(1)}(x_1^j)), \dots, (x_m^j, \hat{f}_i^{(1)}(x_m^j))\}$$

v_r not in the training set $\{x_i^j\} \rightarrow \text{number}$

$$= S_j.$$

$$\sum_{i=1}^T \mathbb{I}\{A(s_i)(v_i) \neq f_i(v_i)\} = \sum_{\ell=1}^{T/2} \mathbb{I}\{A(s_i^\ell)(v_i) \neq \tilde{f}_\ell^{(i)}(v_i)\} + \sum_{\ell=1}^{T/2} \mathbb{I}\{A(s_i^\ell)(v_i) \neq \hat{f}_\ell^{(i)}(v_i)\}$$

$$= \sum_{\ell=1}^{T/2} [\mathbb{I}\{A(s_i^\ell)(v_i) \neq 0\} + \mathbb{I}\{A(s_i^\ell)(v_i) \neq 1\}] = 0$$

$$= \frac{T}{2} \text{ station } \#3 \leftarrow \text{always } \frac{T}{2}.$$

It follows that $\max_{1 \leq i \leq T} \mathbb{E}_{\text{supp}} [L_D(A(s_i))] \geq \frac{1}{2} \min_{1 \leq j \leq T} \min_{1 \leq i \leq p} \frac{1}{T} \cdot \frac{1}{2} = \frac{1}{4}$! \square

2. Error decomposition

let $h_S \in \arg \min_{h \in H} L_S(h)$ where H is some hypothesis class

$$L_S(h_S) = \min_{h \in H} L_S(h) + L_S(h_S) - \min_{h \in H} L_S(h)$$

$\min_{h \in H} L_S(h) = \varepsilon_{\text{app}} = \text{approximation error}$

$$L_S(h_S) - \min_{h \in H} L_S(h) = \varepsilon_{\text{est}} = \text{estimation error}$$

Approximation error: • not connected to any probabilistic statement
• if H is large then we expect that this error is small (could be zero in the realizability case).

Estimation error: Results from replacing the unknown distribution by a training set.

If H is big then we need more examples in the training set.

(e.g. If H is finite we know that $m_H(\varepsilon, \delta)$ of order $\frac{2}{\varepsilon^2} \log\left(\frac{2|H|}{\delta}\right)$)

This called the bias-complexity trade-off.

The VC-dimension (Chapter 6)

Question: Can we find infinite hypothesis classes which PAC learnable?

1. Introductory example

Consider the class of thresholds over \mathbb{R} : $H = \{h_a : a \in \mathbb{R}\}$
and $h_a(x) = \mathbb{I}\{x < a\}$

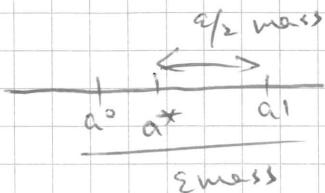
Lemma: H is PAC learnable with an ERM rule. The sample complexity is $m_H(\varepsilon, \delta) \leq \frac{2}{\varepsilon^2} \log\left(\frac{2}{\delta}\right)$ for $(\varepsilon, \delta) \in (0, 1)^2$

Proof: Let a^* be such that h_{a^*} is a perfect classifier, that is,

$$L_D(h^*) \stackrel{\text{def}}{=} P_{x \sim p^*} (h^*(x) \neq y) = 0$$

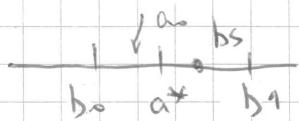
let D_x be the marginal distribution of x , assumed to be continuous.

consider a_0 and a_1 $P_{x \sim p^*} [x \in (a_0, a^*]) = \frac{\varepsilon}{2}$ and $P_{x \sim p^*} [x \in (a^*, a_1)] = \frac{\varepsilon}{2}$



let $S = \{(x_i, y_i), \dots, (x_m, y_m)\}$ be some training set with iid examples.

write $b_0 = \max \{x_i : (x_i, 1) \in S\}$ and $b_1 = \min \{x_i : (x_i, 0) \in S\}$



(we assume that b_0, b_1 are well-defined in \mathbb{R}) Note that $b_0 = b_0(s)$ and $b_1 = b_1(s)$

let b_s be the threshold of an ERM rule $h_s \in \arg\min L_S(h_s)$

Recall that $L_S(h_s) = 0$ with probability 1 (see chapter 2).

$$\Rightarrow \mathbb{1}_{x \in S} h_s(x_i) = y_i \quad \text{for } i \in \{1, \dots, m\}$$

$$\Rightarrow \begin{cases} x_i < b_s \Leftrightarrow y_i = 1 \Rightarrow b_0 < b_s \Leftrightarrow b_0 < b_s \leq b_1 \\ x_i \geq b_s \Leftrightarrow y_i = 0 \Rightarrow b_1 \geq b_s \end{cases}$$

We will show that $\{s : b_0(s) \geq a_0 \text{ and } b_1(s) \leq a_1\} \subset \{s : L_D(b_s) \leq \varepsilon\}$

Assume that the event on the left

$$L_D(h_s) = P_{x \sim p^*} (\mathbb{1}_{x < b_s} + \mathbb{1}_{x > a^*})$$

$$= P_{x \sim p^*} (x < b_s \text{ and } x \geq a^*) + P_{x \sim p^*} (x \geq b_s \text{ and } x < a^*)$$

$$\text{on the left} \leq P(a^* < x < b_1) + P(b_0 < x < a^*)$$

$$\text{above event!} \leq P(a^* < x < a_1) + P(a_0 < x < a^*)$$

$$= \frac{\varepsilon}{2} \quad = \frac{\varepsilon}{2} \quad \text{(by continuity of } D^*)$$

