

Hoffding Impunity

JASD.

Oct 12, 2023.



# Foundations of Machine Learning Hoeffding Inequality

Oct 12, 2023

Lova RALAIOLA l.ralaiola@criteo.com  
Criteo AI LAB.

Created in 2018.

## I. What do we want to achieve?

Uniform generalization bounds, i.e. bounds like:

With high probability,  $\forall f \in \mathcal{F}$

$$\text{Generalization}(f) \leq \text{Empirical error}(f, S) + \mathcal{E}(n, \beta, \dots)$$

"Capacity" of  $\mathcal{F}$

Suppose to  
increase when  $n \rightarrow \infty$ .

## Typical setting of Statistical learning theory.

- $(X, Y) \sim D$ ,  $D$  fixed and unknown distribution

- $S = \{(X_i, Y_i)\}_{i=1}^n$ , where  $(X_i, Y_i)$  are  $n$  independent copies of  $(X, Y)$

In other words, the training sample is an IID sample where IID stands for  
Identically and Independently Distributed.

- Ultimate goal / criterion.

We want to minimize

$$R_\ell(f) := \mathbb{E}_{x,y \sim D} [\ell(f(x), y)]$$

↓  
 ℓ-risk      learned predictor  
 ↑  
 loss function

Examples: binary classifier

- $F \subseteq \{-1, +1\}^X$
- $\ell(f(x), y) := \begin{cases} 1 & f(x) \neq y \\ 0 & f(x) = y \end{cases}$

or

$$\begin{aligned} & F \subseteq \mathbb{R}^X \\ & \ell(f(x), y) := \begin{cases} 1 & yf(x) \leq 0 \\ 0 & yf(x) > 0 \end{cases} \end{aligned}$$

- As we do not know  $D$ , we use  $S$  to "gather" information. The bound we look for is then:

with prob 1 - δ

$$\forall f \in F \quad R_\ell(f) \leq \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) + \varepsilon(n, \delta, F)}_{\downarrow}$$

Empirical ℓ-risk,  
the traditional quantity  
you look at when growing.

- What appears here is the connection between an empirical quantity and its expectation... "concret"

## III. Hoeffding Inequality 1963.

The  $x_1 - x_n$  independent r.v.

such that  $\forall i \in \{1 - n\} \exists a_i, b_i: P(a_i \leq x_i \leq b_i) = 1$ .

let  $S_n := \sum_{i=1}^n x_i$  then

$$P(S_n - E S_n \geq \varepsilon) \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

$$P(E S_n - S_n \geq \varepsilon) \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

Proposition:  $x_1 - x_n$  IID and th:  $P(0 \leq x_i \leq 1) = 1$

If  $\mu := E x_1 (= E x_2 = \dots = E x_n)$

We have.  $P\left(\frac{1}{n} \sum_{i=1}^n x_i - \mu \geq \varepsilon\right) \leq \exp(-2n\varepsilon^2)$

$$P\left(\mu - \frac{1}{n} \sum_{i=1}^n x_i \geq \varepsilon\right) \leq \exp(-2n\varepsilon^2)$$

Ex: Biased coin toss. how to estimate the bias with prob  $1-\delta$ .

$x_i = 1$  if tail  
 $0$  if head

$$\textcircled{2} \quad P\left(\left|\frac{1}{n} \sum_{i=1}^n x_i - \mu\right| \geq \varepsilon\right) \leq 2\exp(-2n\varepsilon^2)$$

bias

② To get our estimate  $\hat{\mu}$ , it suffices that

$$2\exp(-L_n \varepsilon^2) \leq \delta \Leftrightarrow \exp(-L_n \varepsilon^2) \leq \frac{\delta}{2}$$

$$\Leftrightarrow -L_n \varepsilon^2 \leq \log \frac{\delta}{2}$$

$$\Leftrightarrow L_n \varepsilon^2 \geq \log \frac{2}{\delta}$$

$$\Leftrightarrow \varepsilon^2 \geq \frac{1}{2L_n} \log \frac{2}{\delta}$$

$$\Leftrightarrow \varepsilon \geq \sqrt{\frac{1}{2L_n} \log \frac{2}{\delta}} \quad (\varepsilon \geq 0)$$

$$\text{with prob } 1-\delta \quad \mu \in \left[ \frac{1}{n} \sum x_i \pm \sqrt{\frac{1}{2L_n} \log \frac{2}{\delta}} \right]$$

Proof of the prop: do it, it's easy.

Remark: there is another concept note: McDiarmid's inequality that is often used in S2T. It generalizes Hoeffding's Inequality to the case where the function we look at is more complex than the mere average.

