

Solution to “Chapter 2: Basic tail and concentration bounds”

Jiixin Hu

September 28, 2020

1 Summary

Theorem 1.1 (Markov’s inequality). *Let $X \geq 0$ be a random variable with a finite mean. We have*

$$\mathbb{P}(X \geq t) \leq \frac{\mathbb{E}[X]}{t}, \quad \text{for all } t > 0. \quad (1)$$

Theorem 1.2 (Chebyshev’s inequality). *Let $X \geq 0$ be a random variable with a finite mean μ and a finite variance. We have*

$$\mathbb{P}(|X - \mu| \geq t) \leq \frac{\text{var}(X)}{t^2}, \quad \text{for all } t > 0. \quad (2)$$

Theorem 1.3 (Markov’s inequality for polynomial moments). *Let X be a random variable. Suppose that the order k central moment of X exists. Applying Markov’s inequality to the random variable $|X - \mu|^k$ yields*

$$\mathbb{P}(|X - \mu| \geq t) \leq \frac{\mathbb{E}[|X - \mu|^k]}{t^k}, \quad \text{for all } t > 0.$$

Theorem 1.4 (Chernoff bound). *Let X be a random variable. Suppose that the moment generating function of X , denoted $\varphi_X(\lambda)$, exists in the neighborhood of 0; i.e., $\varphi_X(\lambda) = \mathbb{E}[e^{\lambda X}] < +\infty$, for all $\lambda \in (-b, b)$ with some $b > 0$. Applying Markov’s inequality to the random variable $Y = e^{\lambda(X-\mu)}$ yields*

$$\mathbb{P}((X - \mu) \geq t) \leq \frac{\mathbb{E}[e^{\lambda(X-\mu)}]}{e^{\lambda t}}, \quad \text{for all } \lambda \in (0, b).$$

Optimizing the choice of λ for the tightest bound, we obtain the Chernoff bound

$$\mathbb{P}((X - \mu) \geq t) \leq \inf_{\lambda \in [0, b)} \frac{\mathbb{E}[e^{\lambda(X-\mu)}]}{e^{\lambda t}}.$$

Theorem 1.5 (Hoeffding bound for bounded variable). *Let X be a random variable with $\mu = \mathbb{E}(X)$. Suppose that $X \in [a, b]$ almost surely, where $a \leq b \in \mathbb{R}$ are two constants. Then, we have*

$$\mathbb{E}[e^{\lambda X}] \leq e^{\frac{s(b-a)^2}{8}}, \quad \text{for all } \lambda \in \mathbb{R}.$$

Consequently, the variable $X \sim \text{subG}\left(\frac{(b-a)^2}{4}\right)$.

Proof. See Exercise 2.4. □

Theorem 1.6 (Moment of sub-Gaussian variable). *Let $X \sim \text{subG}(\sigma^2)$. For all integer $k \geq 1$, we have*

$$\mathbb{E}[|X|^k] \leq k2^{k/2}\sigma^k\Gamma\left(\frac{k}{2}\right), \quad (3)$$

where the Gamma function is defined as $\Gamma(x) = \int_0^{+\infty} t^{x-1}e^{-t}dt$.

Theorem 1.7 (One-sided Bernstein's inequality). *Let X be a random variable. Suppose $X \leq b$ almost surely. We have*

$$\mathbb{E}\left[e^{\lambda(X-\mathbb{E}[X])}\right] \leq \exp\left\{\frac{\lambda^2\mathbb{E}[X^2]/2}{1-b\lambda/3}\right\}, \quad \text{for all } \lambda \in [0, 3/b).$$

Consequently, let X_i be independent variables, and $X_i \leq b$ almost surely, for all $i \in [n]$. We have

$$\mathbb{P}\left[\sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \geq n\delta\right] \leq \exp\left\{-\frac{n\delta^2}{\sum_{i=1}^n \mathbb{E}[X_i^2]/n + b\delta/3}\right\}, \quad \text{for all } \delta \geq 0. \quad (4)$$

Particularly, let X_i be independent nonnegative variables, for all $i \in [n]$. The equation (4) becomes

$$\mathbb{P}\left[\sum_{i=1}^n (Y_i - \mathbb{E}[Y_i]) \leq n\delta\right] \leq \exp\left\{-\frac{n\delta^2}{\sum_{i=1}^n \mathbb{E}[Y_i^2]/n}\right\}, \quad \text{for all } \delta \geq 0. \quad (5)$$

Definition 1 (Bernstein's condition). Let X be a random variable with mean $\mu = \mathbb{E}[X]$ and variance $\sigma^2 = \text{var}(X)$. We say X satisfies the Bernstein's condition with parameter b if

$$\left|\mathbb{E}[(X - \mu)^k]\right| \leq \frac{1}{2}k!\sigma^2b^{k-2}, \quad \text{for } k = 3, 4, \dots \quad (6)$$

Note that bounded random variables satisfy the Bernstein's condition.

Theorem 1.8 (Bernstein-type bound). *For any variable X satisfying the Bernstein's condition, we have*

$$\mathbb{E}\left[e^{\lambda(X-\mu)}\right] \leq \exp\left\{\frac{\lambda^2\sigma^2}{2(1-b|\lambda|)}\right\}, \quad \text{for all } |\lambda| \leq \frac{1}{b},$$

and the concentration inequality

$$\mathbb{P}[|X - \mu| \geq t] \leq 2\exp\left\{-\frac{t^2}{2(\sigma^2 + bt)}\right\}, \quad \text{for all } t \geq 0. \quad (7)$$

Definition 2 (Bounded difference property). Let $f: \mathbb{R}^n \mapsto \mathbb{R}$ be a function. The function f satisfies the bounded difference property with parameter (L_1, \dots, L_n) if we have

$$\left|f(x^{(k)}) - f(x'^{(k)})\right| \leq L_k, \quad (8)$$

for all $k \in [n]$ and for all $x^{(k)} = (x_1, \dots, x_k, \dots, x_n)$, $x'^{(k)} = (x_1, \dots, x'_k, \dots, x_n) \in \mathbb{R}^n$.

Theorem 1.9 (Bounded differences inequality). *Let $f: \mathbb{R}^n \mapsto \mathbb{R}$ be a function satisfies the bounded difference property (8), and the random variable $X = (X_1, \dots, X_n)$ has independent components. Then,*

$$\mathbb{P}[|f(X) - \mathbb{E}[f(X)]| \geq t] \leq 2e^{-\frac{2t^2}{\sum_{i=1}^n L_i^2}}, \quad \text{for all } t \geq 0.$$

2 Exercises

2.1 Exercise 2.1

(Tightness of inequalities.) The Markov's and Chebyshev's inequalities are not able to be improved in general.

- (a) Provide a random variable $X \geq 0$ that attains the equality in Markov's inequality (1).
- (b) Provide a random variable Y that attains the equality in Chebyshev's inequality (2).

Solution:

- (a) For a given constant $t > 0$, we define a variable $Y_t = X - t\mathbb{1}[X \geq t]$, where $\mathbb{1}$ is the indicator function. Note that Y_t is a nonnegative variable. The Markov's inequality follows by taking the expectation to Y_t ,

$$\mathbb{E}[Y_t] = \mathbb{E}[X] - t\mathbb{P}[X \geq t] \geq 0.$$

Therefore, Markov's inequality meets the equality if and only if the expectation $\mathbb{E}[Y_t] = 0$. Since Y_t is nonnegative, we have $\mathbb{P}(Y_t = 0) = 1$. Note that $Y_t = 0$ if and only if $X = 0$ or $X = t$.

Hence, for the given constant $t > 0$, the nonnegative variable X with distribution $\mathbb{P}(X \in \{0, t\}) = 1$ attains the equality of Markov's inequality.

- (b) Chebyshev's inequality follows by applying Markov's inequality to the nonnegative random variable $Z = (X - \mathbb{E}[X])^2$. Simialrly as in part (a), given a constant $t > 0$, the variable $Z = (X - \mathbb{E}[X])^2$ with distribution $\mathbb{P}(Z \in \{0, t^2\}) = 1$ attains the equality of the Markov's inequality for Z . Consequently, the variable X attains the equality of the Chebyshev's inequality for X . By transformation, the distribution of X satisfies the followings formula,

$$\mathbb{P}(X = x) = \begin{cases} p & \text{if } x = c, \\ \frac{1-p}{2} & \text{if } x = c - t \text{ or } x = c + t, \\ 0 & \text{otherwise,} \end{cases}$$

where $c \in \mathbb{R}$ is a constant and $p \in [0, 1]$.

Remark 1 (Tightness of Markov's inequality). Only a few variables attain the equalities in Markov's and Chebyshev's inequalities. In research, we should pay attention to the concentration bounds tighter than Markov's inequality.

2.2 Exercise 2.2

Lemma 1 (Standard normal distribution). *Let $\phi(z)$ be the density function of a standard normal variable $Z \sim N(0, 1)$. Then,*

$$\phi'(z) + z\phi(z) = 0, \tag{9}$$

and

$$\phi(z) \left(\frac{1}{z} - \frac{1}{z^3} \right) \leq \mathbb{P}(Z \geq z) \leq \phi(z) \left(\frac{1}{z} - \frac{1}{z^3} + \frac{3}{z^5} \right), \quad \text{for all } z > 0. \tag{10}$$

Proof. First, we prove the equation (9).

The pdf of the standard normal distribution is

$$\phi(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right).$$

The equation (9) follows by taking the derivative of $\phi(z)$. Specifically,

$$\phi'(z) = -z \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) = -z\phi(z).$$

Next, we prove the equation (10).

We write the upper tail probability of the standard normal variable as

$$\mathbb{P}(Z \geq z) = \int_z^{+\infty} \phi(t) dt = \int_z^{+\infty} -\frac{1}{t} \phi'(t) dt = \frac{1}{z} \phi(z) - \int_z^{+\infty} \frac{1}{t^2} \phi(t) dt, \quad (11)$$

where the second equality follows by the equation (9). Applying the equation (9) to the last term in equation (11) yields

$$\int_z^{+\infty} \frac{1}{t^2} \phi(t) dt = \int_z^{+\infty} \frac{1}{t^3} \phi'(t) dt = -\frac{1}{z^3} \phi(z) + \int_z^{+\infty} \frac{3}{t^4} \phi(t) dt \geq -\frac{1}{z^3} \phi(z) \quad (12)$$

Plugging the equation (12) into the equation (11), we obtain $\mathbb{P}(Z \geq z) \geq \phi(z) \left(\frac{1}{z} - \frac{1}{z^3}\right)$. Applying the equation (9) again to the equation (12) yields

$$\int_z^{+\infty} \frac{3}{t^4} \phi(t) dt = \int_z^{+\infty} -\frac{3}{t^5} \phi'(t) dt = \frac{3}{z^5} \phi(z) - \int_z^{+\infty} \frac{15}{t^6} \phi(t) dt \leq \frac{3}{z^5} \phi(z). \quad (13)$$

Combing equations (11), (12) and (13), we obtain $\mathbb{P}(Z \geq z) \leq \phi(z) \left(\frac{1}{z} - \frac{1}{z^3} + \frac{3}{z^5}\right)$. \square

Remark 2. Direct calculation of tail probability for a univariate normal variable is hard. Equation (10) provides a numerical approximation to the tail probability. Particularly, the tail probability decays at the rate of $z^{-1}e^{-z^2/2}$ as $z \rightarrow +\infty$. The decay rate is faster than polynomial rate $\mathcal{O}(z^{-\alpha})$, for any $\alpha \geq 1$.