This implies that $P_{S \sim \mathcal{G}^m} (L_S(h_s) > \varepsilon) \leq P_{S \sim \mathcal{G}^m} (b_0 > a_0 \text{ or } b_1 > a_1)$

$$< P_{S \sim \mathcal{G}^m} (b_0 > a_0) + P_{S \sim \mathcal{G}^m} (b_1 > a_1)$$

$$P_{S \sim \mathcal{G}^m} (b_0 > a_0) = P_{S \sim \mathcal{G}^m} (\forall x_i \text{ if } x_i < a^* \text{ then } x_i < a_0)$$

$$= P_{S \sim \mathcal{G}^m} (\forall i, x_i \notin [a_0, a^*])$$

Corollary: Let H be some hypothesis class. Let m be the size of a training set. Assume $\exists C \subset X$ such that $|C| = m$ and the size m is shattered by H . Then, for any learning algorithm A , we can find a distribution D on $X \times \{0, 1\}$ such that

$$1. \exists \text{ classifier } f: L_D(f) = 0$$

$$2. P_{S \sim D^m} (L_D(A(S)) \geq \frac{1}{2}) \geq \frac{1}{2}$$

Definition: (VC dimension)

$\rightarrow \forall C: |C| = d \Rightarrow H \text{ shatters } C$

The VC dimension of some class H , denoted by $\text{VC dim}(H)$, is the maximal size of a set $C \subset X$ that can be shattered by H .

If H can shatter sets of any size, then we say that H has an infinite VC-dimension ($\text{VC dim}(H) = +\infty$)

Theorem Let H be some hypothesis class with $\text{VC dim}(H) = +\infty$.

Then H is not PAC learnable.

Proof: If H is such that $\text{VC dim}(H) = +\infty$, then for any $m \geq 1$,

we can find a set $C \subset X$ such that $|C| = m$ and C is shattered by H . Then, the claim follows from the corollary. \square

Remarks:

- To show that $\text{VC dim}(H) \leq d$, it is enough to show that H does not shatter any set C of size $d+1$.

- To show that $\text{VC dim}(H) \geq d$, it is enough to find a set C of size d such that C is shattered by H .

- To show that $\text{VC dim}(H) = d$, then we need to show
 - $\exists C: |C| = d$ such that H shatters C
 - $\forall C: |C| = d+1$, H does not shatter C

- * If H shatters C , then it shatters any subset $C' \subset C$.

- * If X is finite, then for any $C \subset X$ we have $|H_C| \leq |H|$ if $|C|$ such that $|H| < 2^{|C|}$ ($\Leftrightarrow |C| > \log_2(|H|)$)

- then $|H_C| < 2^{|C|} \Rightarrow H$ does not shatter C .

- Since this true for any $C: |C| > \log_2(|H|)$

$$= \left[P(x_1 \notin [a_0, a^*]) \right]^m \leq \left(1 - P(x_1 \in [a_0, a^*]) \right)^m$$

$$= \left(1 - \frac{\varepsilon}{2} \right)^m$$

Similarly, we can show that $P_{S \sim \mathcal{D}^m}(b_i > a_i) = P(\forall i, x_i \notin [a^*, a_i])$

$$= \left[P(x_i \notin [a^*, a_i]) \right]^m$$

$$\leq P(S \text{ not close to } a^*) = \left(1 - \frac{\varepsilon}{2} \right)^m$$

Therefore $P_{S \sim \mathcal{D}^m}(L_{\mathcal{D}}(h_S) > \varepsilon) \leq 2 \left(1 - \frac{\varepsilon}{2} \right)^m \leq 2 e^{-m \frac{\varepsilon}{2}} \leq 0$

$$\Leftrightarrow e^{-\frac{m\varepsilon}{2}} \leq \frac{\delta}{2} \Leftrightarrow \frac{m\varepsilon}{2} \geq \log\left(\frac{2}{\delta}\right)$$

$$\Leftrightarrow m \geq \frac{\log\left(\frac{2}{\delta}\right)}{\varepsilon}$$

$$\Rightarrow m_H(\varepsilon, \delta) \leq \frac{\log\left(\frac{2}{\delta}\right)}{\varepsilon}$$

as claimed \square

2. The VC dimension

Definition: (restriction of H to C)

let H be a hypothesis class of function from $X \rightarrow \{0, 1\}$. let $C = \{c_1, \dots, c_m\} \subset X$. The restriction of H to C is $H_C = \{h(c_1), \dots, h(c_m)\}$

Definition: (shattering)

A hypothesis class H shatters a finite set $C \subset X$ if $|H_C| = 2^{|C|}$

In other words, H shatters C when H_C is the set of all possible classifiers that can be defined on C .

Example: ① Consider again the class $H = \{x \mapsto 1_{x \geq a} : a \in \mathbb{R}\}$

$C = \{c\}$ with c a fixed real number. There are 2^1 classifiers on C : $f(c) = 0$ and $f(c) = 1$

0 is obtained by $c-1 = a$, 1 is obtained by $a = c+1$

② Take $C = \{c_1, c_2\}$ with $c_1 < c_2$. There are ≥ 4 possible classifiers on C .

if $a \leq c_1$, then $h_a(c_1) = h_a(c_2) = 0$

if $a < c_1 \leq c_2$, then $h_a(c_1) = 1$ and $h_a(c_2) = 0$.

if $c_2 > a$, then $h_a(c_1) = h_a(c_2) = 1$

(hence $|H_C| = 3 \geq 4 \Rightarrow H$ does not shatter $C = \{c_1, c_2\}$ (this is only true for $A \subset C$))

this implies that $Vc \dim(H) \leq \log_2(|H|)$

What if $|H| = 2^{|X|}$? In this case, set $Vc \dim(H) = |X|$. This means that H does not shatter any subset C of size $|X|+1$ (this is true since such a subset does not exist). Note that $\log_2(|H|) = \log_2(2^{|X|}) = |X|$.

