

Let $a_1, \dots, a_n \in \mathbb{R}$ and let $\varepsilon_1, \dots, \varepsilon_n$ be i.i.d. Rademacher random variables: $\mathbb{P}(\varepsilon_i = 1) = \mathbb{P}(\varepsilon_i = -1) = 0.5$.

Theorem 7.1. [Hoeffding] For $t \geq 0$,

$$\mathbb{P}\left(\sum_{i=1}^n \varepsilon_i a_i \geq t\right) \leq \exp\left(-\frac{t^2}{2\sum_{i=1}^n a_i^2}\right).$$

Proof. Similarly to the proof of Bennett's inequality (Lecture 5),

$$\mathbb{P}\left(\sum_{i=1}^n \varepsilon_i a_i \geq t\right) \leq e^{-\lambda t} \mathbb{E} \exp\left(\lambda \sum_{i=1}^n \varepsilon_i a_i\right) = e^{-\lambda t} \prod_{i=1}^n \mathbb{E} \exp(\lambda \varepsilon_i a_i).$$

Using inequality $\frac{e^x + e^{-x}}{2} \leq e^{x^2/2}$ (from Taylor expansion), we get

$$\mathbb{E} \exp(\lambda \varepsilon_i a_i) = \frac{1}{2} e^{\lambda a_i} + \frac{1}{2} e^{-\lambda a_i} \leq e^{\frac{\lambda^2 a_i^2}{2}}.$$

Hence, we need to minimize the bound with respect to $\lambda > 0$:

$$\mathbb{P}\left(\sum_{i=1}^n \varepsilon_i a_i \geq t\right) \leq e^{-\lambda t} e^{\frac{\lambda^2}{2} \sum_{i=1}^n a_i^2}.$$

Setting derivative to zero, we obtain the result. □

Now we change variable: $u = \frac{t^2}{2\sum_{i=1}^n a_i^2}$. Then $t = \sqrt{2u \sum_{i=1}^n a_i^2}$.

$$\mathbb{P}\left(\sum_{i=1}^n \varepsilon_i a_i \geq \sqrt{2u \sum_{i=1}^n a_i^2}\right) \leq e^{-u}$$

and

$$\mathbb{P}\left(\sum_{i=1}^n \varepsilon_i a_i \leq \sqrt{2u \sum_{i=1}^n a_i^2}\right) \geq 1 - e^{-u}.$$

Here $\sum_{i=1}^n a_i^2 = \text{Var}(\sum_{i=1}^n \varepsilon_i a_i)$.

Rademacher sums will play important role in future. Consider again the problem of estimating $\frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f$. We will see that by the Symmetrization technique,

$$\frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f \sim \frac{1}{n} \sum_{i=1}^n f(X_i) - \frac{1}{n} \sum_{i=1}^n f(X'_i).$$

In fact,

$$\mathbb{E} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f \right| \leq \mathbb{E} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \frac{1}{n} \sum_{i=1}^n f(X'_i) \right| \leq 2 \mathbb{E} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f \right|.$$

The second inequality above follows by adding and subtracting $\mathbb{E}f$:

$$\begin{aligned} \mathbb{E} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \frac{1}{n} \sum_{i=1}^n f(X'_i) \right| &\leq \mathbb{E} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f \right| + \mathbb{E} \left| \frac{1}{n} \sum_{i=1}^n f(X'_i) - \mathbb{E}f \right| \\ &= 2 \mathbb{E} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f \right| \end{aligned}$$

while for the first inequality we use Jensen's inequality:

$$\begin{aligned} \mathbb{E} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f \right| &= \mathbb{E} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \frac{1}{n} \sum_{i=1}^n \mathbb{E}f(X'_i) \right| \\ &\leq \mathbb{E}_X \mathbb{E}_{X'} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \frac{1}{n} \sum_{i=1}^n \mathbb{E}f(X'_i) \right|. \end{aligned}$$

Note that $\frac{1}{n} \sum_{i=1}^n f(X_i) - \frac{1}{n} \sum_{i=1}^n \mathbb{E}f(X'_i)$ is equal in distribution to $\frac{1}{n} \sum_{i=1}^n \varepsilon_i (f(X_i) - f(X'_i))$.

We now prove Hoeffding-Chernoff Inequality:

Theorem 7.2. *Assume $0 \leq X_i \leq 1$ and $\mu = \mathbb{E}X$. Then*

$$\mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n X_i - \mu \geq t \right) \leq e^{-n\mathcal{D}(\mu+t, \mu)}$$

where the KL-divergence $\mathcal{D}(p, q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$.

Proof. Note that $\phi(x) = e^{\lambda x}$ is convex and so $e^{\lambda x} = e^{\lambda(x \cdot 1 + (1-x) \cdot 0)} \leq x e^{\lambda} + (1-x) e^{\lambda \cdot 0} = 1 - x + x e^{\lambda}$. Hence,

$$\mathbb{E}e^{\lambda X} = 1 - \mathbb{E}X + \mathbb{E}X e^{\lambda} = 1 - \mu + \mu e^{\lambda}.$$

Again, we minimize the following bound with respect to $\lambda > 0$:

$$\begin{aligned} \mathbb{P} \left(\sum_{i=1}^n X_i \geq n(\mu + t) \right) &\leq e^{-\lambda n(\mu+t)} \mathbb{E}e^{\lambda \sum X_i} \\ &= e^{-\lambda n(\mu+t)} (\mathbb{E}e^{\lambda X})^n \\ &\leq e^{-\lambda n(\mu+t)} (1 - \mu + \mu e^{\lambda})^n \end{aligned}$$

