

Concentration: Lower bounds for deviations, and No Free Lunch theorem

Master 2 Mathematics and Computer Science

Guillaume Aubrun, Aurélien Garivier, Rémi Gribonval
2020-2021



Table of contents

1. Deviation Bound for Bernoulli Variables
2. Kullback-Leibler divergence
3. No-Free-Lunch theorems: when learning is not possible

Deviation Bound for Bernoulli Variables

Chernoff's Bound

Theorem (Chernoff-Hoeffding Deviation Bound)

Let $\mu \in (0, 1)$. $X_1, \dots, X_n \stackrel{iid}{\sim} \mathcal{B}(\mu)$, and let $x \in (\mu, 1]$.

(i) Chernoffs' bound for Bernoulli variables:

$$\mathbb{P}(\bar{X}_n \geq x) \leq \exp(-n \text{kl}(x, \mu)), \quad (1)$$

where $\text{kl}(p, q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$. Same for left deviations.

(ii) If $\phi(x) = \text{kl}(x, \mu)$, then $\phi''(x) = 1/[x(1-x)]$ and

$$\begin{aligned} \text{kl}(x, \mu) &= \frac{(x - \mu)^2}{2} \int_0^1 \phi''(\mu + s(x - \mu)) 2(1-s) ds \\ &\geq \frac{(x - \mu)^2}{2\tilde{x}(1-\tilde{x})} \quad \text{with } \tilde{x} = \frac{2\mu + x}{3} \text{ by Jensen, since } \phi'' \text{ is convex and } \int_0^1 s 2(1-s) ds = \frac{1}{3} \\ &\geq \frac{1}{2 \max_{x \leq u \leq p} u(1-u)} (x - \mu)^2 \geq 2(x - \mu)^2. \end{aligned}$$

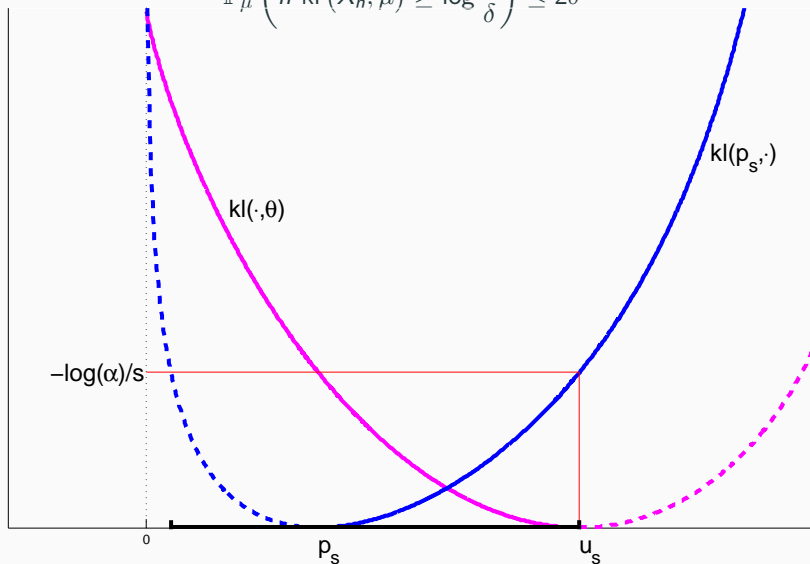
(iii) Hoeffding's bound for Bernoulli variables:

$$\mathbb{P}(\bar{X}_n \geq x) \leq \exp(-2n(x - \mu)^2). \quad (2)$$

(iv) Inequalities (1) and (2) hold for arbitrary independent random variables with range $[0, 1]$ and expectation μ .

A Divergence on the Set of Possible Means

$$\mathbb{P}_\mu \left(n \text{kl}(\bar{X}_n, \mu) \geq \log \frac{1}{\delta} \right) \leq 2\delta$$



Examples

- If $\mu < 1/2$,

$$\mathbb{P}\left(\bar{X}_n > \frac{1}{2}\right) \leq \exp\left(-\frac{n}{2}(1-2\mu)^2\right).$$

(Consequence of Chernoff or direct computation with $(1-u)^n \leq \exp(-nu)$, or of Hoeffding).