3. Examples:

3.1 The class of thresholds:

$$H = \{x \mapsto 1_{x < a} : a \in \mathbb{R}\}$$

We have shown that $\begin{cases} H \text{ shatters } C = \{c\} \text{ (for } c \in \mathbb{R}) \\ H \text{ does not shatter any } C = \{c_1, c_2\} = c_1 \neq c_2 \end{cases}$

This implies that $Vc \dim(H) \geq 1$

3.2 Class of intervals $\underbrace{h_{[a,b]}(x)}$

$$\text{Consider } H = \{x \mapsto 1_{x \in (a,b)} : -\infty < a \leq b < +\infty\}$$

We can show that H shatters $C = \{0, 1\}$ (show this!)

Now, consider $C = \{c_1, c_2, c_3\}$ with $c_1 < c_2 < c_3$

$$\overbrace{c_1 \quad c_2 \quad c_3}$$

If $a < b$ are chosen such that $h_{[a,b]}(c_1) = 1$ and $h_{[a,b]}(c_2) = 0$ then we must have $h_{[a,b]}(c_3) = 0$

$$\overbrace{* \quad | \quad * \quad | \quad |} \quad a \quad c_1 \quad b \quad c_2 \quad c_3$$

This implied that the classifier f such that $f(c_1) = 1$, $f(c_2) = 0$ and $f(c_3) = 1$ is impossible to obtain.

Any $C = \{c_1, c_2, c_3\}$ can not be shattered by H .

This implies that $Vc \dim(H) \geq 2$.

Remark: 3.3 Finite classes:

We have seen that if X is finite (which means that H is finite) $Vc \dim(H) \leq \log_2(|H|)$

But $\log_2(|H|) - Vc \dim(H)$ can be quite large.

$$X = \{1, 2, \dots, k\} \quad H = \{x \mapsto 1_{x \in i} : i \in \{1, \dots, k\}\}$$

$\text{VC dim}(H) = 1$ but $\log_2(|H|) = \log_2(k)$ which $\rightarrow \infty$
if $k \rightarrow \infty$

4. The fundamental Theorem of PAC learning:

We have seen that $\text{VC dim}(H) = +\infty \Rightarrow H$ is not PAC learnable.

Question: Do we have $\text{VC dim}(H) = d \Leftrightarrow$

H is PAC learnable?

Theorem: Let H be some hypothesis class: $X \rightarrow \{0, 1\}$ and let the loss function be the 0-1 loss. Then the following assertions are equivalent:

1. H has the Uniform Convergence property
2. Any ERM rule is a successful agnostic PAC learner for H
3. H is agnostic PAC learnable
4. H is PAC learnable
5. Any ERM rule is a successful PAC learnable for H
6. H is of finite VC dimension.

Proof: (partial)

$$1 \Rightarrow 2 \quad (\text{chapter 4})$$

$$2 \Rightarrow 3 \Rightarrow 4 \quad (\text{trivial})$$

$$2 \Rightarrow 5 \quad (\text{trivial})$$

4 \Rightarrow 6 and 5 \Rightarrow 6 (consequence of the No-free lunch theorem)

If we show 6 \Rightarrow 1 which needs Sauer's lemma.

Then it is proved.

3. Sauer's lemma and the growth of function.

Question: When the size of the growth, how does H_k grow?

$$P[\exists \text{ } f_{\text{fun}}] = P[\exists \int_0^t \epsilon_s d\beta_s \mid f_{\text{fun}}]$$

↓
εC_b. ||
X^a

$A \in \mathcal{F}_t \iff A \cap \mathcal{F}_{t \cup} \quad A \cap (\mathcal{T}_H < t) \in \mathcal{F}_t.$

(VC-dimension). [6.5]

H VC dim(H) is the maximal size of a set C such that H shatters C , $|H_C| = 2^{|C|}$

$|C| = m \text{ if } H \text{ shatters } C \text{最多能shatter的数量.}$

$(m \geq d)$.

(Growth function)

(Sauer lemma)

$$\tau_H(m) = \max_{C \subseteq X, |C|=m} |H_C|$$

$$H, \text{VC dim}(H) \leq d \Rightarrow \tau_H(m) \leq \sum_{i=0}^d \binom{m}{i}$$

if $m \leq \text{VC dim}(H)$, $\tau_H(m) = 2^m$

Proof: $\forall H, |H_C| \leq |\mathcal{B}| \subseteq H \text{ shatters } \mathcal{B} \mid (\forall) \text{ (like } \mathbb{1}_{\mathcal{B}}).$

$$\leq \sum_{i=0}^d \binom{m}{i}$$

To show (*): $m=1, |H_C| = \begin{cases} 1 \\ 2 \end{cases} \quad \mathcal{B} = \begin{cases} \emptyset \\ \{c_1\} \end{cases} \quad |\mathcal{B} \subseteq C : H \text{ shatters } \mathcal{B}| = 1, 2.$

general m , assume it holds for $k \leq m$

給定一個 C 且 $C = \{c_1, c_2, \dots, c_m\}, |H_C| = |\mathcal{Y}_0| + |\mathcal{Y}_1| \rightarrow$ \mathcal{Y}_0 可以被 c_1 補充
shatter 的子集 \mathcal{B} → $\mathcal{B} \subseteq C$ 且 H shatters \mathcal{B} .

