

LECTURE 5

THM (Hoeffding's Inequality) Let X_1, \dots, X_N be iid symmetric Bernoulli r.v's
i.e. $P\{X_i = 1\} = P\{X_i = -1\} = 1/2$. Then

$$P\left\{\frac{1}{\sqrt{N}} \sum_{i=1}^N X_i \geq t\right\} \leq \exp\left(-\frac{t^2}{2}\right) \quad \forall t \geq 0.$$

Proof: "The MGF method" (moment generating function).

• Let $\lambda \geq 0$ be a parameter TBD later ("to be determined").

• Multiply both sides of the inequality by $\lambda \sqrt{N}$ and exponentiate:

$$P\left\{\frac{1}{\sqrt{N}} \sum_{i=1}^N X_i \geq t\right\} = P\left\{\exp\left(\lambda \sum_{i=1}^N X_i\right) \geq \exp(t\lambda \sqrt{N})\right\}$$

Markov's ineq. $\leq \frac{\mathbb{E} \exp\left(\lambda \sum_{i=1}^N X_i\right)}{\exp(t\lambda \sqrt{N})} = e^{-t\lambda \sqrt{N}} \mathbb{E} \prod_{i=1}^N e^{\lambda X_i}$
(independence)

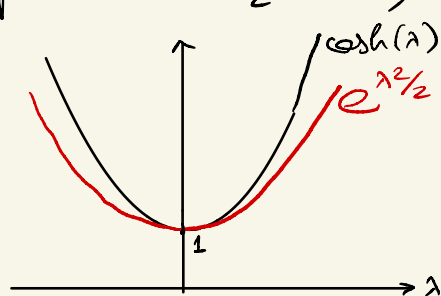
$= e^{-t\lambda \sqrt{N}} \prod_{i=1}^N \mathbb{E} e^{\lambda X_i}$
all terms are the same (identical distr.)

$= e^{-t\lambda \sqrt{N}} \left(\mathbb{E} e^{\lambda X}\right)^N \text{ where } X \sim \text{sym Bernoulli}$

• Since X takes values 1 and -1 with probabilities $1/2$ each,

$$\mathbb{E} e^{\lambda X} = \frac{1}{2} e^{\lambda} + \frac{1}{2} e^{-\lambda}$$

"Moment generating function" (MGF) $= \cosh(\lambda) \stackrel{\text{Taylor}}{\leq} e^{\lambda^2/2}$ HW



$\leq e^{-t\lambda \sqrt{N}} \left(e^{\lambda^2/2}\right)^N = \exp\left(-\underbrace{t\lambda \sqrt{N}}_{\mu} + \frac{\lambda^2 N}{2}\right) = \exp\left(-t\mu + \frac{\mu^2}{2}\right).$

• Minimize in $\mu \Rightarrow$ for $t = \mu$:

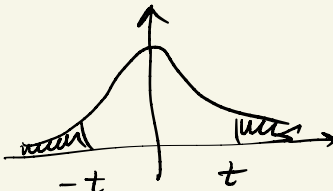
$\leq e^{-t^2/2}$. QED

Example: In N flips of a fair coin, $P\{\text{at least } \frac{3}{4}N \text{ heads}\} = \exp(-\frac{N}{8})$

Proof: #heads = $\sum_{i=1}^N Z_i$ where $Z_i = \begin{cases} 1 & \text{if } i\text{th flip} = \text{head} \\ 0 & \text{if tail} \end{cases} \sim \text{Ber}(1/2)$

Create symmetric r.v.s: $X_i = 2Z_i - 1 = \begin{cases} 1 & \text{if tail} \\ -1 & \text{if head} \end{cases} \sim \text{Sym Ber}$

$$P\{\sum_{i=1}^N Z_i \geq \frac{3N}{4}\} = P\{\underbrace{\frac{1}{\sqrt{N}} \sum_{i=1}^N (2Z_i - 1)}_{\substack{\text{Koeffding} \\ \uparrow \\ t}} \geq \frac{\sqrt{N}}{2}\} \stackrel{\text{Koeffding}}{\leq} \exp(-\frac{(\sqrt{N}/2)^2}{2}) = \exp(-\frac{N}{8}) \quad \square$$

Remark 1 Two sided bound: 

$$P\left\{ \left| \frac{1}{\sqrt{N}} \sum_{i=1}^N X_i \right| \geq t \right\} = \underbrace{P\left\{ \frac{1}{\sqrt{N}} \sum_{i=1}^N X_i > t \right\}}_{\substack{\text{Koeffding} \uparrow \\ e^{-t^2/2}}} + \underbrace{P\left\{ \frac{1}{\sqrt{N}} \sum_{i=1}^N X_i < -t \right\}}_{\substack{\text{Koeffding} \uparrow \\ e^{-t^2/2}}}$$

$$\leq \boxed{2e^{-t^2/2}}$$

$P\left\{ \frac{1}{\sqrt{N}} \sum_{i=1}^N (-X_i) > t \right\}$
Sym Bernoulli, too

Remark 2 Same MGF method \Rightarrow

HW \curvearrowright

THM (General Koeffding) Let X_1, \dots, X_n be independent r.v.s such that $X_i \in [a_i, b_i] \quad \forall i$. Then $S_N = \sum_{i=1}^N X_i$ satisfies

$$P\{S_N - \mathbb{E}S_N \geq t\} \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^N (b_i - a_i)^2}\right)$$