- For all $\mu \in [0, 1]$, Chernoff's bound with $\log(u) \geq (u-1)/u$ yields

$$\mathbb{P}\left(\bar{X}_n < \frac{\mu}{2}\right) \leq \exp\left(-\frac{1-\log(2)}{2} n\mu\right) \approx \exp(-0.153 n\mu) \leq \exp\left(-\frac{n\mu}{7}\right).$$

Hoeffding yields a very poor result, but (ii) gives:

$$\mathbb{P}\left(\bar{X}_n < \frac{\mu}{2}\right) \leq \exp\left(-\frac{3}{20} n\mu\right) = \exp(-0.15 n\mu) \leq \exp\left(-\frac{n\mu}{8}\right).$$

Sub-Gaussian inequalities

Bennett's and Bernstein's inequalities

Let $(X_i)_{1 \leq i \leq n}$ be independent random variables upper-bounded by 1, let $\bar{\mu} = (\mathbb{E}[X_1] + \dots + \mathbb{E}[X_n])/n$, let σ^2 be such that $\mathbb{E}[X_i^2] \leq \sigma^2$ for all i and let $\phi(u) = (1+u) \log(1+u) - u$. Then, for all $x > 0$,

$$\mathbb{P}(\bar{X} \geq \bar{\mu} + x) \leq \exp\left(-n\sigma^2\phi\left(\frac{x}{\sigma^2}\right)\right) \leq \exp\left(-\frac{nx^2/2}{\sigma^2 + x/3}\right).$$

Bernstein from Bennett: $\phi(x) \geq \frac{x^2}{2(1+\frac{x}{3})}$ since $\psi(x) = 2(1+\frac{x}{3})\phi(x) - x^2 \geq 0$.

Extension: if $X_i \leq b$ with $b > 0$,

$$\mathbb{P}(\bar{X}_n \geq \bar{\mu} + x) \leq \exp\left(-\frac{n\sigma^2}{b^2}\phi\left(\frac{bx}{\sigma^2}\right)\right) \leq \exp\left(-\frac{nx^2/2}{\sigma^2 + bx/3}\right).$$

Example: for X with range in $[0, 1]$,

$$\mathbb{P}\left(\bar{X}_n < \frac{\mu}{2}\right) \leq \exp\left(-n\left(\frac{3}{2}\log\frac{3}{2} - \frac{1}{2}\right)\mu\right) \leq \exp\left(-\frac{3n\mu}{28}\right).$$

Parenthesis: a nice proof for the technicalities of Bernstein

From [Pollard, MiniEmpirical ex.14, <http://www.stat.yale.edu/~pollard/Books/Mini/Basic.pdf>]

For any sufficiently smooth real-valued function g defined at least in a neighborhood of 0 let

$$G(x) = \frac{g(x) - g(0) - xg'(0)}{x^2/2} \text{ if } x \neq 0, \text{ and } G(0) = g''(0) .$$

By Taylor's integral formula

$$g(x) - g(0) - xg'(0) = \int_0^x g''(u)(x-u)du = x^2 \int_0^1 g''(sx)(1-s)ds .$$

Thus, $G(x) = \int g''(sx)d\nu(s)$, where $d\nu(s) = 2(1-s)\mathbb{1}\{0 \leq s \leq 1\}ds$.

Hence, if g is convex then $g'' \geq 0$ and $G \geq 0$. Moreover, if g'' is increasing then the functions $x \mapsto g''(sx)$ for $s \in [0, 1]$ are all increasing and G is also increasing as an average of increasing functions. For $g(u) = \exp(u)$, this yields that $(\exp(u) - u - 1)/u^2$ is increasing, as required for the proof of Bernstein's inequality.