且 \mathcal{B} shatters $C \setminus \{c_1\}$

有證明:

only see one value $\mathcal{Y}_0 = \{(y_2, \dots, y_m), (0, y_2, \dots, y_m) \in \mathcal{C} \text{ or } (1, y_2, \dots, y_m) \in \mathcal{H}_0\}$

$$\mathcal{Y}_1 = \{ \quad \text{---} \quad \text{---} \quad \text{and} \quad \text{---} \quad \text{---} \quad \}$$

$y_1 \quad y_m \mid \vee$

$(y_1 \in \mathcal{Y}_0 \rightarrow)$

$$|\mathcal{Y}_0| = |\mathcal{H}_0| = |H_C| \leq |\{\mathcal{B} \subseteq C : H \text{ shatters } \mathcal{B}\}|$$

$$= |\{\mathcal{B} \subseteq C : c_1 \notin \mathcal{B} \text{ and } H \text{ shatters } \mathcal{B}\}|$$

$$H' \subseteq H, H' = \{ h \in H : \exists h' \in H \text{ } (1-h'(c_1), h'(c_2), \dots, h'(c_m)) \\ = (h(c_1), \dots, h(c_m)) \}$$

$$|Y_1| = |H'_C| \leq |B \subseteq C| = H \text{ shatters } B \\ = |\{B \subseteq C : H \text{ shatters } B \cup \{c_i\}\}| \\ = |\{B \subseteq C : \forall q \in B \text{ and } H \text{ shatters } B\}| \\ \leq |\{B \subseteq C : \forall q \in B \text{ and } H \text{ shatters } B\}|$$

$$|H_C| = |Y_1|, |Y_1| \leq |\{B \subseteq C : H \text{ shatters } B\})|$$

□

相当于新增的点能被 shatter.

$$\text{if } m > d+1, T_H(m) \leq \left(\frac{em}{d}\right)^d$$

$$\text{Proof: if } \frac{d}{m} \leq 1 \Rightarrow \left(\frac{d}{m}\right)^d \leq 1$$

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

$$\leq \left(1 + \frac{d}{m}\right)^m \leq \left(1 + \frac{d}{m}\right)^m$$

$$\Rightarrow \sum_{i=0}^d \binom{m}{i} \leq \left(1 + \frac{d}{m}\right)^m \cdot \left(\frac{m}{d}\right)^d$$

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

□

+ 2x slides!

The fundamental theorem of PAC learning

Recall the SGD algorithm runs as follows

• start: $w^{(0)} = 0$

• For $t \in \{1, \dots, T\}$

generate v_t from a distribution s.t.

$$E[v_t | w^{(t)}] \in \delta f(w^{(t)})$$

$$\downarrow \text{update } w^{(t+1)} = w^{(t)} - \eta v_t$$

• Output: $\bar{w} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$.

fixed step size

Theorem (Convergence of SGD)

Let $B > 0$ and $\ell > 0$, let f be a convex function with

$w^* = \operatorname{argmin}_w f(w)$. Assume the SGD algorithm is run for T

iterations with $\eta = \frac{B^2}{\ell^2 T}$. Assume that v_t satisfies $\|v_t\| \leq \rho$

with probability $1 - \sqrt{\epsilon} \in \{1, \dots, T\}$. Then $E[f(\bar{w})] - f(w^*) \leq \frac{B\rho}{\sqrt{T}}$

proof: skipped.

Learning with SGD for risk minimization

Recall that the main goal is to minimize some risk function

consider the risk function $L_D(w) = E_{z \sim D}[\ell(w, z)]$

* Analysing convergence of SGD for convex and p -Lipschitz loss functions:

Assume that $w \mapsto \ell(w, z)$ is convex and p -Lipschitz.

Assume that it is differentiable ($\forall z \in Z$).

Then by the dominated convergence theorem, the function $w \mapsto L_D(w)$ is differentiable if \exists function $z \mapsto k(z)$ s.t. $k \geq 0$ and $E[k(z)] < \infty$, and $\|\nabla \ell(w, z)\| \leq k(z) \quad \forall w, z$.

Furthermore, we have that $\nabla L_D(w) = E_{z \sim D}[\nabla \ell(w, z)]$.

Note that if $w \mapsto \ell(w, z)$ is p -Lipschitz. And differentiable

then $\|\nabla l(u, z)\| \leq p \quad \forall u, z \Rightarrow$ we can take $k(z) = p$. $\forall z$

Define $v_t = \nabla l(u^{(t)}, z)$ where $z \sim \mathcal{D} \perp\!\!\!\perp u^{(t)}$. we know that

$$\begin{aligned} \|v_t\| &\leq p. \text{ let us compute } E[v_t | u^{(t)}] = E[\nabla l(u^{(t)}, z) | u^{(t)}] \\ &\stackrel{z \sim \mathcal{D}}{=} E[\nabla l(u^{(t)}, z)] \\ &\stackrel{z \sim \mathcal{D}}{=} \nabla E[l(u^{(t)}, z)] \\ &= \nabla L_{\mathcal{D}}(u^{(t)}) \quad (\text{the only element in } \nabla L_{\mathcal{D}}(u^{(t)})). \end{aligned}$$