2.3 Exercise 2.3

Lemma 2 (Polynomial bound and Chernoff bound). *Let $X \geq 0$ be a nonnegative variable. Suppose that the moment generating function of X , denoted $\varphi_X(\lambda)$, exists in the neighborhood of $\lambda = 0$. Given some $\delta > 0$, we have*

$$\inf_{k \in \mathbb{Z}_+} \frac{\mathbb{E}[|X|^k]}{\delta^k} \leq \inf_{\lambda > 0} \frac{\mathbb{E}[e^{\lambda X}]}{e^{\lambda \delta}}. \quad (14)$$

Consequently, an optimized bound based on polynomial moments is always at least as good as the Chernoff upper bound.

Proof. By power series, we have

$$e^{\lambda X} = \sum_{k=0}^{+\infty} \frac{X^k \lambda^k}{k!}, \quad \text{for all } \lambda \in \mathbb{R}. \quad (15)$$

Since the moment generating function $\varphi_X(\lambda)$ exists in the neighborhood of $\lambda = 0$, there exists a constant $b > 0$ such that

$$\mathbb{E}[e^{\lambda X}] = \sum_{k=0}^{+\infty} \frac{\mathbb{E}[|X|^k] \lambda^k}{k!} < +\infty, \quad \text{for all } \lambda \in (0, b).$$

Hence, the moment $\mathbb{E}[|X|^k]$ exists, for all $k \in \mathbb{Z}_+$. Applying power series (15) to the right hand side of equation (14) yields

$$\inf_{\lambda > 0} \frac{\mathbb{E}[e^{\lambda X}]}{e^{\lambda \delta}} = \frac{\sum_{k=0}^{+\infty} \frac{\mathbb{E}[|X|^k] \lambda^k}{k!}}{\sum_{k=0}^{+\infty} \frac{\lambda^k \delta^k}{k!}}. \quad (16)$$

By Cauchy's third inequality, we have

$$\frac{\sum_{k=0}^{+\infty} \frac{\mathbb{E}[|X|^k] \lambda^k}{k!}}{\sum_{k=0}^{+\infty} \frac{\lambda^k \delta^k}{k!}} \geq \inf_{k \in \mathbb{Z}_+} \frac{\mathbb{E}[|X|^k]}{\delta^k} \quad (17)$$

Therefore, we obtain the equation (14) by combining the equation (16) with equation (17). \square

Remark 3. Applying different functions $g(X)$ to the Markov's inequality leads to different bounds for the tail probability of variable X . Equation (14) implies that the optimized polynomial bound is at least as tight as the Chernoff bound, provided that the moment generating function of X exists in the neighborhood of 0.

2.4 Exercise 2.4

In Exercise 2.4, we prove Theorem 1.5, the Hoeffding bound for a bounded variable.

Proof. Let X be a bounded random variable, and $X \in [a, b]$ almost surely, where $a \leq b \in \mathbb{R}$ are two constants. Let $\mu = \mathbb{E}[X]$. Define the function

$$g(\lambda) = \log \mathbb{E}[e^{\lambda X}], \quad \text{for all } \lambda \in \mathbb{R}.$$

Applying Taylor Expansion to $g(\lambda)$ at 0, we have

$$g(\lambda) = g(0) + g'(0)\lambda + \frac{g''(\lambda_0)}{2}\lambda^2, \quad \text{where } \lambda_0 = t\lambda, \text{ for some } t \in [0, 1]. \quad (18)$$

In equation (18), the term $g(0) = \log \mathbb{E}[e^0] = 0$. By power series (15), we obtain the first derivative $g'(\lambda)$ as follows,

$$\begin{aligned} g'(\lambda) &= \left(\log \sum_{k=0}^n \frac{\lambda^k}{k!} \mathbb{E}[X^k] \right)' \\ &= \sum_{k=0}^n \frac{\lambda^k}{k!} \mathbb{E}[X^{(k+1)}] \bigg/ \sum_{k=0}^n \frac{\lambda^k}{k!} \mathbb{E}[X^k] \\ &= \frac{\mathbb{E}[X e^{\lambda X}]}{\mathbb{E}[e^{\lambda X}]}. \end{aligned} \quad (19)$$

Therefore, $g'(0) = \mathbb{E}[X] = \mu$. Taking the derivative to equation (19), we obtain the second-order derivative $g''(\lambda)$ as follows,

$$\begin{aligned} g''(\lambda) &= \sum_{k=0}^n \frac{\lambda^k}{k!} \mathbb{E}[X^{(k+2)}] \bigg/ \sum_{k=0}^n \frac{\lambda^k}{k!} \mathbb{E}[X^k] - \left(\sum_{k=0}^n \frac{\lambda^k}{k!} \mathbb{E}[X^{(k+1)}] \bigg/ \sum_{k=0}^n \frac{\lambda^k}{k!} \mathbb{E}[X^k] \right)^2 \\ &= \frac{\mathbb{E}[X^2 e^{\lambda X}]}{\mathbb{E}[e^{\lambda X}]} - \left(\frac{\mathbb{E}[X e^{\lambda X}]}{\mathbb{E}[e^{\lambda X}]} \right)^2. \end{aligned}$$

We interpret the second-order derivative $g''(\lambda)$ as the variance of X with the re-weighted distribution $dP' = e^{\lambda X} / \mathbb{E}[e^{\lambda X}] dP_X$, where P_X is the distribution of X . Taking the integral of 1 with respect to dP' , we have

$$\int_{-\infty}^{+\infty} dP' = \int_{-\infty}^{+\infty} \frac{e^{\lambda X}}{\mathbb{E}[e^{\lambda X}]} dP_X = 1,$$

which implies that the function P' is a valid probability distribution. Under all possible re-weighted distributions, the variance of X is upper bounded as follows,

$$\text{var}(X) = \text{var}\left(X - \frac{a+b}{2}\right) \leq \mathbb{E}\left[\left(X - \frac{a+b}{2}\right)^2\right] \leq \frac{(b-a)^2}{4},$$

where the term $\frac{(b-a)^2}{4}$ follows by letting X supported on the boundaries a and b only. Hence, the second-order derivative $g''(\lambda) \leq \frac{(b-a)^2}{4}$. We plug the results of g' and g'' into the equation (18). Then,

$$g(\lambda) = g(0) + g'(0)\lambda + \frac{g''(\lambda_0)}{2}\lambda^2 \leq 0 + \lambda\mu + \frac{(b-a)^2}{8}\lambda^2. \quad (20)$$

Taking the exponentiation on both sides of the inequality (20), we have

$$\mathbb{E}[e^{\lambda X}] = \exp(g(\lambda)) \leq e^{\mu\lambda + \frac{(b-a)^2}{8}\lambda^2}. \quad (21)$$

The equation (21) implies that X is a sub-Gaussian variable with at most $\sigma = \frac{(b-a)}{2}$. \square

Remark 4. For any bounded random variable X supported on $[a, b]$, X is a sub-gaussian variable with parameter at most $\sigma^2 = (b-a)^2/4$. All the properties for sub-Gaussian variables apply to the bounded variables.

2.5 Exercise 2.5

Lemma 3 (Sub-Gaussian bounds and means/variance). *Let X be a random variable such that*

$$\mathbb{E}[e^{\lambda X}] \leq e^{\frac{\lambda^2 \sigma^2}{2} + \mu\lambda}, \quad \text{for all } \lambda \in \mathbb{R}. \quad (22)$$

Then, $\mathbb{E}[X] = \mu$ and $\text{var}(X) \leq \sigma^2$.

Proof. By equation (22), the moment generating function of X , denoted $\varphi_X(\lambda)$, exists in the neighborhood of $\lambda = 0$. Hence, the mean and variance of X exist. For all λ in the neighborhood of $\lambda = 0$, applying power series on both sides of equation (22) yields

$$\lambda \mathbb{E}[X] + \frac{\lambda^2}{2} \mathbb{E}[X^2] + o(\lambda^2) \leq \mu\lambda + \frac{\lambda^2 \sigma^2 + \lambda^2 \mu^2}{2} + o(\lambda^2). \quad (23)$$

Dividing by $\lambda > 0$ on both sides of equation (23) and letting $\lambda \rightarrow 0^+$, we have $\mathbb{E}(X) \leq \mu$. Dividing by $\lambda < 0$ on both sides of equation (23) and letting $\lambda \rightarrow 0^-$, we have $\mathbb{E}(X) \geq \mu$. Therefore, we obtain the mean $\mathbb{E}[X] = \mu$. Then, we divide $2/\lambda^2$ on both sides of equation (23), for $\lambda \neq 0$. The term $\mathbb{E}[X]\lambda$ and $\mu\lambda$ are cancelled. We have $\mathbb{E}[X^2] \leq \sigma^2 + \mu^2$, and thus the $\text{var}(X) \leq \mathbb{E}[X^2] - (\mathbb{E}[X])^2 = \sigma^2$. \square

Question: Let σ_{\min}^2 denote the smallest possible σ satisfying the inequality (22). Is it true that $\text{var}(X) = \sigma_{\min}^2$?

Solution: The statement that $\text{var}(X) = \sigma_{\min}^2$ is not necessarily true. Recall the function $g(\lambda)$ in Exercise 2.4. By the results in Exercise 2.4, the equation (22) is equal to

$$g''(\lambda) \leq \sigma^2, \quad \text{for all } \lambda \in \mathbb{R},$$

where $g''(\lambda)$ is the variance of X with the re-weighted distribution defined in Exercise 2.4. Therefore, we have $\max_{\lambda} g''(\lambda) = \sigma_{\min}^2$. Note that $g''(0) = \text{var}(X)$. To let the equality $\text{var}(X) = \sigma_{\min}^2$ hold, we need to show that $\max_{\lambda} g''(\lambda) = g''(0)$ holds for X .

However, the statement $\max_{\lambda} g''(\lambda) = g''(0)$ is not necessarily true. A counter example is below. Consider a random variable $Y \sim \text{Ber}(1/3)$. The variance of Y is $\text{var}(Y) = 2/9$. Let $\lambda = 1$. The re-weighted distribution dP' is

$$P'(Y = 0) = \frac{2}{3\mathbb{E}[e^Y]} \quad \text{and} \quad P'(Y = 1) = \frac{e}{3\mathbb{E}[e^Y]}, \quad \text{where } \mathbb{E}[e^Y] = \frac{2}{3} + \frac{e}{3}.$$

The variance of Y with dP' is $2/3\mathbb{E}[e^Y] \times e/3\mathbb{E}[e^Y] = 0.2442 > 2/9$. Therefore, we have $\text{var}(Y) < g''(1) \leq \max_{\lambda} g''(\lambda) = \sigma_{\min}^2$. The statement $\max_{\lambda} g''(\lambda) = g''(0)$ is not true for this variable Y .

Remark 5. Parameters of a sub-Gaussian distribution provide the exact value of the mean and an upper bound of the variance; i.e., $\mathbb{E}[X] = \mu$ and $\text{var}(X) \leq \sigma^2$. Suppose the moment generating function of variable X exists over the entire real interval. Then, the tail distribution of X is bounded by a sub-Gaussian distribution with a proper choice of σ^2 .