- This contains Koeffding we proved earlier (check!)
- And it works for a biased coin, i.e. for $X_i \sim \text{Ber}(p) \Rightarrow S_N = \text{Binom}(N, p)$

Ex. (a) $P\{\frac{1}{4}\text{-biased coin comes up } \geq \frac{N}{2} \text{ heads}\} \leq \exp(-\frac{N}{8})$
 (b) More generally, $S_N \sim \text{Binom}(N, p)$ satisfies $P\{|S_N - pN| \geq \delta N\} \leq 2\exp(-2\delta^2 N)$ (HW)
 $\forall \delta \geq 0$

APPLICATION: MEAN ESTIMATION

- Problem: Estimate the mean μ of a distribution from a sample X_1, \dots, X_n , which is drawn independently from this distribution.

- Classical estimator: the sample mean

$$\hat{\mu} := \frac{1}{N} \sum_{i=1}^N X_i$$

- $E\hat{\mu} = \mu$, i.e. unbiased. Mean-squared error:

$$E(\hat{\mu} - \mu)^2 = \text{Var}(\hat{\mu}) = \frac{1}{N^2} \sum_{i=1}^N \text{Var}(X_i) = \frac{\sigma^2}{N}$$

\Rightarrow the error is $\frac{\sigma}{\sqrt{N}}$ on average

- Confidence interval? Chebyshev's inequality:

$$P\left\{|\hat{\mu} - \mu| \geq \frac{t\sigma}{\sqrt{N}}\right\} \leq \frac{\sigma^2/N}{(t\sigma/\sqrt{N})^2} = \frac{1}{t^2} \quad \text{Meh.}$$

- For Gaussian distribution, we can get a much higher confidence:

$$\text{if } X_i \sim N(\mu, \sigma^2) \Rightarrow \hat{\mu} \sim N\left(\mu, \frac{\sigma^2}{N}\right).$$

$$\text{Gaussian tail} \Rightarrow P\left\{|\hat{\mu} - \mu| \geq \frac{t\sigma}{\sqrt{N}}\right\} \leq \boxed{\exp(-t^2/2)} \quad \text{☺}$$

↑ Gaussian tail (lec. 4)

- Can we get a similar confidence for general distr.?

SURPRISINGLY, YES!

Thm $\forall 0 \leq t \leq \sqrt{N}$, \exists estimator $\tilde{\mu} = \tilde{\mu}_t(X_1, \dots, X_N)$ that satisfies

$$P\left\{|\tilde{\mu} - \mu| \geq \frac{t\sigma}{\sqrt{N}}\right\} \leq 2\exp(-ct^2)$$

if X_i 's are drawn from a distribution with mean μ and variance σ^2 .
Here $c > 0$ is an absolute constant.

• With only finite variance assumption, Gaussian power is surprising!

Proof $\tilde{\mu} =$ "Median-of-Means" (MoM):

- Partition the sample into K blocks of equal size M $N = MK$
(For simplicity, assume N is divisible by M . General case \rightarrow KW)

$$\underbrace{X_1 X_2 \dots X_M}_{B_1} \quad \underbrace{X_{M+1} \dots X_{2M}}_{B_2} \quad \dots \quad \underbrace{\dots X_N}_{B_K}$$

- Take the mean of each block; then take the median of those means:

$$(*) \quad \hat{\mu}_j := \frac{1}{M} \sum_{i \in B_j} X_i \quad ; \quad \tilde{\mu} := \text{Med}(\hat{\mu}_1, \dots, \hat{\mu}_K)$$

(i.e. at least $\frac{K}{2}$ of $\hat{\mu}_j$ are $\leq \tilde{\mu}$ and at least $\frac{K}{2}$ are $\geq \tilde{\mu}$)

- Accuracy of each $\hat{\mu}_j$? Mean μ , variance $\sigma^2/M \Rightarrow$

$$P\left\{\hat{\mu}_j \geq \mu + \frac{t\sigma}{\sqrt{N}}\right\} \leq \frac{\sigma^2/M}{\left(\frac{t\sigma}{\sqrt{N}}\right)^2} = \frac{N}{t^2 M} = \frac{K}{t^2} = \frac{1}{4}$$

Chebyshev

if we set $K := t^2/4$; note that $K \leq N$ by assumption

- Therefore

def of median

(For simplicity, assume $t/4$ is an integer.)
General case \rightarrow KW

$$P\left\{\hat{\mu} \geq \mu + \frac{t\sigma}{\sqrt{N}}\right\} \stackrel{\text{def of median}}{=} P\left\{\text{at least } \frac{K}{2} \text{ of } \hat{\mu}_j \text{'s are } \geq \mu + \frac{t\sigma}{\sqrt{N}}\right\}$$

$$= P\left\{\text{at least a half of } \frac{1}{4}\text{-biased coins come up heads}\right\}$$

$$\leq \exp(-K/8) = \exp(-t^2/32)$$

(Ex. \uparrow on the bottom of p.2) Q.E.D.

Remarks: The restriction $t \leq \sqrt{N}$ is necessary.

\hookrightarrow KW

- ② Can one construct the same estimator $\tilde{\mu}$ for all t simultaneously?
Not always, see [Gabor Lugosi: median-of-means tournaments].