Take derivative w.r.t. λ :

$$-n(\mu + t)e^{-\lambda n(\mu+t)}(1 - \mu + \mu e^{\lambda})^n + n(1 - \mu + \mu e^{\lambda})^{n-1} \mu e^{\lambda} e^{-\lambda n(\mu+t)} = 0$$

$$-(\mu + t)(1 - \mu + \mu e^{\lambda}) + \mu e^{\lambda} = 0$$

$$e^{\lambda} = \frac{(1 - \mu)(\mu + t)}{\mu(1 - \mu - t)}.$$

Substituting,

$$\begin{aligned} \mathbb{P} \left(\sum_{i=1}^n X_i \geq n(\mu + t) \right) &\leq \left(\left(\frac{\mu(1 - \mu - t)}{(1 - \mu)(\mu + t)} \right)^{\mu+t} \left(1 - \mu + \frac{(1 - \mu)(\mu + t)}{1 - \mu - t} \right) \right)^n \\ &= \left(\left(\frac{\mu}{\mu + t} \right)^{\mu+t} \left(\frac{1 - \mu}{1 - \mu - t} \right)^{1 - \mu - t} \right)^n \\ &= \exp \left(-n \left((\mu + t) \log \frac{\mu + t}{\mu} + (1 - \mu - t) \log \frac{1 - \mu - t}{1 - \mu} \right) \right), \end{aligned}$$

completing the proof. Moreover,

$$\mathbb{P}\left(\mu - \frac{1}{n} \sum_{i=1}^n X_i \geq t\right) = \mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n Z_i - \mu_Z \geq t\right) \leq e^{-n\mathcal{D}(\mu_Z+t, \mu_Z)} = e^{-n\mathcal{D}(1-\mu_X+t, 1-\mu_X)}$$

where $Z_i = 1 - X_i$ (and thus $\mu_Z = 1 - \mu_X$). □

If $0 < \mu \leq 1/2$,

$$\mathcal{D}(1 - \mu + t, 1 - \mu) \geq \frac{t^2}{2\mu(1 - \mu)}.$$

Hence, we get

$$\mathbb{P}\left(\mu - \frac{1}{n} \sum_{i=1}^n X_i \geq t\right) \leq e^{-\frac{nt^2}{2\mu(1-\mu)}} = e^{-u}.$$

Solving for t ,

$$\mathbb{P}\left(\mu - \frac{1}{n} \sum_{i=1}^n X_i \geq \sqrt{\frac{2\mu(1-\mu)u}{n}}\right) \leq e^{-u}.$$

If $X_i \in \{0, 1\}$ are i.i.d. Bernoulli trials, then $\mu = \mathbb{E}X = \mathbb{P}(X = 1)$, $\text{Var}(X) = \mu(1-\mu)$, and $\mathbb{P}\left(\mu - \frac{1}{n} \sum_{i=1}^n X_i \geq t\right) \leq e^{-\frac{nt^2}{2\text{Var}(X)}}$.

The following inequality says that if we pick n reals $a_1, \dots, a_n \in \mathbb{R}$ and add them up each multiplied by a random sign ± 1 , then the expected value of the sum should not be far off from $\sqrt{\sum |a_i|^2}$.

Theorem 7.3. [Khinchine inequality] Let $a_1, \dots, a_n \in \mathbb{R}$, $\epsilon_1, \dots, \epsilon_n$ be i.i.d. Rademacher random variables: $\mathbb{P}(\epsilon_i = 1) = \mathbb{P}(\epsilon_i = -1) = 0.5$, and $0 < p < \infty$. Then

$$A_p \cdot \left(\sum_{i=1}^n |a_i|^2\right)^{1/2} \leq \left(\mathbb{E} \left| \sum_{i=1}^n a_i \epsilon_i \right|^p\right)^{1/p} \leq B_p \cdot \left(\sum_{i=1}^n |a_i|^2\right)^{1/2}$$

for some constants A_p and B_p depending on p .

Proof. Let $\sum |a_i|^2 = 1$ without lossing generality. Then

$$\begin{aligned} \mathbb{E} \left| \sum a_i \epsilon_i \right|^p &= \int_0^\infty \mathbb{P} \left(\left| \sum a_i \epsilon_i \right|^p \geq s^p \right) ds^p \\ &= \int_0^\infty \mathbb{P} \left(\left| \sum a_i \epsilon_i \right| \geq s \right) \cdot ps^{p-1} ds^p \\ &= \int_0^\infty \mathbb{P} \left(\left| \sum a_i \epsilon_i \right| \geq s \right) \cdot ps^{p-1} ds^p \\ &\leq \int_0^\infty 2 \exp\left(-\frac{s^2}{2}\right) \cdot ps^{p-1} ds^p, \text{ Hoeffding's inequality} \\ &= (B_p)^p, \text{ when } p \geq 2. \end{aligned}$$

When $0 < p < 2$,

$$\begin{aligned}\mathbb{E} \left| \sum a_i \epsilon_i \right|^p &\leq \mathbb{E} \left| \sum a_i \epsilon_i \right|^2 \\ &= \mathbb{E} \left| \sum a_i \epsilon_i \right|^{\frac{2}{3}p + (2 - \frac{2}{3}p)} \\ &\leq \left(\mathbb{E} \left| \sum a_i \epsilon_i \right|^p \right)^{\frac{2}{3}} \left(\mathbb{E} \left| \sum a_i \epsilon_i \right|^{6-2p} \right)^{\frac{1}{3}}, \text{ Holder's inequality} \\ &\leq (B_{6-2p})^{2 - \frac{2}{3}p} \cdot \left(\mathbb{E} \left| \sum a_i \epsilon_i \right|^p \right)^{\frac{2}{3}}.\end{aligned}$$

Thus $\mathbb{E} \left| \sum a_i \epsilon_i \right|^p \leq (B_{6-2p})^{6-2p}$, completing the proof. \square