Similarly, if g'' is convex then G is also convex as an average of convex functions ($x \mapsto g''(sx)$). Moreover, by Jensen's inequality applied to convex function $\psi(s) = g''(xs)$ with the probability measure $d\nu(s) = 2(1-s)\mathbb{1}\{0 \leq s \leq 1\}ds$

$$G(x) = \int_0^1 g''(xs) 2(1-s)ds \geq g'' \left(x \int_0^1 s \times 2(1-s)ds \right) = g'' \left(\frac{x}{3} \right) .$$

For $g(u) = (1+u) \log(1+u) - u$, $g''(u) = 1/(1+u)$ and this yields:

$$\frac{g(u)}{u^2/2} \geq g'' \left(\frac{u}{3} \right) = \frac{1}{1+u/3} .$$

Exercise: for $X_i \stackrel{iid}{\sim} \mathcal{B}(\mu)$, $\mathbb{P}(\bar{X}_n \geq 2\mu) \leq \exp(-n \times ?)$

Chernoff + Taylor: since $\log(u) \geq (u - 1)/u$,

$$\text{kl}(2\mu, \mu) = 2\mu \log(2) + (1 - 2\mu) \log \frac{1 - 2\mu}{1 - 2\mu} \geq 2\mu \log(2) - \mu = \mu(2 \log(2) - 1) \approx 0.386 \mu .$$

Chernoff with convexity:

$$\text{kl}(2\mu, \mu) \geq \frac{(2\mu - \mu)^2/2}{4/3\mu} = \frac{3}{8} \mu = 0.375 \mu .$$

Improved Hoeffding:

$$\text{kl}(2\mu, \mu) \geq \frac{(2\mu - \mu)^2/2}{\max_{\mu \leq u \leq 2\mu} u(1-u)} \geq \frac{\mu^2/2}{2\mu} = \frac{1}{4} \mu = 0.25 \mu .$$

Bennett:

$$2\mu \log \frac{2\mu}{\mu} - (2\mu - \mu) = \mu(2 \log(2) - 1) \approx 0.386 \mu .$$

Bernstein:

$$\frac{(2\mu - \mu)^2/2}{\mu(1 - \mu) + (2\mu - \mu)/3} \geq \frac{\mu^2/2}{\mu + \mu/3} \frac{3}{8} \mu = 0.375 \mu .$$

Hoeffding: $2(2\mu - \mu)^2 = 2\mu^2$, very poor (as expected) when μ is small.

Kullback-Leibler divergence

Kullback-Leibler divergence

Definition

Let P and Q be two probability distributions on a measurable set Ω . The Kullback-Leibler divergence from Q to P is defined as follows:

- if P is not absolutely continuous with respect to Q , then $\text{KL}(P, Q) = +\infty$;
- otherwise, let $\frac{dP}{dQ}$ be the Radon-Nikodym derivative of P with respect to Q . Then

$$\text{KL}(P, Q) = \int_{\Omega} \log \frac{dP}{dQ} dP = \int_{\Omega} \frac{dP}{dQ} \log \frac{dP}{dQ} dQ .$$

Property: $0 \leq \text{KL}(P, Q) \leq +\infty$, $\text{KL}(P, Q) = 0$ iff $P = Q$.


If $P \ll Q$ and $f = \frac{dP}{dQ}$, $\int_{\Omega} f \log(f) dQ = \int_{\Omega} [f \log(f)]_+ dQ - \int_{\Omega} [f \log(f)]_- dQ$, the later is finite since $[f \log(f)]_- \leq 1/e$.

