

14 Mar 2022

# Chernoff Bound Applications

## Announcements:

1. Next quiz is a week from Weds., i.e. March 23.
2. Practice problems to be distributed by Thurs.
3. Topics: Gaussians and Chernoff bounds.
4. Lectures from last week have been typeset; see "Lectures" page on course website.

## Applications of Chernoff Bounds

Suppose we can draw i.i.d. samples from  $[0, M]$

call them  $Y_1, Y_2, \dots$

and we want an estimate  $\hat{Y}$  s.t.

$$\Pr\left(|\hat{Y} - \mathbb{E}Y| > \epsilon\right) < \delta \quad \leftarrow \begin{array}{l} \text{"probably approximately} \\ \text{correct"} \\ \text{(PAC)} \end{array}$$

$\mathbb{E}Y = \mathbb{E}[Y_i]$  for any  $i$

If we just use  $\hat{Y} = \frac{1}{N}(Y_1 + \dots + Y_N)$  how large must  $N$  be to guarantee this?

"Sample complexity of  $(\epsilon, \delta)$ -PAC estimation"

First set  $X_i = Y_i/M$ .  $\hat{X} = \hat{Y}/M$ .

$$|\hat{X} - \mathbb{E}X| > \epsilon/M \iff |\hat{Y} - \mathbb{E}Y| > \epsilon$$

$$\mathbb{E}X_i \iff |X_1 + \dots + X_N - \mathbb{E}(X_1 + \dots + X_N)| > \frac{\epsilon N}{M}$$

Aside: Another inequality called Hoeffding bound, with almost same proof, says

$$\Pr(|X_1 + \dots + X_n - \mathbb{E}(X_1 + \dots + X_n)| > \lambda) < 2 \exp\left(-\frac{2\lambda^2}{n}\right)$$

assuming  $X_i$  is  $[0,1]$ -valued

We're applying Hoeffding with  $\lambda = \frac{\epsilon N}{M}$  and we want

$$2 \exp\left(-\frac{2\lambda^2}{n}\right) < \delta$$

$\Leftrightarrow$

$$-\frac{2\lambda^2}{n} < \ln\left(\frac{\delta}{2}\right)$$

$\Leftrightarrow$

(multiply by  $-\frac{1}{2}$ )

$$\frac{\lambda^2}{n} > \frac{1}{2} \ln\left(\frac{2}{\delta}\right)$$

$\Leftrightarrow$

$$\frac{\epsilon^2 N}{M^2} > \frac{1}{2} \ln\left(\frac{2}{\delta}\right)$$

"relative range" squared

$\Leftrightarrow$

$$N >$$

$$\frac{M^2}{2\epsilon^2} \ln\left(\frac{2}{\delta}\right)$$

log of inverse failure probability

To become 10 times more certain of success, average in another 1.2  $M^2/\epsilon^2$  samples.

# Empirical Risk Minimization

We have data samples  $z_1, \dots, z_N$ .

hypotheses  $h_1, \dots, h_M$ .

loss function  $L(h, z) \in [0, 1]$

$L(h, z)$  specifies how badly hypothesis  $h$  explains data point  $z$ .

Example:  $z_i = (\phi, y)$   $\phi$  = feature vector  
 $y$  = label

$h$ : a function from feature vectors to labels

$$L(h, z) = \begin{cases} 0 & \text{if } y = h(\phi) \\ 1 & \text{if } y \neq h(\phi) \end{cases}$$

ERM is the algorithm that outputs hypothesis

$h_{\text{ERM}} \in \{h_1, \dots, h_M\}$  that minimizes

$$\frac{1}{N} \sum_{j=1}^N L(h, z_{ij})$$

Goal. Analyze generalization error of ERM.

IF  $z_1, \dots, z_N$  are i.i.d. samples from a distribution on  $Z$ 's, and  $Z$  is one more sample (unseen during training)

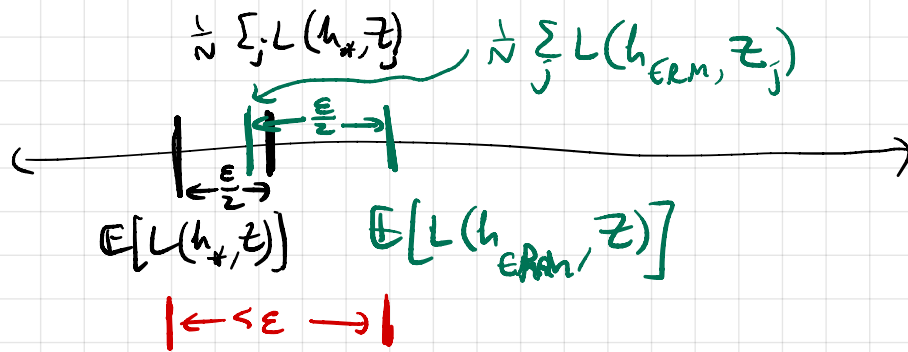
how large must  $N$  be so that with probability at least  $1 - \delta$ ,

$$(*) \quad \mathbb{E}[L(h_{ERM}, Z)] \leq \epsilon + \min_{i \in [m]} \mathbb{E}[L(h_i, Z)].$$

Let  $h_*$  be the minimizer on the RHS above.

One way to guarantee that inequality  $(*)$  happens is to insist that  $\forall i \in [m]$ ,

$$(**) \quad \left| \frac{1}{n} \sum_{j=1}^n L(h_i, Z_j) - \mathbb{E}[L(h_i, Z)] \right| < \frac{\epsilon}{2}.$$



How large must  $n$  be so that  $(**)$  holds with probability  $\geq 1 - \delta$  for any specific  $i$ ?

$$\frac{1}{2\epsilon^2} \ln\left(\frac{2M}{\delta}\right) = \frac{2}{\epsilon^2} \ln\left(\frac{2M}{\delta}\right).$$

Drive failure probability down to  $\frac{\delta}{m}$  for each  $i$  by using  $\frac{2}{\epsilon^2} \ln\left(\frac{2M}{\delta}\right)$  samples.

$$\begin{aligned} \text{Now, } & \mathbb{E}[\#i \text{ for which } (**) \text{ fails to hold}] \\ & \leq \sum_{i=1}^m \Pr(**) \text{ violated for } i \\ & \leq m \cdot (\delta/m) = \delta \end{aligned}$$