Define $v_t \in \partial l(u^{(t)}, z)$ with $z \sim \mathcal{D} \perp\!\!\!\perp u^{(t)}$. we have again $\|v_t\| \leq p$. I shall applying properties of subgradient for p -Lipschitz function.)

$$\begin{aligned} \text{fixed } u \quad l(u, z) &\geq l(u^{(t)}, z) + \langle u - u^{(t)}, v_t \rangle. \\ \Rightarrow E[l(u, z) | u^{(t)}] &\geq E[l(u^{(t)}, z) | u^{(t)}] + E[\langle u - u^{(t)}, v_t \rangle | u^{(t)}] \\ &\stackrel{z \sim \mathcal{D}}{\geq} E[l(u, z)] \geq E[\nabla l(u^{(t)}, z)] + \langle u - u^{(t)}, E[v_t | u^{(t)}] \rangle \\ &\Leftarrow L_{\mathcal{D}}(u) \geq L_{\mathcal{D}}(u^{(t)}) + \langle u - u^{(t)}, E[v_t | u^{(t)}] \rangle. \end{aligned}$$

The previous inequality holds $\forall u$ and hence $E[v_t | u^{(t)}] \in \partial L_{\mathcal{D}}(u^{(t)})$. Therefore we can give the SGD algorithm more precisely by the following steps.

- * Start: $u^{(0)} = 0$ generate $z \sim \mathcal{D} \perp\!\!\!\perp u^{(0)}$
- * For $t = 1, \dots, T$ $z \sim \mathcal{D}$
 - pick $v_t \in \partial l(u^{(t)}, z)$
 - update $u^{(t+1)} = u^{(t)} + \gamma v_t$
- * Output $\bar{u} = \frac{1}{T} \sum_{t=1}^T u^{(t)}$.

Corollary: Consider a convex - Lipschitz - bounded learning problem with parameters $p > 0$ and $B > 0$ ($H = \{w \in \mathbb{R}^d : \|w\| \leq B\}$) and

$w \mapsto l(w, z)$ is convex p -Lipschitz.)

→ If we run the SGD algorithm for minimizing $w \mapsto l(w)$ (con A) with $T \geq \frac{B^2 p^2}{\varepsilon^2}$ and $\eta = \sqrt{\frac{B^2}{p^2 T}}$, then

$$\mathbb{E}[l(\bar{w})] \leq \min_{w \in \mathcal{H}} L_\phi(w) + \varepsilon.$$

• Analyzing convergence of SGD for convex-smooth loss functions.

Theorem Assume that $\forall z \in \mathcal{H} : w \mapsto l(w, z)$ is convex and β -smooth for some $\beta > 0$. Then if we run the SGD algorithm, we have that and if $\eta < \frac{1}{\beta}$, then

$$\mathbb{E}[L_\phi(\bar{w})] \leq \frac{1}{1-\eta\beta} \left(L_\phi(w^*) + \frac{\|w^*\|^2}{2\eta T} \right)$$

Proof: Recall that if f is a β -smooth and non-negative function, then

$\| \nabla f(w) \|^2 \leq \beta f(w)$. For $t \in \{1, \dots, T\}$, let us denote $z_t = z \sim \mathcal{D}$ at iteration t s.t. $z_t \perp \!\!\! \perp w^{(t)}$, $f_t(w) = l(w, z_t)$.

Then, $v_t = \nabla f_t(w^{(t)})$. By convexity of f_t , we have that $f_t(w^{(t)}) \leq f_t(w^*) + \langle w^{(t)} - w^*, v_t \rangle$

$$\begin{aligned} \text{Summing over } t = 1, \dots, T \text{ yields } & \sum_{t=1}^T f_t(w^{(t)}) \leq \sum_{t=1}^T f_t(w^*) \\ & + \underbrace{\sum_{t=1}^T \langle w^{(t)} - w^*, v_t \rangle}_{\text{at } w^{(t)}} \\ & \leq \frac{\|w^*\|^2}{2\eta} + \frac{1}{2} \sum_{t=1}^T \|v_t\|^2 \end{aligned}$$

using Lemma 14.1
Since $\# \geq \eta$ (weak #11)

Using self-boundness of f_t , it follows that $\sum_{t=1}^T \|v_t\|^2 \leq \beta \sum_{t=1}^T f_t(w^{(t)})$. Hence, $\frac{1}{T} \sum_{t=1}^T f_t(w^*) \leq \frac{1}{T} \sum_{t=1}^T f_t(w^*) + \frac{\|w^*\|^2}{2\eta T} + \eta \beta \frac{1}{T} \sum_{t=1}^T f_t(w^{(t)})$.

$$\begin{aligned} & \Leftrightarrow (1-\eta\beta) \frac{1}{T} \sum_{t=1}^T f_t(w^{(t)}) \leq \frac{1}{T} \sum_{t=1}^T f_t(w^*) + \frac{\|w^*\|^2}{2\eta T} \\ & \Leftrightarrow \frac{1}{T} \sum_{t=1}^T f_t(w^{(t)}) \leq \frac{1}{1-\eta\beta} \left(\frac{1}{T} \sum_{t=1}^T f_t(w^*) + \frac{\|w^*\|^2}{2\eta T} \right) \\ & \Rightarrow \frac{1}{T} \sum_{t=1}^T \mathbb{E}[f_t(w^{(t)})] \leq \frac{1}{1-\eta\beta} \left(\frac{1}{T} \sum_{t=1}^T \mathbb{E}[f_t(w^*)] + \frac{\|w^*\|^2}{2\eta T} \right). \end{aligned}$$