Examples:

$$\text{KL}(\mathcal{B}(p), \mathcal{B}(q)) = \text{kl}(p, q), \text{KL}(\mathcal{N}(\mu_1, \sigma^2), \mathcal{N}(\mu_2, \sigma^2)) = \frac{(\mu_1 - \mu_2)^2}{2\sigma^2} .$$

Lower Bound: Change of Measure

For all $\epsilon > 0$ and all $\alpha > 0$,


$$\begin{aligned}\mathbb{P}_\mu(\bar{X}_n \geq x) &= \mathbb{E}_\mu[\mathbb{1}\{\bar{X}_n \geq x\}] \\ &= \mathbb{E}_{x+\epsilon} \left[\mathbb{1}\{\bar{X}_n \geq x\} \times \frac{d\mathbb{P}_\mu}{d\mathbb{P}_{x+\epsilon}}(X_1, \dots, X_n) \right] \\ &= \mathbb{E}_{x+\epsilon} \left[\mathbb{1}\{\bar{X}_n \geq x\} \times e^{-\sum_{i=1}^n \log \frac{dP_{x+\epsilon}}{dP_\mu}(X_i)} \right] \\ &\geq \mathbb{E}_{x+\epsilon} \left[\mathbb{1}\{\bar{X}_n \geq x\} \mathbb{1}\left\{ \frac{1}{n} \sum_{i=1}^n \log \frac{dP_{x+\epsilon}}{dP_\mu}(X_i) \leq \mathbb{E}_{x+\epsilon} \left[\log \frac{dP_{x+\epsilon}}{dP_\mu}(X_1) \right] + \alpha \right\} \right. \\ &\quad \left. \times e^{-\sum_{i=1}^n \log \frac{dP_{x+\epsilon}}{dP_\mu}(X_i)} \right] \\ &\geq e^{-n \left\{ \mathbb{E}_{x+\epsilon} \left[\log \frac{dP_{x+\epsilon}}{dP_\mu}(X_1) \right] + \alpha \right\}} \left[1 - \mathbb{P}_{x+\epsilon}(\bar{X}_n < x) \right. \\ &\quad \left. - \mathbb{P}_{x+\epsilon} \left(\frac{1}{n} \sum_{i=1}^n \log \frac{dP_{x+\epsilon}}{dP_\mu}(X_i) > \mathbb{E}_{x+\epsilon} \left[\log \frac{dP_{x+\epsilon}}{dP_\mu}(X_1) \right] + \alpha \right) \right] \\ &= e^{-n \{ \text{kl}(x+\epsilon, \mu) + \alpha \}} (1 - o_n(1)).\end{aligned}$$

Lower Bound: Change of Measure

For all $\epsilon > 0$ and all $\alpha > 0$,

$$\begin{aligned}\mathbb{P}_\mu(\bar{X}_n \geq x) &= \mathbb{E}_\mu[\mathbb{1}\{\bar{X}_n \geq x\}] \\ &\geq \mathbb{E}_{x+\epsilon} \left[\mathbb{1}\{\bar{X}_n \geq x\} \mathbb{1}\left\{ \frac{1}{n} \sum_{i=1}^n \log \frac{dP_{x+\epsilon}}{dP_\mu}(X_i) \leq \mathbb{E}_{x+\epsilon} \left[\log \frac{dP_{x+\epsilon}}{dP_\mu}(X_1) \right] + \alpha \right\} \right. \\ &\quad \left. \times e^{-\sum_{i=1}^n \log \frac{dP_{x+\epsilon}}{dP_\mu}(X_i)} \right] \\ &\geq e^{-n \left\{ \mathbb{E}_{x+\epsilon} \left[\log \frac{dP_{x+\epsilon}}{dP_\mu}(X_1) \right] + \alpha \right\}} \left[1 - \mathbb{P}_{x+\epsilon}(\bar{X}_n < x) \right. \\ &\quad \left. - \mathbb{P}_{x+\epsilon} \left(\frac{1}{n} \sum_{i=1}^n \log \frac{dP_{x+\epsilon}}{dP_\mu}(X_i) > \mathbb{E}_{x+\epsilon} \left[\log \frac{dP_{x+\epsilon}}{dP_\mu}(X_1) \right] + \alpha \right) \right] \\ &= e^{-n \{ \text{kl}(x+\epsilon, \mu) + \alpha \}} (1 - o_n(1)).\end{aligned}$$