$\Rightarrow$  It provides an inequality in the form

$$P(\Phi(x_1, \dots, x_n) - E\Phi(x_1, \dots, x_n) \geq \varepsilon) \leq \exp\left(-\frac{2\varepsilon^2}{\sum c_i^2}\right)$$

Proof of Hoeffding Inequality:

- Markov inequality
- Laplace transform / exp transform
- Hoeffding lemma
- Close out (easy).

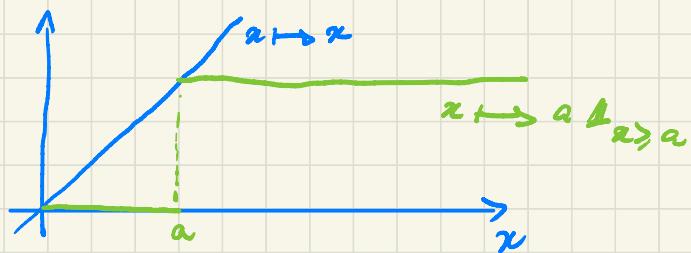
$$\begin{aligned} & \text{where } b_i \\ & \sup_{x_1, x_1'} |\Phi(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n) \\ & \quad - \Phi(x_1, \dots, x_{i-1}', x_i, x_{i+1}, \dots, x_n)| \\ & \leq c_i \end{aligned}$$

## Lemma: Markov Inequality

Let  $X$  be a rv taking nonnegative values ( $P(X \geq 0) = 1$ ) such that  $E[X] < +\infty$

$$\forall a > 0 \quad P(X \geq a) \leq \frac{E[X]}{a}$$

Proof:



$$\forall x \quad x \geq a \cdot 1_{x \geq a}$$

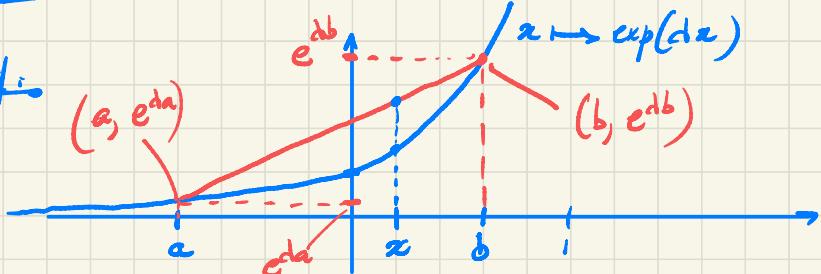
$$\begin{aligned} \Rightarrow E[X] &\geq E[a \cdot 1_{X \geq a}] = a \cdot E[1_{X \geq a}] \\ &= a \cdot P(X \geq a) \quad \blacksquare \end{aligned}$$

## Hoeffding's lemma

$X$  rv /  $E[X=0]$ ,  $\exists a, b$ ,  $P(a \leq X \leq b) = 1$

$$\forall d \in \mathbb{R}^+ \quad E[e^{dX}] \leq \exp\left(\frac{d^2(b-a)^2}{8}\right)$$

Proof:



$$\forall x \in [a, b] \quad e^{dx} \leq \frac{b-x}{b-a} e^{da} + \frac{x-a}{b-a} e^{db}$$

$$\Rightarrow \mathbb{E}[e^{dx}] \leq \mathbb{E}\left[\frac{b-a}{b-a} e^{da} + \frac{x-a}{b-a} e^{db}\right]$$

$$\Leftrightarrow \mathbb{E}[e^{dx}] \leq \frac{b}{b-a} e^{da} - \frac{a}{b-a} e^{db}$$

!

$$\begin{aligned} \Leftrightarrow \mathbb{E}[e^{dx}] &\leq e^{d(b-a)} \frac{a}{b-a} \left(1 + \frac{a}{b-a} - \frac{a}{b-a} e^{d(b-a)}\right) \\ &= e^{dp} \left(1 + p - pe^h\right) \text{ where } \begin{cases} h := d(b-a) \\ p := \frac{a}{b-a} \end{cases} \\ &= e^{L(h)} \text{ where } L(h) := -hp + \log(1 + p + pe^h) \end{aligned}$$

Taylor expansion:  $\exists t \in [0, 1] \quad L(h) = L(0) + hL'(0) + \frac{1}{2} h^2 L''(t)$

with:  $L(0) = 0$

$L'(0) = 0$

$L''(t) = \dots = t(1-t)$  with  $t \in [0, 1]$

$\leq \frac{1}{4}$ .