2.6 Exercise 2.6

Lemma 4 (Lower bounds on squared sub-Gaussians). *Let $\{X_i\}_{i=1}^n$ be an i.i.d. sequence of zero-mean sub-Gaussian variables with parameter σ . The normalized sum $Z_n = \frac{1}{n} \sum_{i=1}^n X_i^2$ satisfies*

$$\mathbb{P}[Z_n - \mathbb{E}[Z_n] \leq \sigma^2 \delta] \leq e^{-n\delta^2/16}, \quad \text{for all } \delta \geq 0. \quad (24)$$

The equation (24) implies that the lower tail of the sum of squared sub-Gaussian variables behaves in a sub-Gaussian way.

Proof. Since X_i^2 are i.i.d. nonnegative variables, we apply the equation (5) to the variables $\{X_i^2\}_{i=1}^n$. Then, we have

$$\mathbb{P}\left[\sum_{i=1}^n (X_i^2 - \mathbb{E}[X_i^2]) \leq n\sigma^2 \delta\right] \leq \exp\left\{-\frac{n\delta^2 \sigma^4}{\mathbb{E}[X_1^4]}\right\}, \quad \text{for all } \delta \geq 0. \quad (25)$$

By equation (3), we have

$$\mathbb{E}[X_1^4] \leq 16\sigma^4. \quad (26)$$

Combing equations (25), (26) and the definition of Z_n , we obtain

$$\mathbb{P}[Z_n - \mathbb{E}[Z_n] \leq \sigma^2 \delta] \leq \exp \left\{ -\frac{n\delta^2}{16} \right\}, \quad \text{for all } \delta \geq 0.$$

□

Remark 6. Equation (24) implies that the lower tail of the sum of squared sub-Gaussian variables behaves in a sub-Gaussian way. In following sections, we will show that the variable $Z_n - \mathbb{E}[Z_n]$ in Lemma 4 is a sub-exponential variable.

2.7 Exercise 2.7

Lemma 5 (Bennett's inequality). *Let X_1, \dots, X_n be a sequence of independent zero-mean random variables with $|X_i| \leq b$ and $\text{var}(X_i) = \sigma_i^2$, for all $i \in [n]$. Then, we have the Bennett's inequality*

$$\mathbb{P} \left[\sum_{i=1}^n X_i \geq n\delta \right] \leq \exp \left\{ -\frac{n\sigma^2}{b^2} h \left(\frac{b\delta}{\sigma^2} \right) \right\}, \quad \text{for all } \delta \geq 0,$$

where $\sigma^2 = \frac{1}{n} \sum_{i=1}^n \sigma_i^2$ and $h(t) := (1+t) \log(1+t) - t$ for $t \geq 0$.

Proof. First, we consider the moment generating function of X_i , for all $i \in [n]$.

By power series, for all $i \in [n]$, we have

$$\mathbb{E} \left[e^{\lambda X_i} \right] = \sum_{k=0}^{+\infty} \frac{\lambda^k \mathbb{E}[X_i^k]}{k!} = 1 + 0 + \sum_{k=2}^{+\infty} \frac{\lambda^k \mathbb{E}[X_i^k]}{k!} \leq \exp \left\{ \sum_{k=2}^{+\infty} \frac{\lambda^k \mathbb{E}[X_i^k]}{k!} \right\}, \quad (27)$$

where the 0 comes from the fact that $\mathbb{E}[X_i] = 0$, and the last inequality follows from $1 + x \leq e^x$. By $|X_i| < b$, we bound the last term in equation (27) as follows

$$\sum_{k=2}^{+\infty} \frac{\lambda^k \mathbb{E}[X_i^k]}{k!} \leq \sum_{k=2}^{+\infty} \frac{\lambda^k \mathbb{E}[X_i^2 |X_i|^{k-2}]}{k!} \leq \sum_{k=2}^{+\infty} \frac{\lambda^k \sigma_i^2 b^{k-2}}{k!} = \sigma_i^2 \left(\frac{e^{\lambda b} - 1 - \lambda b}{b^2} \right). \quad (28)$$

Combing the equation (27) with equation (28), we obtain the following upper bound of the moment generating function of $\sum_{i=1}^n X_i$.

$$\mathbb{E} \left[e^{\lambda \sum_{i=1}^n X_i} \right] = \prod_{i=1}^n \mathbb{E} \left[e^{\lambda X_i} \right] \leq \exp \left\{ n\sigma^2 \left(\frac{e^{\lambda b} - 1 - \lambda b}{b^2} \right) \right\}, \quad (29)$$

where $\sigma^2 = \frac{1}{n} \sum_{i=1}^n \sigma_i^2$. Combing the Chernoff bound with equation (29), the upper tail of $\sum_{i=1}^n X_i$ follows

$$\begin{aligned} \mathbb{P} \left[\sum_{i=1}^n X_i \geq n\delta \right] &\leq \exp \left\{ n\sigma^2 \left(\frac{e^{\lambda b} - 1 - \lambda b}{b^2} \right) - \lambda n\delta \right\} \\ &= \exp \left\{ \frac{n\sigma^2}{b^2} \left(e^{\lambda b} - \lambda b - \lambda \frac{\delta b^2}{\sigma^2} - 1 \right) \right\}, \quad \text{for all } \delta \geq 0. \end{aligned} \quad (30)$$

The upper bound (30) achieves the minimum when $\lambda = b^{-1} \log(1 + \frac{\delta b}{\sigma^2})$ by the first-order condition of minimization. Plugging $\lambda = b^{-1} \log(1 + \frac{\delta b}{\sigma^2})$ into the equation (30), we obtain the Bennett's inequality

$$\mathbb{P}\left[\sum_{i=1}^n X_i \geq n\delta\right] \leq \exp\left\{-\frac{n\sigma^2}{b^2} h\left(\frac{b\delta}{\sigma^2}\right)\right\}, \quad \text{for all } \delta \geq 0, \quad (31)$$

where $h(t) := (1+t) \log(1+t) - t$ for $t \geq 0$.

Further, we show that the Bennett's inequality is at least as good as the Bernstein's inequality.

The Bernstein's inequality for $\sum_{i=1}^n X_i$ is

$$\mathbb{P}\left[\sum_{i=1}^n X_i \geq n\delta\right] \leq \exp\left\{\frac{-3n\delta^2}{(2b\delta + 6\sigma^2)}\right\} = \exp\left\{-\frac{n\sigma^2}{b^2} g\left(\frac{b\delta}{\sigma^2}\right)\right\}, \quad \text{for all } \delta \geq 0, \quad (32)$$

where $g(t) := \frac{3t^2}{2t+6}$ for $t \geq 0$. Since $g(t) \leq h(t)$ holds for all $t \geq 0$, we conclude that the Bennett's inequality (31) is at least as good as Bernstein's inequality (32). \square

Remark 7. So far, we have three inequalities controlling the tail of bounded variables: Hoeffding's inequality, Bernstein's inequality, and Bennett's inequality. Particularly, Hoeffding's inequality implies the sub-Gaussianity of bounded variables. As the proof for Lemma 5 shows, Bennett's inequality is at least as good as the Bernstein's inequality, for bounded random variables.

2.8 Exercise 2.8

Lemma 6 (Bernstein and expectation). *Let Z be a nonnegative random variable satisfying the following concentration inequality*

$$\mathbb{P}[Z \geq t] \leq Ce^{-\frac{t^2}{2(\nu^2+Bt)}}, \quad \text{for all } t \geq 0, \quad (33)$$

where (ν, B) are two positive constants and $C \geq 1$. Then, the expectation of Z satisfies

$$\mathbb{E}[Z] \leq 2\nu(\sqrt{\pi} + \sqrt{\log C}) + 4B(1 + \log C). \quad (34)$$

Further, let $\{X_i\}_{i=1}^n$ be a sequence of i.i.d. zero-mean variables satisfying the Bernstein condition (6). The sample mean of $\{X_i\}_{i=1}^n$ satisfies

$$\mathbb{E}\left[\left|\frac{1}{n} \sum_{i=1}^n X_i\right|\right] \leq 2\sigma(\sqrt{\pi} + \sqrt{\log 2}) + 4b(1 + \log 2). \quad (35)$$

Proof. First, we prove the equation (34).

By equation (33), we have

$$\begin{aligned} \mathbb{P}[Z \geq t] &\leq C \max\left\{\exp\left(-\frac{t^2}{4\nu^2}\right), \exp\left(-\frac{t^2}{4Bt}\right)\right\} \\ &\leq C \exp\left(-\frac{t^2}{4\nu^2}\right) + C \exp\left(-\frac{t^2}{4Bt}\right). \end{aligned} \quad (36)$$

Plugging the inequality (36) to $\mathbb{E}[Z] = \int_0^{+\infty} \mathbb{P}[Z \geq t] dt$, we have

$$\mathbb{E}[Z] = \int_0^{+\infty} \min\left\{1, C \exp\left(-\frac{t^2}{4\nu^2}\right)\right\} dt + \int_0^{+\infty} \min\left\{1, C \exp\left(-\frac{t^2}{4Bt}\right)\right\} dt =: I_1 + I_2.$$

To evaluate I_1 , we solve $1 = C \exp\left(-\frac{t^2}{4\nu^2}\right)$, and the minimization term becomes

$$\min \left\{ 1, C \exp\left(-\frac{t^2}{4\nu^2}\right) \right\} = \begin{cases} 1 & \text{when } t < 2\nu\sqrt{\log C}, \\ C \exp\left(-\frac{t^2}{4\nu^2}\right) & \text{when } t \geq 2\nu\sqrt{\log C}. \end{cases}$$

Therefore, let $y = \frac{t}{2\nu} - \sqrt{\log C}$, we have

$$\begin{aligned} I_1 &= \int_0^{2\nu\sqrt{\log C}} 1 dt + \int_{2\nu\sqrt{\log C}}^{+\infty} C \exp\left(-\frac{t^2}{4\nu^2}\right) dt \\ &= 2\nu\sqrt{\log C} + 2\nu \int_0^{+\infty} \exp\left(-y^2 - 2y\sqrt{\log C}\right) dy \\ &\leq 2\nu\sqrt{\log C} + 2\nu \int_0^{+\infty} \exp\left(-y^2\right) dy. \end{aligned}$$

By Gaussian integral $\int_0^{+\infty} e^{-x^2} dx = \frac{\sqrt{\pi}}{2}$, we obtain $I_1 \leq 2\nu(\sqrt{\pi} + \sqrt{\log C})$. Similarly, we evaluate I_2 as follows

$$\begin{aligned} I_2 &= \int_0^{4B\log C} 1 dt + \int_{4B\log C}^{+\infty} C \exp\left(-\frac{t}{4B}\right) dt \\ &= 4B(\log C + 1). \end{aligned}$$

Hence, we obtain the expectation of Z ,

$$\mathbb{E}[Z] = I_1 + I_2 \leq 2\nu(\sqrt{\pi} + \sqrt{\log C}) + 4B(\log C + 1).$$

Next, we prove the equation (35).