Now, $\mathbb{E}[f_t(w^*)] = \mathbb{E}_{z \sim \mathcal{D}}[l(w^*, z)] = L_\phi(w^*) \quad \forall t$.

$$\Rightarrow \frac{1}{T} \sum_{t=1}^T \mathbb{E}[f_t(w^{(t)})] = L_\phi(w^*)$$

$$\mathbb{E}[f_t(w^{(t)})] = \mathbb{E}[l(w^{(t)}, z^{(t)})] = \mathbb{E}[L_\phi(w^{(t)})]$$

$$\text{Hence, } \frac{1}{T} \sum_{t=1}^T \mathbb{E}[f_t(w^{(t)})] = \frac{1}{T} \sum_{t=1}^T [L_\phi(w^{(t)})].$$

□

$$\text{Now, } \frac{1}{T} \sum_{t=1}^T \mathbb{E}[L_D(w^{(t)})] = \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T L_D(w^{(t)})\right]$$

The function $w \mapsto L_D(w) = \mathbb{E}_{z \sim D}[l(w, z)]$ is convex (by convexity of $w \mapsto l(w, z)$ for $\forall z$).

Hence, $\frac{1}{T} \sum_{t=1}^T L_D(w^{(t)}) \geq L_D\left(\frac{1}{T} \sum_{t=1}^T w^{(t)}\right)$ (Tensen's inequality).

We conclude that $\mathbb{E}[L_D(\bar{w})] \leq \frac{1}{T} \mathbb{E}\left[L_D(w^*) + \frac{\|w^*\|^2}{2T}\right]$ as claimed.

Corollary: Consider a convex-smooth-bounded learning problem with parameters $\beta > 0$ and $B > 0$. Assume in addition that $l(0, z) \leq 1 \quad \forall z \in \mathcal{Z}$

[Exco. 1]. Set $\eta = \frac{1}{\beta(1+\frac{\varepsilon}{2})}$. Then running SGD with $T \geq \frac{12B^2}{\varepsilon^2}$

yields $\mathbb{E}[L_D(\bar{w})] \leq \min_{w \in H} L_D(w) + \varepsilon$ ($H = \{w \in \mathbb{R}^d : \|w\| \leq B\}$)

Proof: $1 - \eta\beta = \frac{3}{\varepsilon+3} \iff \frac{1}{1-\eta\beta} = 1 + \frac{\varepsilon}{3}$

By the previous theorem, we have that

$$\begin{aligned} \mathbb{E}[L_D(\bar{w})] &\leq \left(1 + \frac{\varepsilon}{3}\right) \left(L_D(w^*) + \frac{B^2}{2T}\right), \text{ where } w^* = \underset{w \in H}{\operatorname{argmin}} L_D(w) \\ &\leq \left(1 + \frac{\varepsilon}{3}\right) L_D(w^*) + \left(1 + \frac{\varepsilon}{3}\right) \frac{\varepsilon^2 + 3\varepsilon}{24} \end{aligned}$$

$$L_D(w^*) = L_D(0) = \mathbb{E}_{z \sim D}[l(0, z)] \leq 1$$

$$\leq L_D(w^*) + \frac{\varepsilon}{3} + \left(1 + \frac{1}{\varepsilon}\right) \frac{4}{24} \varepsilon$$

$$\leq L_D(w^*) + \varepsilon.$$

□.

Regularization and stability (chapter 13).

1 Regularized loss minimization (RLM)

RLM is a learning rule which minimize the criterion $w \mapsto l(w) + R(w)$ for a given regularization function R .

Example: If $R(w) = \lambda \|w\|^2 = \lambda \sum_{i=1}^d w_i^2$, then we talk about

when applying Tikhonov regularization to linear regression with

$\log(w) = (\langle w, x \rangle - y)^2$, we obtain the learning rate.

$$\hat{w} = \underset{w}{\operatorname{argmin}} \left(\lambda \|w\|^2 + \frac{1}{m} \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2 \right)$$

$$\text{Put } f(w) = \lambda \|w\|^2 + \frac{1}{m} \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2 \text{ for } w \in \mathbb{R}^d$$

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

$$\nabla f(w) = 2\lambda w + \frac{1}{m} \sum_{i=1}^m x_i (x_i^\top w - y_i) = 0.$$

$$\Leftrightarrow \left[\underbrace{\sum_{i=1}^m x_i x_i^\top}_{A} + \lambda m I \right] w = \sum_{i=1}^m y_i x_i$$

The matrix A is semi-positive definite because for $\forall a \in \mathbb{R}^d$,

$$a^\top A a = a^\top \sum_{i=1}^m x_i x_i^\top a = \sum_{i=1}^m (x_i^\top a)^2 \geq 0$$

$\Rightarrow \exists$ orthogonal matrix $Q \in \mathbb{R}^{d \times d}$ such that

$$A = Q^\top \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_d \end{pmatrix} Q \text{ with } \lambda_i \geq 0, \forall i \in \{1, \dots, d\}.$$

$$\Rightarrow A + \lambda m I = Q^\top \begin{pmatrix} \lambda_1 + \lambda m & & \\ & \ddots & \\ & & \lambda_d + \lambda m \end{pmatrix} Q$$

Since $\lambda_i + \lambda m > 0 \quad \forall i \in \{1, \dots, d\}$, $A + \lambda m I$ is invertible

$$\text{and } \hat{w} = [A + \lambda m I]^{-1} \left(\sum_{i=1}^m y_i x_i \right)$$

2. Stable rules do not overfit.

Let A be a learning algorithm and $S = \{z_1, \dots, z_m\}$ be some training set with m iid example $\sim D$ (unknown distribution)

Denote $A(S)$ the output of A when it is run on S .