Asymptotic Optimality (Large Deviation Lower Bound)

$$\liminf_n \frac{1}{n} \log \mathbb{P}_\mu(\bar{X}_n \geq x) \geq -\text{kl}(x, \mu).$$

Lower Bound: Change of Measure

For all $\epsilon > 0$ and all $\alpha > 0$,

$$\begin{aligned} \mathbb{P}_\mu (\bar{X}_n \geq x) &= \mathbb{E}_\mu \left[\mathbb{1} \{ \bar{X}_n \geq x \} \right] \\ &\geq \mathbb{E}_{x+\epsilon} \left[\mathbb{1} \{ \bar{X}_n \geq x \} \mathbb{1} \left\{ \frac{1}{n} \sum_{i=1}^n \log \frac{dP_{x+\epsilon}}{dP_\mu} (X_i) \leq \mathbb{E}_{x+\epsilon} \left[\log \frac{dP_{x+\epsilon}}{dP_\mu} (X_1) \right] + \alpha \right\} \right. \\ &\quad \left. \times e^{-\sum_{i=1}^n \log \frac{dP_{x+\epsilon}}{dP_\mu} (X_i)} \right] \\ &\geq e^{-n \left\{ \mathbb{E}_{x+\epsilon} \left[\log \frac{dP_{x+\epsilon}}{dP_\mu} (X_1) \right] + \alpha \right\}} \left[1 - \mathbb{P}_{x+\epsilon} (\bar{X}_n < x) \right. \\ &\quad \left. - \mathbb{P}_{x+\epsilon} \left(\frac{1}{n} \sum_{i=1}^n \log \frac{dP_{x+\epsilon}}{dP_\mu} (X_i) > \mathbb{E}_{x+\epsilon} \left[\log \frac{dP_{x+\epsilon}}{dP_\mu} (X_1) \right] + \alpha \right) \right] \\ &= e^{-n \{ \text{kl}(x+\epsilon, \mu) + \alpha \}} (1 - o_n(1)) . \end{aligned}$$

Asymptotic Optimality (Large Deviation Principle)

$$\frac{1}{n} \log \mathbb{P}_\mu (\bar{X}_n \geq x) \xrightarrow[n \rightarrow \infty]{} -\text{kl}(x, \mu) .$$

Properties of KL divergence

Tensorization of entropy:

If $P = P_1 \otimes P_2$ and $Q = Q_1 \otimes Q_2$, then

$$\text{KL}(P, Q) = \text{KL}(P_1, Q_1) + \text{KL}(P_2, Q_2) .$$

Contraction of entropy data-processing inequality:

Let (Ω, \mathcal{A}) be a measurable space, and let P and Q be two probability measures on (Ω, \mathcal{A}) . Let $X : \Omega \rightarrow (\mathcal{X}, \mathcal{B})$ be a random variable, and let P^X (resp. Q^X) be the push-forward measures, ie the laws of X wrt P (resp. Q). Then

$$\text{KL}(P^X, Q^X) \leq \text{KL}(P, Q) .$$

Pinsker's inequality:

Let $P, Q \in \mathfrak{M}_1(\Omega, \mathcal{A})$. Then

$$\|P - Q\|_{\text{TV}} \stackrel{\text{def}}{=} \sup_{A \in \mathcal{A}} |P(A) - Q(A)| \leq \sqrt{\frac{\text{KL}(P, Q)}{2}} .$$

Proof: contraction

Contraction: if $\text{KL}(P, Q) = +\infty$, the result is obvious. Otherwise, $P \ll Q$ and there exists $\frac{dP}{dQ} : \Omega \rightarrow \mathbb{R}$ such that for all measurable $f : \Omega \rightarrow \mathbb{R}$, $\int_{\Omega} f dP = \int_{\Omega} f \frac{dP}{dQ} dQ$.