For all $i \in [n]$, since X_i satisfies the Bernstein condition with parameter (σ, b) , the variable X_i satisfies the concentration bound (7),

$$\mathbb{P}[|X_i| \geq t] \leq 2 \exp\left\{-\frac{t^2}{2(\sigma^2 + bt)}\right\}, \quad \text{for all } t \geq 0.$$

By equation (34), we have

$$\mathbb{E}[|X_i|] \leq 2\sigma(\sqrt{\pi} + \sqrt{\log 2}) + 4b(1 + \log 2).$$

Therefore, the expectation of the sample mean satisfies

$$\mathbb{E}\left[\left|\frac{1}{n} \sum_{i=1}^n X_i\right|\right] \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}[|X_i|] \leq 2\sigma(\sqrt{\pi} + \sqrt{\log 2}) + 4b(1 + \log 2).$$

□

Remark 8. For a nonnegative random variables satisfying the Bernstein-type inequality (33), the expectation of the variable is upper bounded by a function of the parameter (ν, B, C) . Particularly, for a zero-mean random variable X satisfying the Bernstein condition with parameter b , $|X|$ satisfies the inequality (33) with parameters $(\sigma, b, 2)$, where $\sigma^2 = \text{var}(X)$. Then, the expectation of the absolute variable $|X|$ is upper bounded by a function of (σ^2, b) .

2.9 Exercise 2.9

Lemma 7 (Sharp upper bounds on binomial tails). *Let $\{X_i\}_{i=1}^n$ be a sequence of i.i.d. Bernoulli variables with parameter $\alpha \in (0, 1/2]$. Consider the binomial random variable $Z_n = \sum_{i=1}^n X_i$. The left tail probability of Z_n is upper bounded as*

$$\mathbb{P}[Z_n \leq \delta n] \leq e^{-nD(\delta\|\alpha)}, \quad \text{for all } \delta \in (0, \alpha), \quad (37)$$

where the quantity

$$D(\delta\|\alpha) := \delta \log \frac{\delta}{\alpha} + (1 - \delta) \log \frac{1 - \delta}{1 - \alpha} \quad (38)$$

is the KL divergence between the Bernoulli distributions with parameters δ and α , respectively. Further, the upper bound (41) is strictly tighter than the Hoeffding bound for all $\delta \in (0, \alpha)$.

Proof. Applying the Chernoff inequality to the random variable $-Z_n$, we have

$$\mathbb{P}[-Z_n \geq -\delta n] = \mathbb{P}[Z_n \leq \delta n] \leq \frac{\mathbb{E}[e^{-\lambda Z_n}]}{e^{-\lambda \delta n}}, \quad \text{for all } \lambda > 0.$$

Therefore, the left tail probability of random variable Z_n satisfies

$$\mathbb{P}[Z_n \leq \delta n] \leq \frac{\mathbb{E}[e^{-\lambda Z_n}]}{e^{-\lambda \delta n}} = \left\{ \exp \left\{ \log \left(1 - \alpha + \alpha e^{-\lambda} \right) + \lambda \delta \right\} \right\}^n, \quad \text{for all } \lambda > 0, \quad (39)$$

where the second equality follows from the moment generating function of Z_n ,

$$\mathbb{E}[e^{-\lambda Z_n}] = \left(1 - \alpha + \alpha e^{-\lambda} \right)^n.$$

The upper bound (39) achieves the minimum when $\lambda = -\log \frac{\delta(1-\alpha)}{\alpha(1-\delta)}$ by the first-order condition. Plugging the minimizer into equation (39), we obtain the result

$$\mathbb{P}[Z_n \leq \delta n] \leq e^{-nD(\delta\|\alpha)}, \quad \text{for all } \delta \in (0, \alpha).$$

Next, we show that equation (41) is strictly tighter than the Hoeffding bound of Z_n , for all $\delta \in (0, \alpha)$.

Hoeffding bound of random variable Z_n is

$$\mathbb{P}[Z_n \leq \delta n] \leq e^{-2n(\delta - \alpha)^2}, \quad \text{for all } \delta \in (0, \alpha). \quad (40)$$

Note that parameter α is fixed. Consider the function $g(\delta) = 2(\delta - \alpha)^2 - D(\delta\|\alpha)$ for $\delta \in (0, \alpha)$. The first-order derivative and second-order derivative of g are

$$g'(\delta) = 4(\delta - \alpha) - \log \frac{\delta}{\alpha} + \log \frac{1 - \delta}{1 - \alpha} \quad \text{and} \quad g''(\delta) = 4 - \frac{1}{\delta} - \frac{1}{1 - \delta}.$$

Note that $g(\alpha) = g'(\alpha) = 0$, and $g''(\delta) < 0$ for all $0 < \delta < 1/2$. Then, we have $g(\delta) < 0$, for all $\delta \in (0, \alpha)$. Hence, the upper bound (41) is strictly tighter than Hoeffding bound (40). \square

Corollary 1 (Sharp upper bounds on binomial right tails). *Let $\{X_i\}_{i=1}^n$ be a sequence of i.i.d. Bernoulli variables with parameter $\alpha \in [1/2, 1)$. Consider the binomial random variable $Z_n = \sum_{i=1}^n X_i$. The left tail probability of Z_n is upper bounded as*

$$\mathbb{P}[Z_n \geq \delta n] \leq e^{-nD(\delta\|\alpha)}, \quad \text{for all } \delta \in (\alpha, 1), \quad (41)$$

where $D(\delta\|\alpha)$ is the KL divergence defined in Lemma 7.

Proof. Applying the Chernoff inequality directly to the random variable Z_n , we have

$$\mathbb{P}[Z_n \geq \delta n] \leq \frac{\mathbb{E}[e^{\lambda Z_n}]}{e^{\lambda \delta n}} = \left\{ \exp \left\{ \log \left(1 - \alpha + \alpha e^\lambda \right) - \lambda \delta \right\} \right\}^n, \quad \text{for all } \lambda > 0.$$

Plugging the minimizer $\lambda = \log \frac{\delta(1-\alpha)}{\alpha(1-\delta)}$ into the above inequality, we obtain the results

$$\mathbb{P}[Z_n \geq \delta n] \leq e^{-nD(\delta||\alpha)}, \quad \text{for all } \delta \in (\alpha, 1).$$

□

Remark 9. The upper bound on a binomial tail is a function of the KL divergence between Bernoulli distributions with distinct parameters. The bound using KL divergence is strictly tighter than the Hoeffding bound. The underperformance of Hoeffding bound may attribute to the utilization of the sub-Gaussian parameter σ^2 . The sub-Gaussian parameter σ^2 is not necessarily the optimal choice to describe the tail performance of a random variable with good properties.

2.10 Exercise 2.10

Lemma 8 (Lower bounds on binomial tails). *Let $\{X_i\}_{i=1}^n$ be a sequence of i.i.d. Bernoulli variables with parameter $\alpha \in (0, 1/2]$. Consider the binomial random variable $Z_n := \sum_{i=1}^n X_i$. For some fixed $\delta \in (0, \alpha)$, let $m = \lfloor n\delta \rfloor$; i.e., m is the largest integer less or equal to $n\delta$. Let $\tilde{\delta} = m/n$. We have*

$$\frac{1}{n} \log \mathbb{P}[Z_n \leq n\delta] \geq \frac{1}{n} \log \binom{n}{m} + \tilde{\delta} \log \alpha + (1 - \tilde{\delta}) \log(1 - \alpha). \quad (42)$$

Further, the binomial coefficient satisfies

$$\frac{1}{n} \log \binom{n}{m} \geq \phi(\tilde{\delta}) - \frac{\log(n+1)}{n}. \quad (43)$$

Consequently, the lower tail of binomial variable Z_n satisfies

$$\mathbb{P}[Z_n \leq n\delta] \geq \frac{1}{n+1} e^{-nD(\delta||\alpha)}, \quad (44)$$

where $D(\delta||\alpha)$ is the KL divergence defined in equation (38).

Proof. First, we prove equation (42).

Since Z_n is a binomial variable with size n and probability α , we have

$$\begin{aligned} \mathbb{P}[Z_n \leq n\delta] &= \sum_{k=1}^m \binom{n}{k} \alpha^k (1-\alpha)^{n-k} \\ &\geq \binom{n}{m} \alpha^m (1-\alpha)^{n-m}. \end{aligned} \quad (45)$$

Taking the log and dividing by n on the both sides of the inequality (45), we obtain

$$\frac{1}{n} \log \mathbb{P}[Z_n \leq n\delta] \geq \frac{1}{n} \log \binom{n}{m} + \tilde{\delta} \log \alpha + (1 - \tilde{\delta}) \log(1 - \alpha).$$

Next, we prove equation (43).

Consider a binomial variable $Y \sim \text{Bin}(n, \tilde{\delta})$. For all $k = 0, 1, \dots, (n-1)$, the ratio between $\mathbb{P}[Y = k]$ and $\mathbb{P}[Y = k+1]$ is

$$\frac{\mathbb{P}[Y = k+1]}{\mathbb{P}[Y = k]} = \frac{\binom{n}{k+1} \tilde{\delta}^{k+1} (1-\tilde{\delta})^{n-k-1}}{\binom{n}{k} \tilde{\delta}^k (1-\tilde{\delta})^{n-k}} = \frac{(n-k)\tilde{\delta}}{(k+1)(1-\tilde{\delta})}.$$

To let the ratio $\frac{\mathbb{P}[Y=k+1]}{\mathbb{P}[Y=k]} \geq 1$, we need $(k+1) \leq (n+1)\tilde{\delta}$. Thus, the probability $\mathbb{P}[Y = l]$ achieves the maximum when $l = n\tilde{\delta}$. Consequently, we have

$$(n+1)\mathbb{P}[Y = n\tilde{\delta}] \geq 1 \quad \Leftrightarrow \quad \mathbb{P}[Y = n\tilde{\delta}] \geq \frac{1}{n+1}. \quad (46)$$

Taking the log and dividing by n on the both sides of the inequality (46), we obtain

$$\frac{1}{n} \log \binom{n}{m} \geq \tilde{\delta} \log(\alpha) - (1-\tilde{\delta}) \log(1-\alpha) - \frac{\log(n+1)}{n} = \phi(\tilde{\delta}) - \frac{\log(n+1)}{n}.$$

Last, we prove the lower bound equation (44).

Plugging equation (43) into equation (42), we have

$$\frac{1}{n} \log \mathbb{P}[Z_n \leq n\delta] \geq \phi(\tilde{\delta}) + \tilde{\delta} \log \alpha + (1-\tilde{\delta}) \log(1-\alpha) - \frac{\log(n+1)}{n} \geq -D(\delta \parallel \alpha) - \frac{\log(n+1)}{n}. \quad (47)$$

Multiplying by n and take exponential on both sides of the inequality (47), we obtain the result

$$\mathbb{P}[Z_n \leq n\delta] \geq \frac{1}{n+1} e^{-nD(\delta \parallel \alpha)}.$$

□

Remark 10. The lower bound on a binomial tail is also a function of the KL divergence between Bernoulli distributions with distinct parameters. Combining the upper bound (41) and lower bound (44), we conclude that the binomial tail is (upper and lower) bounded by the functions of KL divergence between Bernoulli distributions.