To assess stability of A , we look at the influence of replacing an example z_i by some $z' \sim D$.

Given S and an additional example $z' \sim D \setminus S$, let

$$S^{(i)} = \{z_1, \dots, z_{i-1}, z', z_{i+1}, \dots, z_m\}.$$

We investigate $l(A(S^{(i)}), z_i) - l(A(S), z_i)$

\downarrow \downarrow
does not observe z_i observe z_i .

Theorem Let D be some distribution, $S = \{z_1, \dots, z_m\}$ a training set where z_1, \dots, z_m are iid $\sim D$ and $z' \sim D \setminus S$.

Then, for any learning algorithm we have that

$$E_{\text{sign}} [l_D(A(S)) - l_S(A(S))] = E_{\substack{(S, z', i) \sim D^{m+1} \times U[m]} \left[l(A(S^{(i)}), z_i) - l(A(S), z_i) \right]} \quad (*)$$

where $[m] = \{1, 2, \dots, m\}$.

~~$$E_{\text{sign}} [l_D(A(S))] = E_{\substack{S \sim D^m}} [E_{z' \sim D} (l(A(S), z'))] \quad (z' \setminus S)$$~~

$$= \mathbb{E}_{(s, z) \sim \mathcal{D}^{m+1}} [l(A(s), z)]$$

$$\{z_1, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_m\}$$

$$= \mathbb{E}_{(s, z) \sim \mathcal{D}^m} [l(A(\{z_1, \dots, z_{i-1}, z_m\}), z)]$$

$$= \mathbb{E}_{(s, z) \sim \mathcal{D}^m} [l(A(s')), z)].$$

Since the last expectation does not depend on the index i , it follows that

$$\mathbb{E}_{(s, z) \sim \mathcal{D}^m} [z] = \mathbb{E}_{(s, z) \sim \mathcal{D}^m \times \mathcal{U}^m} [z].$$

On the other hand, we have that

$$\begin{aligned} \mathbb{E}_{\mathcal{S}^m} [L_s(A(s))] &= \mathbb{E}_{\mathcal{S}^m} \left[\frac{1}{m} \sum_{i=1}^m l(A(s), z_i) \right] \text{ by definition of} \\ &= \mathbb{E}_{(s, z) \sim \mathcal{D}^m \times \mathcal{U}^m} [l(A(s), z)] \text{ the training error/rate} \end{aligned}$$

Same proof as before. \square .

When the right term of (x) is small, this is a indication that replacing z_i by "something else" does not influence the algorithm "too much". In other words, we say that the algorithm is stable (a change in a single example does not result in a big change on average.)

Definition (on-average-replace-one-stable)

Let $\varepsilon: \mathbb{N} \rightarrow \mathbb{R}$ be a decreasing function, we say that some learning algorithm A is on-average-replace-one-stable with rate $\varepsilon(m)$ if \forall distribution \mathcal{D} $\mathbb{E}_{(s, z) \sim \mathcal{D}^{m+1} \times \mathcal{U}^m} [l(A(s'), z) - l(A(s), z)] \leq \varepsilon(m)$

Remark: In view of the previous theorem, when this property holds then \forall distribution \mathcal{D}

$$\mathbb{E}_{\mathcal{S}^m} [L_{\mathcal{D}}(A_{m+1}) - L_s(A(s))] \leq \varepsilon(m) \leftarrow \text{means you one out overfitting}$$

This means that A does not overfit \rightarrow A "good" learning algorithm should balance between fitting and being stable.

3. Tikhonov regularization as a stabilizer:

In the following, we will try to apply Tikhonov regularization to convex and p -Lipschitz loss function.

definition: (Strongly convex functions) A function f is said to be λ -strongly convex if $\forall u, v \in \mathbb{R}^n$ and $\alpha \in [0, 1]$

$$f(\alpha u + (1-\alpha)v) \leq \alpha f(u) + (1-\alpha)f(v) - \frac{\lambda}{2}\alpha(1-\alpha)\|u-v\|^2.$$

- Lemma
1. The function $f(u) = \lambda\|u\|^2$ is $(>\lambda)$ -strongly convex.
 2. If f is λ -strongly convex and g is convex, then fg is λ -strongly convex

3. If f is λ -strongly convex and u is the minimizer of f , then $f(u) - f(v) \geq \frac{\lambda}{2}\|u-v\|^2$.

Proof 1 and 2 are easy to show.

3. Let u be the minimizer of f . For $\alpha \in [0, 1]$,

$$f(\alpha u + (1-\alpha)v) = f(\alpha u + (1-\alpha)v) \leq \alpha f(u) + (1-\alpha)f(v) - \frac{\lambda}{2}\alpha(1-\alpha)\|u-v\|^2$$
$$\Leftrightarrow \underbrace{f(u + \alpha(v-u)) - f(u)}_{\alpha} \leq f(v) - f(u) - \frac{\lambda}{2}(1-\alpha)\|u-v\|^2$$

Since $f(u) = \inf_v f(v)$, then $f(u + \alpha(v-u)) \geq f(u)$, and hence

$$0 \leq f(u) - f(u) - \frac{\lambda}{2}(1-\alpha)\|u-v\|^2$$

$$\Leftrightarrow f(u) - f(v) \geq \frac{\lambda}{2}(1-\alpha)\|u-v\|^2 \quad \forall \alpha \in [0, 1]$$