Hence:  $\forall h, L(h) = 0 + 0 + \frac{1}{2} h^2 L''(t) \leq \frac{1}{4}$

$\leq \frac{1}{8} h^2$

$= \frac{1}{8} d^2 (b-a)^2$

■

Lemma:  $X_1, \dots, X_n$  independent r.v.

$$\mathbb{E}\left[\prod_{i=1}^n X_i\right] = \prod_{i=1}^n \mathbb{E}X_i$$

• Proof of Hoeffding's inequality.

$$\text{P} \left( S_n - \mathbb{E}[S_n] \geq \varepsilon \right) = \text{P} \left( \exp \left( d(S_n - \mathbb{E}S_n) \right) \geq \exp(d\varepsilon) \right)$$

$$= \text{P} \left( \exp \left( d \left( \sum_{i=1}^n X_i - \mathbb{E} \sum_{i=1}^n X_i \right) \right) \geq \exp(d\varepsilon) \right)$$

$$= \text{P} \left( \exp \left[ d \sum_{i=1}^n (X_i - \mathbb{E}X_i) \right] \geq \exp(d\varepsilon) \right)$$

$$\left| \begin{array}{l} Z_i := X_i - \mathbb{E}X_i \\ \mu_i := \mathbb{E}X_i \end{array} \quad \text{P} \left( Z_i \in \left[ \underbrace{a_i - \mu_i}_{< 0}, \underbrace{b_i - \mu_i}_{> 0 \atop \text{}} \right] \right) = 1 \right\}$$

$$= \text{P} \left( \exp \left( d \sum_{i=1}^n Z_i \right) \geq \exp(d\varepsilon) \right)$$

$$\leq \mathbb{E} \left[ \exp \left[ d \sum_{i=1}^n Z_i \right] \right] / \exp(d\varepsilon) \quad (\text{Markov})$$

$$= \mathbb{E} \left[ \prod_{i=1}^n \exp(dZ_i) \right] \times \exp(-d\varepsilon)$$

Independent

$$= \prod_{i=1}^n \mathbb{E} \exp(dZ_i) \cdot \exp(-d\varepsilon)$$

$$\prod_{i=1}^n \exp \left( \frac{d^2}{8} ((b_i - \mu_i) - (a_i - \mu_i))^2 \right) \exp(-d\varepsilon)$$

or, equivalently:

$$= \exp \left( \frac{d^2}{8} \sum (b_i - a_i)^2 - d\varepsilon \right)$$

$$Z_i := \mathbb{E}[\phi(\cdot) | X_1 \dots X_i]$$

$$- \mathbb{E}[\phi(\cdot) | X_1 \dots X_{i-1}] \quad \text{minimizing wrt } d$$

$$g(d) = \frac{d^2}{8} \sum_{i=1}^n (b_i - a_i)^2 - d\varepsilon$$

$$g'(d) = \frac{1}{4} d \sum_{i=1}^n (b_i - a_i)^2 - \varepsilon$$

$$g'(d) = 0 \Leftrightarrow d = \frac{4\varepsilon}{\sum (b_i - a_i)^2}$$

$$= \exp \left( - \frac{2\varepsilon^2}{\sum (b_i - a_i)^2} \right)$$

In the case of McDiarmid's Inequality

We use:

$$\mathbb{E}[\phi(x_1 \dots x_n)]$$

$$= \mathbb{E}_{X_1 \dots X_{n-1}} \mathbb{E}_{x_n} [\phi(\cdot) | X_1 \dots X_{n-1}]$$

$$= \exp \left( \frac{d^2}{8} \sum (b_i - a_i)^2 - d\varepsilon \right)$$

$$Z_i := \mathbb{E}[\phi(\cdot) | X_1 \dots X_i]$$

$$- \mathbb{E}[\phi(\cdot) | X_1 \dots X_{i-1}] \quad \text{minimizing wrt } d$$

$$g(d) = \frac{d^2}{8} \sum_{i=1}^n (b_i - a_i)^2 - d\varepsilon$$

$$g'(d) = \frac{1}{4} d \sum_{i=1}^n (b_i - a_i)^2 - \varepsilon$$

$$g'(d) = 0 \Leftrightarrow d = \frac{4\varepsilon}{\sum (b_i - a_i)^2}$$

such that

$$\frac{\phi(x_1 \dots x_n)}{\sum_{i=1}^n Z_i} = \mathbb{E}[\phi(\cdot) | X_1 \dots X_n]$$

$$- \mathbb{E}[\phi(x_1 \dots x_n)]$$

$$= \exp \left( - \frac{2\varepsilon^2}{\sum (b_i - a_i)^2} \right)$$