2.11 Exercise 2.11

Lemma 9 (Gaussian maxima). *Let $\{X_i\}_{i=1}^n$ be a sequence of i.i.d. normal random variables following $N(0, \sigma^2)$. Consider the random variable $Z_n := \max_{i \in [n]} |X_i|$. The expectation of Z_n is upper bounded as follows*

$$\mathbb{E}[Z_n] \leq \sqrt{2\sigma^2 \log n} + \frac{4\sigma}{\sqrt{2 \log n}}, \quad \text{for all } n \geq 2. \quad (48)$$

The lower bound of $\mathbb{E}[Z_n]$ is

$$\mathbb{E}[Z_n] \geq (1 - 1/e) \sqrt{\sigma^2 \log n}, \quad \text{for all } n \geq 6. \quad (49)$$

Moreover, we have

$$\frac{\mathbb{E}[Z_n]}{\sqrt{2\sigma^2 \log n}} \rightarrow 1, \quad \text{as } n \rightarrow +\infty. \quad (50)$$

The lower bound (49) is different to the original Exercise up to the constant.

Proof. First, we prove the equation (48).

Consider the tail probability of Z_n . We have

$$\mathbb{P}[Z_n \geq t] = \mathbb{P}[\max_{i \in [n]} |X_i| \geq t] = 1 - \mathbb{P}[\max_{i \in [n]} |X_i| < t] = 1 - (1 - \mathbb{P}[|X_1| \geq t])^n,$$

where the last equality follows from the independence of $\{X_i\}_{i=1}^n$. By Bernoulli's inequality, we have

$$(1 - \mathbb{P}[|X_1| \geq t])^n \geq 1 - n\mathbb{P}[|X_1| \geq t], \quad \text{for all } t > 0.$$

Therefore, the expectation of Z_n follows

$$\mathbb{E}[Z_n] = \int_0^{+\infty} \mathbb{P}[Z_n \geq t] dt \leq c + \int_c^{+\infty} n\mathbb{P}[|X_1| \geq t] dt, \quad \text{for all } c > 0. \quad (51)$$

Since the tail bound $\mathbb{P}[U \geq t] \leq \sqrt{\frac{2}{\pi}} \frac{1}{t} e^{-t^2/2}$ holds for all standard normal random variable U , the tail bound of $|X_1|$ satisfies

$$\mathbb{P}[|X_1| \geq t] \leq 2\sqrt{\frac{2}{\pi}} \frac{\sigma}{t} e^{-\frac{t^2}{2\sigma^2}}.$$

Let $u = \frac{t}{\sigma}$. Hence, for $c > \sigma$, the last integral in equation (51) is upper bounded as

$$\begin{aligned} \int_c^{+\infty} n\mathbb{P}[|X_1| \geq t] dt &\leq \frac{2n\sigma}{c} \sqrt{\frac{2}{\pi}} \int_c^{+\infty} e^{-\frac{t^2}{2\sigma^2}} dt \\ &\leq \frac{2n\sigma^2}{c} \sqrt{\frac{2}{\pi}} \int_{\frac{c}{\sigma}}^{+\infty} u e^{-\frac{u^2}{2}} du \\ &= \frac{2n\sigma^2}{c} \sqrt{\frac{2}{\pi}} e^{-\frac{c^2}{2\sigma^2}}. \end{aligned} \quad (52)$$

For all $n \geq 2$, let $c = \sqrt{2\sigma^2 \log n}$ and plug the c into equations (51) and (52). Then, we obtain the result

$$\mathbb{E}[Z_n] \leq \sqrt{2\sigma^2 \log n} + \frac{2\sigma}{\sqrt{2 \log n}} \sqrt{\frac{2}{\pi}} \leq \sqrt{2\sigma^2 \log n} + \frac{4\sigma}{\sqrt{2 \log n}}, \quad \text{for all } n \geq 2.$$

Alternative Way: An alternative proof for equation (48) is below.

Since $\{X_i\}_{i=1}^n$ are i.i.d. random normal variables following $N(0, \sigma^2)$, by Exercise 2.12, we have

$$\mathbb{E}[\max_{i \in [n]} X_i] \leq \sqrt{2\sigma^2 \log n}, \quad \text{for all } n \geq 1.$$

Note that $\max_{i \in [n]} |X_i| = \max_{i \in [2n]} X_i$, where $X_{n+i} = -X_i$, for all $i \in [n]$. Then, we have $\max_{i \in [n]} |X_i| \leq \sqrt{2\sigma^2 \log(2n)}$. Note that

$$\left(\sqrt{2\sigma^2 \log n} + \frac{4\sigma}{\sqrt{2 \log n}} \right)^2 = 2\sigma^2 \log n + 8\sigma^2 \left(1 + \frac{1}{\log n} \right) \geq 2\sigma^2 \log n + 2\sigma^2 \log 2 = (\sqrt{2\sigma^2 \log 2n})^2.$$

Therefore, we obtain the upper bound

$$\mathbb{E}[Z_n] \leq \sqrt{2\sigma^2 \log(2n)} \leq \sqrt{2\sigma^2 \log n} + \frac{4\sigma}{\sqrt{2 \log n}}.$$

Next, we prove the equation (49).

The expectation of Z_n is lower bounded as

$$\mathbb{E}[Z_n] \geq \sigma \sqrt{\log n} \mathbb{P}[Z_n > \sigma \sqrt{\log n}] = \sigma \sqrt{\log n} \left[1 - \left(1 - \mathbb{P}[|X_i| > \sigma \sqrt{\log n}] \right)^n \right].$$

Since $|X_i|$ follows the half-normal distribution, the upper tail probability of $|X_i|$ is

$$\mathbb{P}[|X_i| > \sigma \sqrt{\log n}] = 1 - \operatorname{erf}\left(\frac{\sqrt{\log n}}{\sqrt{2}}\right) \geq 1 - \sqrt{1 - n^{-\frac{2}{\pi}}},$$

where $\operatorname{erf}(\cdot)$ is the CDF function of half-normal distribution, and the last inequality follows from the fact that $\operatorname{erf}(x) \leq \sqrt{1 - \exp(-\frac{4}{\pi}x^2)}$. When $n \geq 6$, we have $1 - \sqrt{1 - n^{-\frac{2}{\pi}}} \geq \frac{1}{n}$, and thus the expectation satisfies

$$\mathbb{E}[Z_n] \geq \sigma \sqrt{\log n} \left[1 - \left(1 - \frac{1}{n} \right)^n \right] \geq \left(1 - \frac{1}{e} \right) \sigma \sqrt{\log n}, \quad \text{for all } n \geq 6.$$

Last, we prove the convergence property (50).

My thoughts: To prove $\lim_{n \rightarrow +\infty} \frac{\mathbb{E}[Z_n]}{\sqrt{2\sigma^2 \log n}} = 1$, we need to prove $\limsup_{n \rightarrow +\infty} \frac{\mathbb{E}[Z_n]}{\sqrt{2\sigma^2 \log n}} = \liminf_{n \rightarrow +\infty} \frac{\mathbb{E}[Z_n]}{\sqrt{2\sigma^2 \log n}} = 1$. The lim sup part is easy because of the upper bound (48). Then, we only need to prove the lim inf part.

Note that

$$\mathbb{E}[Z_n] \geq \sqrt{2\sigma^2 \log n} \mathbb{P}[Z_n > \sigma \sqrt{2 \log n}] = \sigma \sqrt{2 \log n} \left[1 - \left(1 - \mathbb{P}[|X_i| > \sigma \sqrt{2 \log n}] \right)^n \right].$$

If I can prove that $\lim_{n \rightarrow +\infty} \left(1 - \mathbb{P}[|X_i| > \sigma \sqrt{2 \log n}] \right)^n = 0$, the problem is done.

Note that $\mathbb{P}[|X_i| > \sigma \sqrt{2 \log n}] = 2\mathbb{P}[X_i > \sigma \sqrt{2 \log n}] \leq 2\sqrt{\frac{2}{\pi}} \frac{\sigma}{\sqrt{2\sigma^2 \log n}} e^{-\log n} = \frac{2}{\pi n \sqrt{\log n}}$. Then, the limit is upper bounded

$$\lim_{n \rightarrow +\infty} \left(1 - \mathbb{P}[|X_i| > \sigma \sqrt{2 \log n}] \right)^n \geq \lim_{n \rightarrow +\infty} \left(1 - \frac{2}{\pi n \sqrt{\log n}} \right)^n = 1.$$

I can not get the result I want.

□

Remark 11. The absolute Gaussian maxima Z_n concentrates around $\sigma \sqrt{\log n}$ for large n . Standard Gaussian variables concentrate around 0, while the maxima of the n standard Gaussian variables is no longer a Gaussian variable. The scale of Gaussian maxima scales as \sqrt{n} , however, the sum of n i.i.d. Gaussian variables scales as n .

2.12 Exercise 2.12

Lemma 10 (Sharp upper bounds for sub-Gaussian maxima). *Let $\{X_i\}_{i=1}^n$ be a sequence of zero-mean sub-Gaussian variables with parameter σ . Then, we have*

$$\mathbb{E} \left[\max_{i \in [n]} X_i \right] \leq \sqrt{2\sigma^2 \log n}, \quad \text{for all } n \geq 1.$$

Note that independence assumptions are unnecessary.

Proof. For all $t > 0$, the function $f(x) = \exp(tx)$ is a convex function. Apply the function $f(x)$ to $\mathbb{E}[\max_{i \in [n]} X_i]$. By Jensen's inequality, we have

$$\exp \left\{ t \mathbb{E} \left[\max_{i \in [n]} X_i \right] \right\} \leq \mathbb{E} \left[\exp \left(t \max_{i \in [n]} X_i \right) \right] \leq \sum_{i=1}^n \mathbb{E} [\exp(tX_i)] \leq n e^{\frac{t^2 \sigma^2}{2}}, \quad \text{for all } t > 0, \quad (53)$$

where the last inequality follows from the sub-Gaussianity of X_i s. Take the log of both sides of equation (53). We have

$$\mathbb{E} \left[\max_{i \in [n]} X_i \right] \leq \frac{\log n}{t} + \frac{t \sigma^2}{2}. \quad (54)$$

Plugging $t = \frac{\sqrt{2 \log n}}{\sigma}$ into the equation (54), we obtain

$$\mathbb{E} \left[\max_{i \in [n]} X_i \right] \leq \sqrt{2 \sigma^2 \log n}.$$

□

Remark 12. The maxima for n i.i.d. sub-Gaussian variables is upper bounded by $\sigma \sqrt{\log n}$. The independence assumption is not necessary for the upper bound of sub-Gaussian maxima.

2.13 Exercise 2.13

Lemma 11 (Operations on sub-Gaussian variables). *Let X_1 and X_2 be two zero-mean sub-Gaussian variables with parameters σ_1 and σ_2 respectively.*

- (a). *If X_1 and X_2 are independent, then random variable $X_1 + X_2$ is sub-Gaussian with parameter $\sqrt{\sigma_1^2 + \sigma_2^2}$.*
- (b). *Without independence assumptions, then random variable $X_1 + X_2$ is sub-Gaussian with parameter at most $\sigma_1 + \sigma_2$.*
- (c). *If X_1 and X_2 are independent, then random variable $X_1 X_2$ is sub-exponential with parameter $(\nu, b) = \left(\sqrt{2} \sigma_1 \sigma_2, \frac{1}{\sqrt{2} \sigma_1 \sigma_2} \right)$.*

Proof. First, we prove (a).

Since $X_1 \sim \text{subG}(\sigma_1)$ and $X_2 \sim \text{subG}(\sigma_2)$ are independent, we have

$$\mathbb{E} \left[e^{t(X_1 + X_2)} \right] = \mathbb{E} [e^{tX_1}] \mathbb{E} [e^{tX_2}] \leq e^{\frac{t^2(\sigma_1^2 + \sigma_2^2)}{2}}, \quad \text{for all } t \in \mathbb{R}.$$

Therefore, $X_1 + X_2$ is sub-Gaussian with parameter $\sqrt{\sigma_1^2 + \sigma_2^2}$.