Taking the limit as $\alpha \downarrow 0$, we get that $f(u) - f(v) \geq \frac{\lambda}{2}\|u-v\|^2$ \square

for a training set $S = \{z_1, \dots, z_n\}$, consider the RLM Tikhonov $A(s) = \arg \min_u$

$(L_S(u) + \lambda\|u\|^2)$. Let us write $f_S(u) = L_S(u) + \lambda\|u\|^2$. Under convexity of $u \mapsto L_S(u)$, $\forall z \in S$ and using the previous lemma, we see that f_S is

$(>\lambda)$ -strongly convex. Also, by the ^(①+②) same lemma (③), we have $\forall v$,

$$f_S(u) - f_S(A(s)) \geq \lambda\|A(s) - v\|^2 \quad (*)$$

for any v, u and index $i \in \{1, \dots, n\}$

$$\begin{aligned} f_s(v) - f_s(u) &= L_s(v) + \lambda \|v\|^2 - L_s(u) - \lambda \|u\|^2 = L_{s^{(i)}}(v) - \frac{1}{m} l(v, z_i) + \frac{1}{m} l(v, z_i) \\ &\quad + \lambda \|v\|^2 - (L_{s^{(i)}}(u) - \frac{1}{m} l(u, z_i) + \frac{1}{m} l(u, z_i)) \\ &= L_{s^{(i)}}(v) + \lambda \|v\|^2 - (L_{s^{(i)}}(u) + \lambda \|u\|^2) + \frac{l(v, z_i) - l(u, z_i)}{m} + \lambda \|u\|^2 \end{aligned}$$

Now take $u = A(s)$ and $v = A(s^{(i)})$ since v minimizes $\sum_i l(v, z_i) + \lambda \|v\|^2$, it follows that,

$$f_s(A(s^{(i)})) - f_s(A(s)) \leq \frac{l(A(s^{(i)}), z_i) - l(A(s), z_i)}{m} + \frac{l(A(s), z_i) - l(A(s^{(i)}), z_i)}{m} - \lambda \|A(s^{(i)}) - A(s)\|^2$$

If the loss function is also p -Lipschitz (in addition to convexity), then

$$\text{Term} \leq |\text{Term}| \leq \frac{2\ell}{m} \|A(s^{(i)}) - A(s)\|.$$

$$\begin{aligned} \text{Therefore we get } \lambda \|A(s^{(i)}) - A(s)\|^2 &\leq \frac{2\ell}{m} \|A(s^{(i)}) - A(s)\| \\ \Leftrightarrow \|A(s^{(i)}) - A(s)\| &\leq \frac{2\ell}{m\lambda} \end{aligned}$$

ℓ -Lipschitz.

$$\text{This also implies that } |l(A(s^{(i)}), z_i) - l(A(s), z_i)| \leq 1 \quad 1 \leq \frac{2\ell^2}{\lambda m}.$$

$\forall s, i, (\forall z \in \mathcal{D})$.

□

Corollary: If the loss function is convex and p -Lipschitz, then the Tikhonov RLM rule is on-average-replace-one-stable with rate $\mathbb{E}_{\text{avg}}[\dots] \leq \frac{2\ell^2}{\lambda m}$. We have $\mathbb{E}_{\text{avg}}[L_D(A(s)) - L_S(A(s))] \leq \frac{2\ell^2}{\lambda m}$.

4. Controlling the fitting-stability trade-off (for convex and p -Lipschitz loss functions).

$$\begin{aligned} \text{Let us start with } \mathbb{E}_{\text{avg}}[L_D(A(s))] &= \mathbb{E}_{\text{avg}}[L_S(A(s))] + \mathbb{E}_{\text{avg}}[L_D(A(s)) - L_S(A(s))] \\ &\quad \text{fit.} \quad \text{stability.} \end{aligned}$$

let $A(s) = \arg\min_u (L_S(u) + \lambda \|u\|^2)$, we have $\forall w^*, L_S(A(s)) \leq L_S(w^*) + \lambda \|w^*\|^2$.

$$\mathbb{E}_{\text{avg}}[L_S(w^*)] = L_D(w^*). \text{ Hence } \mathbb{E}_{\text{avg}}[L_S(A(s))] \leq L_D(w^*) + \lambda \|w^*\|^2.$$

$$\Rightarrow \mathbb{E}_{\text{avg}}[L_D(A(s))] \leq (L_D(w^*) + \lambda \|w^*\|^2) + \mathbb{E}_{\text{avg}}[L_D(A(s)) - L_S(A(s))]$$

corollary: If the loss function and p -Lipschitz, then the Tikhonov PLM rule satisfies

$$\mathbb{E}_{\text{sup}}[L_D(A(s))] \leq L_D(w^*) + \lambda \|w^*\|^2 + \frac{\rho^2}{\lambda m}.$$

we will take $w^* = \operatorname{argmin}_w L_D(w)$

$$w : \|w\| \leq B$$

corollary: Let (H, \mathcal{L}, l) be a convex, Lipschitz bounded learning problem with parameters ρ and B . Also, set $\lambda = \sqrt{\frac{2\rho}{B^2m}}$. Then, then Tikhonov PLM rule satisfies

$$\mathbb{E}_{\text{sup}}[L_D(A(s))] \leq \underbrace{\min_{w \in H} L_D(w) + \rho B \sqrt{\frac{8}{m}}}_{L_D(w^*)} \quad (H = \{w \in \mathbb{R}^d : \|w\| \leq B\}).$$