- We first prove that $P^X \ll Q^X$ and, if $\gamma(x) := \mathbb{E}_Q \left[\frac{dP}{dQ} \mid X = x \right]$ is the Q -a.s. unique function such that $\mathbb{E}_Q \left[\frac{dP}{dQ} \mid X \right] = \gamma(X)$, then $\gamma = \frac{dP^X}{dQ^X}$. Indeed, for all $B \in \mathcal{B}$,

$$\begin{aligned} P^X(B) &= P(X \in B) = \int_{X \in B} \frac{dP}{dQ} dQ = \mathbb{E}_Q \left[\frac{dP}{dQ} \mathbb{1}\{X \in B\} \right] \\ &= \mathbb{E}_Q \left[\mathbb{E}_Q \left[\frac{dP}{dQ} \mathbb{1}\{X \in B\} \mid X \right] \right] = \mathbb{E}_Q \left[\mathbb{1}\{X \in B\} \mathbb{E}_Q \left[\frac{dP}{dQ} \mid X \right] \right] \\ &= \mathbb{E}_Q \left[\mathbb{1}\{X \in B\} \gamma(X) \right] = \int_{X \in B} \gamma(X) dQ = \int_B \gamma dQ^X \end{aligned}$$

and hence $P^X \ll Q^X$ and $\frac{dP^X}{dQ^X} = \gamma$.

- Now,

$$\begin{aligned} \text{KL}(P^X, Q^X) &= \int_{\mathcal{X}} \gamma \log \gamma dQ^X = \int_{\Omega} \gamma(X) \log \gamma(X) dQ \\ &= \mathbb{E}_Q \left[\phi \left(\mathbb{E}_Q \left[\frac{dP}{dQ} \mid X \right] \right) \right] \quad \text{where } \phi := x \mapsto x \log(x) \text{ is convex} \\ &\leq \mathbb{E}_Q \left[\mathbb{E}_Q \left[\phi \left(\frac{dP}{dQ} \right) \mid X \right] \right] \quad \text{by (conditional) Jensen's inequality} \\ &= \mathbb{E}_Q \left[\phi \left(\frac{dP}{dQ} \right) \right] = \text{KL}(P, Q). \end{aligned}$$

Let $A \in \mathcal{A}$, $p = P(A)$ and $q = Q(A)$. By contraction,

$$\text{KL}(P, Q) \geq \text{KL}(P^{\mathbb{1}_A}, Q^{\mathbb{1}_A}) = \text{KL}(\mathcal{B}(P(A)), \mathcal{B}(Q(A))) = \text{kl}(P(A), Q(A)) \geq 2(P(A) - Q(A))^2.$$

Lower Bound: the Entropic Way

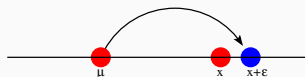
Let $\Omega = \{0, 1\}^n$, $X_i(\omega) = \omega_i$

Probability laws on Ω : $\mathbb{P}_p = \mathcal{B}(p)^{\otimes n}$.

For all $\epsilon > 0$,

$$\begin{aligned} n \text{kl}(x + \epsilon, \mu) &= \text{KL}(\mathbb{P}_{x+\epsilon}, \mathbb{P}_\mu) & \text{KL}(P \otimes P', Q \otimes Q') &= \text{KL}(P, Q) + \text{KL}(P', Q') \\ &\geq \text{KL}\left(\mathbb{P}_{x+\epsilon}^1\{\bar{X}_n \geq x\}, \mathbb{P}_\mu^1\{\bar{X}_n \geq x\}\right) & \text{KL}(P, Q) &\geq \text{KL}(P^X, Q^X) \\ &= \text{kl}\left(\mathbb{P}_{x+\epsilon}(\bar{X}_n \geq x), \mathbb{P}_\mu(\bar{X}_n \geq x)\right) & \text{contraction of entropy} \\ &\geq \mathbb{P}_{x+\epsilon}(\bar{X}_n \geq x) \log \frac{1}{\mathbb{P}_\mu(\bar{X}_n \geq x)} - \log(2) & = \text{data-processing inequality} \end{aligned}$$