Next, we prove (b).

For all $t \in \mathbb{R}$, we have

$$\mathbb{E} \left[e^{t(X_1 + X_2)} \right] = \mathbb{E} [e^{tX_1} e^{tX_2}]. \quad (55)$$

Introduce $p, q \geq 1$ such that $\frac{1}{p} + \frac{1}{q} = 1$. By Holder's inequality, the equation (55) becomes

$$\mathbb{E} \left[e^{t(X_1 + X_2)} \right] \leq \left(\mathbb{E} [e^{ptX_1}] \right)^{\frac{1}{p}} \left(\mathbb{E} [e^{qtX_2}] \right)^{\frac{1}{q}} \leq e^{\frac{pt^2 \sigma_1^2 + qt^2 \sigma_2^2}{2}}, \quad (56)$$

where the last inequality follows from the sub-Gaussianity of X_1 and X_2 . The inequality (56) achieves its minimum when $p = \sigma_2/\sigma_1 + 1$. Plugging the minimizer to inequality (56), we obtain

$$\mathbb{E} \left[e^{t(X_1+X_2)} \right] \leq e^{\frac{t^2(\sigma_1+\sigma_2)^2}{2}}.$$

Therefore, the random variable $X_1 + X_2$ is sub-Gaussian with parameter at most $\sigma_1 + \sigma_2$.

Last, we prove (c).

My thoughts: Note that $X_1 X_2 = \frac{1}{4} ((X_1 + X_2)^2 - (X_1 - X_2)^2)$. By part (a), we have $X_1 + X_2$ and $X_1 - X_2$ are sub-Gaussian with the parameter $\sqrt{\sigma_1^2 + \sigma_2^2}$. Note that for a variable $X \sim \text{subG}(\sigma^2)$, the variable $Z = X^2 - \mathbb{E}[X^2] \sim \text{subE}(16\sigma^2, 16\sigma^2)$.

The characteristic function of $X_1 X_2$ is

$$\mathbb{E}[e^{tX_1 X_2}] = \mathbb{E}[e^{\frac{t}{4}((X_1+X_2)^2 - (X_1-X_2)^2)}] = \mathbb{E}[e^{\frac{t}{4}(X_1+X_2)^2} e^{-\frac{t}{4}(X_1-X_2)^2}], \quad \text{for all } t \in \mathbb{R}. \quad (57)$$

By Cauchy-Schwartz inequality, the expectation (57) satisfies

$$\mathbb{E}[e^{\frac{t}{4}(X_1+X_2)^2} e^{-\frac{t}{4}(X_1-X_2)^2}] \leq \left(\mathbb{E}[e^{\frac{t}{2}(X_1+X_2)^2}] \right)^{1/2} \left(\mathbb{E}[e^{-\frac{t}{2}(X_1-X_2)^2}] \right)^{1/2}. \quad (58)$$

Let $Z_1 = (X_1 + X_2)^2 - \mathbb{E}[(X_1 + X_2)^2]$. Note that $\mathbb{E}[X_i^2] = \text{var}(X_i) \leq \sigma_i^2$, for $i = 1, 2$, and the variable $Z_1 \sim \text{subE}(16(\sigma_1^2 + \sigma_2^2), 16(\sigma_1^2 + \sigma_2^2))$. We have

$$\mathbb{E}[e^{\frac{t}{2}(X_1+X_2)^2}] = \mathbb{E}[e^{\frac{t}{2}Z_1 + \frac{t}{2}\mathbb{E}[(X_1+X_2)^2]}] \leq \mathbb{E}[e^{\frac{t}{2}Z_1}] e^{\frac{t}{2}(\sigma_1^2 + \sigma_2^2)} \leq e^{32t^2(\sigma_1^2 + \sigma_2^2)^2 + \frac{t}{2}(\sigma_1^2 + \sigma_2^2)}, \quad (59)$$

for all $|t| \leq \frac{1}{8(\sigma_1^2 + \sigma_2^2)}$. Similarly, we have

$$\mathbb{E}[e^{-\frac{t}{2}(X_1-X_2)^2}] \leq e^{32t^2(\sigma_1^2 + \sigma_2^2)^2 - \frac{t}{2}(\sigma_1^2 + \sigma_2^2)}, \quad (60)$$

for all $|t| \leq \frac{1}{8(\sigma_1^2 + \sigma_2^2)}$. Plug the inequalities (59) and (60) into the inequality (58). We have

$$\mathbb{E}[e^{\frac{t}{4}(X_1+X_2)^2} e^{-\frac{t}{4}(X_1-X_2)^2}] \leq e^{32t^2(\sigma_1^2 + \sigma_2^2)^2}, \quad \text{for all } |t| \leq \frac{1}{8(\sigma_1^2 + \sigma_2^2)}.$$

Therefore, the variable $X_1 X_2 \sim \text{subE}(8(\sigma_1^2 + \sigma_2^2), 8(\sigma_1^2 + \sigma_2^2))$. □

Remark 13. Without the independence assumption, the sum of sub-Gaussian variables is still a sub-Gaussian variable. Under the independence assumption, the product of two sub-Gaussian variables is a sub-exponential variable.

2.14 Exercise 2.14

Lemma 12. Let X be a scalar random variable. Suppose there exist two positive constant c_1, c_2 such that

$$\mathbb{P}[|X - \mathbb{E}[X]| \geq t] \leq c_1 e^{-c_2 t^2}, \quad \text{for all } t \geq 0. \quad (61)$$

Then, the variance $\text{var}(X) \leq c_1/c_2$.

Let m_X be the median of X ; i.e., $\mathbb{P}[X \geq m_X] \geq 1/2$ and $\mathbb{P}[X \leq m_X] \geq 1/2$. Note that the median m_X is not necessarily unique. Given mean concentration (61), for all median m_X , we have

$$\mathbb{P}[|X - m_X| \geq t] \leq c_3 e^{-c_4 t^2}, \quad \text{for all } t \geq 0, \quad (62)$$

where $c_3 := 4c_1$ and $c_4 := c_2/8$.

Conversely, if the equation (62) holds, the mean concentration (61) also holds with parameter $c_1 = 2c_3$ and $c_2 = c_4/4$.

Proof. First, we show the variance $\text{var}(X) \leq c_1/c_2$.

By the definition of variance, we have

$$\text{var}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2] = \int_0^{+\infty} \mathbb{P}[|X - \mathbb{E}[X]| \geq \sqrt{t}] dt. \quad (63)$$

Since X satisfies the mean concentration equation (61), for all $t > 0$, we have

$$\mathbb{P}[|X - \mathbb{E}[X]| \geq \sqrt{t}] \leq c_1 e^{-c_2 t}. \quad (64)$$

Plugging the equation (64) into equation (63), we obtain

$$\text{var}(X) \leq \frac{c_1}{c_2}.$$

Next, we show by an example that the median of X , m_X , is not necessarily unique.

Consider a random variable $Y \sim \text{Bin}(5, 1/2)$. Then, we have $\mathbb{P}[Y \geq 3] = 1/2$ and $\mathbb{P}[Y \leq 3] > 1/2$. Meanwhile, we also have $\mathbb{P}[Y \geq 2] > 1/2$ and $\mathbb{P}[Y \leq 2] = 1/2$. By the definition of median, the number 2 and 3 both are the median of Y . Therefore, the median of a random variable is not necessarily unique.

Then, we prove the median concentration (62), provided that random variable X satisfies the mean concentration (61).

Let $\Delta = |\mathbb{E}[X] - m_X|$. Consider the following two cases.

1. Case 1: Suppose $t \geq 2\Delta$.

Note the triangle inequality $|X - \mathbb{E}[X]| \geq |X - m_X| - \Delta$ and the assumption $\frac{t}{2} \geq \Delta$. We have

$$\mathbb{P}[|X - m_X| \geq t] \leq \mathbb{P}\left[|X - m_X| \geq \frac{t}{2} + \Delta\right] \leq \mathbb{P}\left[|X - \mathbb{E}[X]| \geq \frac{t}{2}\right].$$

By equation (61), we obtain

$$\mathbb{P}[|X - m_X| \geq t] \leq c_1 e^{-\frac{c_2 t^2}{4}}. \quad (65)$$

2. Case 2: Suppose $t < 2\Delta$.

By the definition of median, we have

$$\mathbb{P}[|X - \mathbb{E}[X]| \geq \Delta] \geq \mathbb{P}[X \geq m_X] \geq \frac{1}{2}. \quad (66)$$

Meanwhile, by equation (61), we have

$$\mathbb{P}[|X - \mathbb{E}[X]| \geq \Delta] \leq c_1 e^{-c_2 \Delta^2}. \quad (67)$$

Combing the assumption $\Delta > \frac{t}{2}$ and inequalities (66) (67), we obtain

$$\mathbb{P}[|X - m_X| \geq t] \leq 1 \leq 2c_1 e^{-c_2 \Delta^2} \leq 2c_1 e^{-\frac{c_2 t^2}{4}}. \quad (68)$$

Hence, combining the inequality (65) and inequality (68), we conclude that

$$\mathbb{P}[|X - m_X| \geq t] \leq c_3 e^{-c_4 t^2},$$

where $c_3 = 2c_1$ and $c_4 = c_2/4$. The upper bound also holds when $c_3 = 4c_1$ and $c_4 = c_2/8$.

Last, we prove the mean concentration (61), provided that the random variable X satisfies the median concentration (62). □

Remark 14. If a variable concentrates around mean, the variable also concentrates around the median, vice versa.

2.15 Exercise 2.15

Lemma 13. Let $\{X_i\}_{i=1}^n$ be i.i.d. sample of random variables with density f on the real line. A standard estimate of f is the kernel density estimate

$$\hat{f}_n(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right),$$

where $K: \mathbb{R} \mapsto [0, +\infty)$ is a kernel function satisfying $\int_{-\infty}^{+\infty} K(x)dx = 1$, and $h > 0$ is a bandwidth parameter. Suppose we assess the estimate $\hat{f}_n(x)$ using the L^1 norm; i.e.,

$$\|\hat{f}_n - f\|_1 = \int_{-\infty}^{+\infty} |\hat{f}_n(t) - f(t)| dt.$$

The upper tail probability of the L^1 norm satisfies

$$\mathbb{P}\left[\|\hat{f}_n - f\|_1 - \mathbb{E}\left[\|\hat{f}_n - f\|_1\right] \geq \delta\right] \leq e^{-\frac{n\delta^2}{8}}, \quad \text{for all } \delta \geq 0.$$

Proof. Note that $\|\hat{f}_n - f\|_1$ is a function of $X = (X_1, \dots, X_n)$, denoted $g(X) = g(X_1, \dots, X_n) = \|\hat{f}_n - f\|_1$ and $g: \mathbb{R}^n \mapsto \mathbb{R}$. Let $x^{(k)} = (x_1, \dots, x_k, \dots, x_n) \in \mathbb{R}^n$ and $x'^{(k)} = (x_1, \dots, x'_k, \dots, x_n) \in \mathbb{R}^n$ be two vectors. The absolute difference between the $g(x^{(k)})$ and $g(x'^{(k)})$ is

$$\begin{aligned} |g(x^{(k)}) - g(x'^{(k)})| &= \left| \int_{-\infty}^{+\infty} \left| \frac{1}{nh} \sum_{i=1}^n K\left(\frac{t - x_i}{h}\right) - f(t) \right| dt - \int_{-\infty}^{+\infty} \left| \frac{1}{nh} \sum_{i \neq k}^n K\left(\frac{t - x_i}{h}\right) - f(t) + \frac{1}{nh} K\left(\frac{t - x'_k}{h}\right) \right| dt \right| \\ &\leq \int_{-\infty}^{+\infty} \left| \frac{1}{nh} \left(K\left(\frac{t - x_k}{h}\right) - K\left(\frac{t - x'_k}{h}\right) \right) \right| dt \\ &\leq \frac{1}{nh} \left(\int_{-\infty}^{+\infty} K\left(\frac{t - x_k}{h}\right) dt + \int_{-\infty}^{+\infty} K\left(\frac{t - x'_k}{h}\right) dt \right). \end{aligned} \tag{69}$$

Let $t_1 = \frac{t - x_k}{h}$. By the definition of K , the first integral in equation (69) are

$$\int_{-\infty}^{+\infty} K\left(\frac{t - x_k}{h}\right) dt = h \int_{-\infty}^{+\infty} K(t_1) dt_1 = h.$$

Similarly, the second integral $\int_{-\infty}^{+\infty} K\left(\frac{t-x'_k}{h}\right) dt = h$. Then, we conclude the function g satisfies

$$|g(x^{(k)}) - g(x'^{(k)})| \leq \frac{2}{n}. \quad (70)$$

The inequality (70) also holds for all $k \in [n]$ and for all $x^{(k)}, x'^{(k)} \in \mathbb{R}^n$. Therefore, $g(X)$ satisfies the bounded difference property (8) with parameters $(\frac{2}{n}, \dots, \frac{2}{n})$. By Theorem 1.9, we obtain

$$\mathbb{P}\left[\left\|\hat{f}_n - f\right\|_1 - \mathbb{E}\left[\left\|\hat{f}_n - f\right\|_1\right] \geq \delta\right] \leq e^{-\frac{n\delta^2}{2}} \leq e^{-\frac{n\delta^2}{8}}.$$

□

Remark 15. Kernel density estimator satisfies the bounded difference property and possesses a concentration property. Kernel density estimator converges to its mean in probability at the speed of e^{-n} when n goes to infinity.

2.16 Exercise 2.16

Lemma 14. Let $\{X_i\}_{i=1}^n$ be a sequence of independent variables taking values in a Hilbert space \mathbb{H} . Suppose that $\|X_i\|_{\mathbb{H}} \leq b_i$ almost surely, for all $i \in [n]$. Consider the real valued random variable $S_n = \|\sum_{i=1}^n X_i\|_{\mathbb{H}}$. The concentration bound of S_n is

$$\mathbb{P}[|S_n - \mathbb{E}[S_n]| \geq n\delta] \leq 2e^{-\frac{n\delta^2}{8b^2}}, \quad \text{for all } \delta \geq 0, \quad (71)$$

where $b^2 = \frac{1}{n} \sum_{i=1}^n b_i^2$. The upper tail probability bound of S_n is

$$\mathbb{P}\left[\frac{S_n}{n} \geq a + \delta\right] \leq e^{-\frac{n\delta^2}{8b^2}}, \quad \text{for all } \delta \geq 0, \quad (72)$$

where $a = \sqrt{\frac{1}{n^2} \sum_{i=1}^n \mathbb{E}[\|X_i\|_{\mathbb{H}}^2]}$.

Proof. First, we prove the concentration bound (71).

Note that S_n is a function of $X = (X_1, \dots, X_n)$, denoted $S_n(X) = S_n(X_1, \dots, X_n)$, and $S_n(X): \mathbb{R}^n \mapsto \mathbb{R}$. Let $x^{(k)} = (x_1, \dots, x_k, \dots, x_n) \in \mathbb{R}^n$ and $x'^{(k)} = (x_1, \dots, x'_k, \dots, x_n) \in \mathbb{R}^n$ be two vectors. The absolute difference between $S_n(x^{(k)})$ and $S_n(x'^{(k)})$ is

$$\begin{aligned} \left|S_n(x^{(k)}) - S_n(x'^{(k)})\right| &= \left|\|x_1 + \dots + x_k + \dots + x_n\|_{\mathbb{H}} - \|x_1 + \dots + x'_k + \dots + x_n\|_{\mathbb{H}}\right| \\ &\leq \|x_k - x'_k\|_{\mathbb{H}} \\ &\leq 2b_k, \end{aligned} \quad (73)$$

where the last inequality follows by the boundedness $\|X_k\|_{\mathbb{H}} \leq b_k$, for all $k \in [n]$. The inequality (73) also holds for all $k \in [n]$ and for all $x^{(k)}, x'^{(k)} \in \mathbb{R}^n$. Therefore, $S_n(X)$ satisfies the bounded property (8) with parameters $(2b_1, \dots, 2b_n)$. By Theorem 1.9, we obtain

$$\mathbb{P}[|S_n - \mathbb{E}[S_n]| \geq n\delta] \leq 2e^{-\frac{n\delta^2}{2b^2}} \leq 2e^{-\frac{n\delta^2}{8b^2}}, \quad \text{for all } \delta \geq 0.$$

Next, we prove the upper tail probability bound (72).

The expectation of $S_n(X) = \mathbb{E}[\sum_{i=1}^n \|X_i\|_{\mathbb{H}}]$ satisfies

$$\mathbb{E} \left[\left\| \sum_{i=1}^n X_i \right\|_{\mathbb{H}} \right] \leq \sqrt{\mathbb{E} \left[\left\| \sum_{i=1}^n X_i \right\|_{\mathbb{H}}^2 \right]} = \sqrt{\sum_{i=1}^n \mathbb{E} [\|X_i\|_{\mathbb{H}}^2]} = na, \quad (74)$$

where the first inequality follows by the fact that $\mathbb{E}[X^2] \geq (\mathbb{E}[X])^2$ and the second equation follows by the independence of $\{X_i\}_{i=1}^n$. Therefore, combining the inequality (74) with the concentration bound (71), we obtain

$$\mathbb{P} \left[\frac{S_n}{n} \geq a + \delta \right] = \mathbb{P} [S_n - na \geq n\delta] \leq \mathbb{P} [S_n - \mathbb{E}[S_n] \geq n\delta] \leq e^{-\frac{n\delta^2}{8b^2}}.$$

□

Remark 16. The sample mean also for n variables in a Hilbert space converge to its expectation in probability at the speed of e^{-n} .

2.17 Exercise 2.17

Lemma 15 (Hanson-Wright Inequality). *Let $\mathbf{Q} = [\mathbf{Q}_{ij}] \in \mathbb{R}^{n \times n}$ be a positive semi-definite matrix, and $\{X_i\}_{i=1}^n$ be i.i.d. random variables with mean zero, variance 1 and σ -sub-Gaussian. Consider the random variable $Z = \sum_{i=1}^n \sum_{j=1}^n \mathbf{Q}_{ij} X_i X_j$. The Hanson-Wright inequality for Z is*

$$\mathbb{P} [Z \geq \text{Trace}(\mathbf{Q}) + \sigma t] \leq 2e^{-\min \left\{ \frac{c_1 t}{\|\mathbf{Q}\|_{op}}, \frac{c_2 t^2}{\|\mathbf{Q}\|_F^2} \right\}}, \quad \text{for all } t \geq 0,$$

where $\|\mathbf{Q}\|_{op}$ is the operator norm of matrix \mathbf{Q} , $\|\mathbf{Q}\|_F$ is the Frobenius norm of matrix \mathbf{Q} , and c_1, c_2 are two universal constants independent with n and t .

The proof for Hanson-Wright Inequality under the special case $X_i \stackrel{i.i.d.}{\sim} N(0, 1)$ is below. Let $\sigma = 1$.

Proof. Since the random variables $\{X_i\}_{i=1}^n$ are independent and $\text{var}(X_i) = 1$ for all $i \in [n]$, the expectation $\mathbb{E}[Z] = \sum_{i=1}^n \mathbf{Q}_{ii} = \text{Trace}(\mathbf{Q})$. Then, the upper tail probability of Z becomes

$$\mathbb{P} [Z \geq \text{Trace}(\mathbf{Q}) + t] = \mathbb{P} \left[\left(\sum_{i=1}^n \mathbf{Q}_{ii} (X_i^2 - 1) + \sum_{i \neq j} \mathbf{Q}_{ij} X_i X_j \right) \geq t \right]$$

For two random variables X, Y and a constant t , we have the inequality $\mathbb{P}[X + Y \geq t] \leq \mathbb{P}[X \geq t/2] + \mathbb{P}[Y \geq t/2]$. Hence, the tail probability is upper bounded as

$$\mathbb{P} [Z \geq \text{Trace}(\mathbf{Q}) + t] \leq \mathbb{P} \left[\sum_{i=1}^n \mathbf{Q}_{ii} (X_i^2 - 1) \geq \frac{t}{2} \right] + \mathbb{P} \left[\sum_{i \neq j} \mathbf{Q}_{ij} X_i X_j \geq \frac{t}{2} \right] =: p_1 + p_2.$$

Consider p_1 . For a standard normal variable X , the squared variable $X^2 \sim \chi_1^2$ and thus X^2 is sub-exponential with parameter $(\nu, b) = (2, 4)$. Given a constant c , cX^2 is still sub-exponential with parameters $(2c, 4c^2)$. Hence, for all $i \in [n]$, the variable $\mathbf{Q}_{ii} X_i^2$ is sub-exponential with parameters

$(2\mathbf{Q}_{ii}, 4\mathbf{Q}_{ii}^2)$. By the Bernstein-type inequalities (see Step 1 in reference), the probability p_1 is upper bounded as following,

$$p_1 = \mathbb{P} \left[\sum_{i=1}^n (\mathbf{Q}_{ii} X_i^2 - \mathbf{Q}_{ii} \mathbb{E}[X_i^2]) \geq \frac{t}{2} \right] \leq \exp \left(- \min \left\{ \frac{ct^2}{\sum_i \mathbf{Q}_{ii}^2}, \frac{c't}{\max_{i \in [n]} |\mathbf{Q}_{ii}|} \right\} \right),$$

where c and c' are two constants independent with n and t .

Consider p_2 . For two i.i.d. standard normal variables X, Y , the product $XY = \frac{1}{4}(X+Y)^2 - \frac{1}{4}(X-Y)^2$, where $\frac{1}{2}(X+Y)^2$ and $\frac{1}{2}(X-Y)^2$ follow the chi-square distribution χ_1^2 independently. Then, XY is the subtraction of two independent sub-exponential variables. Let D_1 and D_2 be two identical and independent sub-exponential variables with parameters (ν, b) . We have

$$\mathbb{E}[e^{\lambda(D_1-D_2)}] = \mathbb{E}[e^{\lambda D_1}] \mathbb{E}[e^{-\lambda D_2}] \leq e^{2\nu^2 \lambda^2}, \quad \text{for all } |\lambda| < \frac{1}{b}.$$

Therefore, the subtraction of two identical and independent sub-exponential variables $D_1 - D_2$ is still sub-exponential with parameters $(\sqrt{2}\nu, b)$.