$\text{kl}(p, q) \geq p \log \frac{1}{q} - \log 2$



A non-asymptotic lower bound

$$\forall \epsilon > 0, \quad \mathbb{P}_\mu(\bar{X}_n \geq x) \geq e^{-\frac{n \text{kl}(x+\epsilon, \mu) + \log(2)}{1 - e^{-2n\epsilon^2}}}.$$

No-Free-Lunch theorems: when learning is not possible

The No-Free-Lunch theorem

A learning algorithm A for binary classification maps a sample $S \sim \mathcal{D}^{\otimes n}$ to a decision rule \hat{h}_n .

Theorem

Let A be any learning algorithm for binary classification over a domain \mathcal{X} . If the training set size is $n \leq |\mathcal{X}|/2$, then there exists a distribution \mathcal{D} over $\mathcal{X} \times \{0, 1\}$ such that:

- there exists a function $f : \mathcal{X} \rightarrow \{0, 1\}$ with $L_{\mathcal{D}}(f) = 0$;
- with probability at least $1/7$ over the choice of $S \sim \mathcal{D}^{\otimes n}$,

$$L_{\mathcal{D}}(A(S)) \geq \frac{1}{8} .$$

Note that the ERM over $\mathcal{H} = \{f\}$, or over any set \mathcal{H} such that $n \geq 8 \log(7|\mathcal{H}|/6)$, is a successful learner in that setting.

Proof

Take $C \subset \mathcal{X}$ of cardinality $2n$, and $\{0, 1\}^C = \{f_1, \dots, f_T\}$ where $T = 2^{2n}$. For each $1 \leq i \leq T$, we denote by D_i the probability distribution on $C \times \{0, 1\}$ defined by

$$D_i(\{x, y\}) = \begin{cases} \frac{1}{2n} & \text{if } y = f_i(x), \\ 0 & \text{otherwise.} \end{cases}$$

We will show that $\max_{1 \leq i \leq T} \mathbb{E}[L_{D_i}(A(S))] \geq 1/4$, which entails the result thanks to the small lemma: if $P(0 \leq Z \leq 1) = 1$ and $\mathbb{E}[Z] \geq 1/4$, then $\mathbb{P}(Z \geq 1/8) \geq 1/7$. Indeed, $1/4 \leq \mathbb{E}[Z] \leq \mathbb{P}(Z < 1/8)/8 + \mathbb{P}(Z \geq 1/8) = 1/8 - 7\mathbb{P}(Z \geq 1/8)/8$.

All the X -samples S_1^X, \dots, S_k^X , for $k = (2n)^n$, are equally likely. For $1 \leq j \leq k$, if $S_j^X = (x_1, \dots, x_n)$ we denote by $S_j^i = ((x_1, f_i(x_1)), \dots, (x_n, f_i(x_n)))$, and $\hat{f}_j^i = A(S_j^i)$.

$$\begin{aligned} \max_{1 \leq i \leq T} \mathbb{E}[L_{D_i}(A(S))] &= \max_{1 \leq i \leq T} \frac{1}{k} \sum_{j=1}^k L_{D_i}(\hat{f}_j^i) \geq \frac{1}{T} \sum_{i=1}^T \frac{1}{k} \sum_{j=1}^k L_{D_i}(\hat{f}_j^i) \\ &= \frac{1}{k} \sum_{j=1}^k \frac{1}{T} \sum_{i=1}^T L_{D_i}(\hat{f}_j^i) \geq \frac{1}{k} \sum_{j=1}^k \min_{1 \leq i \leq T} \frac{1}{T} \sum_{i=1}^T L_{D_i}(\hat{f}_j^i). \end{aligned}$$

Fix $1 \leq j \leq k$, denote $S_j^X = (x_1, \dots, x_n)$ and define $\{v_1, \dots, v_p\} = C \setminus \{x_1, \dots, x_n\}$, where $p \geq n$. Then

$$L_{D_i}(\hat{f}_j^i) = \frac{1}{2n} \sum_{x \in C} \mathbb{1}\{\hat{f}_j^i(x) \neq f_i(x)\} \geq \frac{1}{2p} \sum_{r=1}^p \mathbb{1}\{\hat{f}_j^i(v_r) \neq f_i(v_r)\}$$

and hence

$$\frac{1}{T} \sum_{i=1}^T L_{D_i}(\hat{f}_j^i) \geq \frac{1}{T} \sum_{i=1}^T \frac{1}{2p} \sum_{r=1}^p \mathbb{1}\{\hat{f}_j^i(v_r) \neq f_i(v_r)\} \geq \frac{1}{2} \min_{1 \leq r \leq p} \frac{1}{T} \sum_{i=1}^T \mathbb{1}\{\hat{f}_j^i(v_r) \neq f_i(v_r)\}.$$

Fix $1 \leq r \leq p$. Then the functions $\{f_i : 1 \leq i \leq T\}$ can be grouped into $T/2$ pairs of functions $(\tilde{f}_i^0, \tilde{f}_i^1)$, $1 \leq i \leq T/2$ which agree on all $x \in C$ except on v_r , and for all $1 \leq i \leq T/2$ it holds that $\mathbb{1}\{\hat{f}_j^i(v_r) \neq \tilde{f}_i^0(v_r)\} + \mathbb{1}\{\hat{f}_j^i(v_r) \neq \tilde{f}_i^1(v_r)\} = 1$. Hence,

$$\sum_{i=1}^T \mathbb{1}\{\hat{f}_j^i(v_r) \neq f_i(v_r)\} = \sum_{i=1}^{T/2} \mathbb{1}\{\hat{f}_j^i(v_r) \neq \tilde{f}_i^0(v_r)\} + \mathbb{1}\{\hat{f}_j^i(v_r) \neq \tilde{f}_i^1(v_r)\} = T/2, \text{ which concludes the proof.}$$

Consequence: infinite VC-dimension \implies no learnability

Recall that a hypothesis class \mathcal{H} is *agnostic PAC learnable* if there exists a function $n_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$ and a learning algorithm $S \mapsto \hat{h}_n$ such that for every $\epsilon, \delta \in (0, 1)$, for every distribution D on $\mathcal{X} \times \mathcal{Y}$ when $S = ((X_1, Y_1), \dots, (X_n, Y_n)) \stackrel{iid}{\sim} D$,

$$\mathbb{P}\left(L_D(\hat{h}_n) \geq \min_{h' \in \mathcal{H}} L_D(h') + \epsilon\right) \leq \delta$$

for all $n \geq n_{\mathcal{H}}(\epsilon, \delta)$.

Theorem

Let \mathcal{H} be a class of infinite VC-dimension. Then \mathcal{H} is not PAC-learnable.

Proof: for every training size n , there exists a set $C \subset \mathcal{X}$ of size $2n$ that is shattered by \mathcal{H} . By the NFL theorem, for every learning algorithm A there exists a probability distribution D over $\mathcal{X} \times \{0, 1\}$ and $h : \mathcal{X} \rightarrow \{0, 1\}$ such that $L_D(h) = 0$ but with probability at least $1/7$ over the training set, we have $L_D(A(S)) \geq 1/8$.

Consequence: Curse of Dimensionality

Theorem

Let $c > 1$ be a Lipschitz constant. Let A be any learning algorithm for binary classification over a domain $\mathcal{X} = [0, 1]^d$. If the training set size is $n \leq (c + 1)^d/2$, then there exists a distribution \mathcal{D} over $[0, 1]^d \times \{0, 1\}$ such that:

- $\eta(x) = \mathbb{P}(Y = 1|X = x)$ is c -Lipschitz;
- the Bayes error of the distribution is 0;
- with probability at least $1/7$ over the choice of $S \sim \mathcal{D}^{\otimes n}$,

$$L_{\mathcal{D}}(A(S)) \geq \frac{1}{8}.$$