Back to $\mathbf{Q}_{ij} X_i X_j$. By the discussion above, for all $i \neq j$, the variable $\mathbf{Q}_{ij} X_i X_j$ is sub-exponential with parameters $(\sqrt{2}\mathbf{Q}_{ij}, \mathbf{Q}_{ij}^2)$. By the Bernstein-type inequalities, the probability p_2 is upper bounded as following,

$$p_2 = \mathbb{P} \left[\sum_{i \neq j} \mathbf{Q}_{ij} X_i X_j \geq \frac{t}{2} \right] \leq \exp \left(- \min \left\{ \frac{c''t^2}{\sum_{i \neq j} \mathbf{Q}_{ij}^2}, \frac{c'''t}{\max_{i \neq j} |\mathbf{Q}_{ij}|} \right\} \right),$$

where c'' and c''' are two constants independent with n and t .

Therefore, by the definition of Frobenius norm and the fact that $\max_{i,j \in [n]} |\mathbf{Q}_{ij}| \leq \|\mathbf{Q}\|_{op}$, we obtain the goal inequality

$$\mathbb{P}[Z \geq \text{Trace}(\mathbf{Q}) + t] \leq p_1 + p_2 \leq 2 \exp \left(- \min \left\{ \frac{c_2 t^2}{\|\mathbf{Q}\|_F^2}, \frac{c_1 t}{\|\mathbf{Q}\|_{op}} \right\} \right),$$

for some universal constants c_1, c_2 independent with n and t . □

Reference: <http://www-personal.umich.edu/~rudelson/papers/rv-Hanson-Wright.pdf>.

Remark 17. For a random vector X , Hanson-Wright Inequality implies the tail probability of matrix product $X^T \mathbf{Q} X$ is upper bounded at the level of $e^{-\min\{\|\mathbf{Q}\|_F^{-2}, \|\mathbf{Q}\|_{op}^{-1}\}}$, where \mathbf{Q} is a positive semi-definite matrix.

2.18 Exercise 2.18

Definition 3 (Orlicz Norms). Let $\psi: \mathbb{R}_+ \mapsto \mathbb{R}_+$ be a strictly increasing convex function satisfying $\psi(0) = 0$. The ψ -Orlicz norm of a random variable X is defined as

$$\|X\|_\psi := \inf \{t > 0: \mathbb{E}[\psi(t^{-1}|X|)] \leq 1\}.$$

The norm $\|X\|_\psi$ is infinite if there is no finite t such that the expectation $\mathbb{E}[\psi(t^{-1}|X|)]$ exists.

Given the function $\psi_q(u) = u^q$ for some $q \in [1, +\infty)$, the Orlicz norm is simply the l_q -norm $\|X\|_q = (\mathbb{E}[|X|^q])^{1/q}$.

Lemma 16. Let $\psi_q(u) = \exp(u^q) - 1$ be a function for $u \in \mathbb{R}_+$ with positive parameter $q \in [1, +\infty)$, and $\|\cdot\|_{\psi_q}$ denote the ψ_q -Orlicz norm. Then,

(a.) if $\|X\|_{\psi_q} < +\infty$, there exist two positive constants c_1, c_2 such that

$$\mathbb{P}[|X| > t] \leq c_1 \exp(-c_2 t^q), \quad \text{for all } t > 0. \quad (75)$$

In particular, the constants $c_1 = 2$ and $c_2 = \|X\|_{\psi_q}$.

(b.) if the random variable X satisfies the tail bound (75), the Orlicz norm $\|X\|_{\psi_q} < +\infty$.

Proof. First, we prove (a).

Let t_0 denote the Orlicz norm $\|X\|_{\psi_q}$. Since $t_0 < +\infty$, the expectation $\mathbb{E}[\exp(\lambda|X|^q)]$ exists for $\lambda \in [0, t_0^{-q}]$. Applying Chernoff's bound to $|X|^q$, the upper tail bound of $|X|$ is

$$\mathbb{P}[|X| > t] = \mathbb{P}\left[\frac{|X|^q}{t_0^q} > \frac{t^q}{t_0^q}\right] \leq e^{-\lambda t^q} \mathbb{E}[\exp(\lambda|X|^q)], \quad \text{for all } t \geq 0, \lambda \in [0, t_0^{-q}]. \quad (76)$$

By the definition of Orlicz norm, the expectation satisfies

$$\mathbb{E}\left[\exp\left(t_0^{-q}|X|^q\right)\right] = \mathbb{E}\left[\psi(t_0^{-1}|X|)\right] + 1 \leq 2. \quad (77)$$

Let $\lambda = t_0^{-q}$ and plug the expectation (77) to the inequality (76). Then, we obtain the tail bound

$$\mathbb{P}[|X| > t] \leq 2e^{-\|X\|_{\psi_q}^{-q} t^q}, \quad \text{for all } t \geq 0.$$

Next, we prove (b).

Let t_1 be a positive constant. Note the fact that $\mathbb{E}[X] = \int_{-\infty}^{+\infty} \mathbb{P}(X > t) dt$. The expectation of $\mathbb{E}[\exp(t_1^{-q}|X|^q)]$ is

$$\begin{aligned} \mathbb{E}[\exp(t_1^{-q}|X|^q)] &= \int_0^{+\infty} \mathbb{P}\left[\exp\left(\frac{|X|^q}{t_1^q}\right) > \exp\left(\frac{t^q}{t_1^q}\right)\right] dt \\ &= \int_0^{+\infty} \mathbb{P}[|X| > t] dt \\ &\leq \int_0^{+\infty} c_1 \exp(-c_2 t^q) dt, \end{aligned} \quad (78)$$

where the last inequality follows from the tail bound (75). Let $u = c_2 t^q$. The integral in equation (78) becomes

$$\int_0^{+\infty} c_1 \exp(-c_2 t^q) dt = \frac{c_1}{qc_2^{1/q}} \int_0^{+\infty} u^{1/q-1} e^{-u} du = \frac{c_1}{qc_2^{1/q}} \Gamma(1/q) < +\infty.$$

Therefore, the expectation

$$\mathbb{E}\left[\psi(t_1^{-1}|X|)\right] = \mathbb{E}[\exp(t_1^{-q}|X|^q)] - 1 < +\infty, \quad \text{for all } t_1 > 0,$$

which implies the Orlicz norm $\|X\|_{\psi_q}$ is definite. \square

Remark 18. For a random variable X , the existence of X 's Orlicz norm implies the tail probability of X is upper bounded at the level of $e^{-\|X\|_{\psi_q}^{-q}}$.

2.19 Exercise 2.19

Lemma 17 (Maxima of Orlicz variables). *Recall the definition of Orlicz norm 3. Let $\{X\}_{i=1}^n$ be a sequence of i.i.d. zero-mean random variables. Given a positive, increasing and convex function $\psi: \mathbb{R}_+ \mapsto \mathbb{R}_+$, the random variable X_i has a finite Orlicz norm $\sigma = \|X_i\|_\psi$, for all $i \in [n]$. Then, the maxima of X_i satisfies*

$$\mathbb{E} \left[\max_{i \in [n]} |X_i| \right] < \sigma \psi^{-1}(n).$$

Proof. Since ψ is a convex function, by Jensen's inequality, we have

$$\psi \left(\mathbb{E} \left[\max_{i \in [n]} \frac{|X_i|}{\sigma} \right] \right) \leq \mathbb{E} \left[\max_{i \in [n]} \psi \left(\frac{|X_i|}{\sigma} \right) \right] \leq n \mathbb{E} \left[\psi \left(\frac{|X_1|}{\sigma} \right) \right]. \quad (79)$$

By the definition of Orlicz norm, we have $\mathbb{E} \left[\psi \left(\frac{|X_1|}{\sigma} \right) \right] \leq 1$. Therefore, taking the inverse function ψ^{-1} on both sides of inequality (79), we obtain

$$\mathbb{E} \left[\max_{i \in [n]} |X_i| \right] \leq \sigma \psi^{-1}(n).$$

□

Remark 19. The maxima of Orlicz variables is upper bounded by the Orlicz norm σ and scales as $\psi^{-1}(n)$. In previous exercises, the maxima of sub-Gaussian variables concentrates around the sub-Gaussian variable σ and scales at $\log(n)$.

2.20 Exercise 2.20

Lemma 18. *Let $\{X_i\}_{i=1}^n$ be a sequence of independent zero-mean random variables. Suppose that for some fixed integer $m \geq 1$ the random variables satisfy the following moment bound,*

$$\|X_i\|_{2m} := (\mathbb{E}[X_i^{2m}])^{\frac{1}{2m}} \leq C_m, \quad \text{for all } i \in [n].$$

Then, we have the tail bound

$$\mathbb{P} \left[\left| \frac{1}{n} \sum_{i=1}^n X_i \right| \geq \delta \right] \leq B_m \left(\frac{1}{\sqrt{n}\delta} \right)^{2m}, \quad \text{for all } \delta > 0,$$

where B_m is a universal constant depending on C_m and m only.

Hint: Let $\{X_i\}_{i=1}^n$ be a sequence of independent zero-mean random variables. The Rosenthal's inequality is

$$\mathbb{E} \left[\left(\sum_{i=1}^n X_i \right)^{2m} \right] \leq R_m \left\{ \sum_{i=1}^n \mathbb{E}[X_i^{2m}] + \left(\sum_{i=1}^n \mathbb{E}[X_i^2] \right)^m \right\}, \quad (80)$$

where R_m is a universal constant depending on m only.

Proof. Applying the $2m$ -order central moment version of the Markov's inequality to $\left| \frac{1}{n} \sum_{i=1}^n X_i \right|$, we have the tail bound

$$\mathbb{P} \left[\left| \frac{1}{n} \sum_{i=1}^n X_i \right| \geq \delta \right] = \mathbb{P} \left[\left| \sum_{i=1}^n X_i \right| \geq n\delta \right] \leq (n\delta)^{-2m} \mathbb{E} \left[\left(\sum_{i=1}^n X_i \right)^{2m} \right]. \quad (81)$$

By the moment bound, for all $i \in [n]$, we have

$$\mathbb{E}[X_i^2] \leq \mathbb{E}[X_i^{2m}] \leq C_m^{2m}. \quad (82)$$

Plugging the inequality (82) to the Rosenthal's inequality (80), we obtain

$$\mathbb{E} \left[\left(\sum_{i=1}^n X_i \right)^{2m} \right] \leq R_m \left(n C_m^{2m} + n^m C_m^{2m^2} \right). \quad (83)$$

Combining the inequality (83) with tai bound (81), we have

$$\mathbb{P} \left[\left| \frac{1}{n} \sum_{i=1}^n X_i \right| \geq \delta \right] \leq \frac{R_m \left(n C_m^{2m} + n^m C_m^{2m^2} \right)}{n^{2m} \delta^{2m}} = B_m \left(\frac{1}{\sqrt{n} \delta} \right)^{2m}, \quad \text{for all } \delta > 0,$$

where B_m is a universal constant depending on C_m and m only. \square

Remark 20. For n independent zero-mean variables X_i , the order $2m$ moment condition for X_i implies that the sample mean $\frac{1}{n} \sum_{i=1}^n X_i$ converges to 0 at the speed of n^{-m} .

2.21 Exercise 2.21

Definition 4. Let $X = (X_1, \dots, X_n)$ be a sequence of i.i.d. Bernoulli variables with parameter $1/2$. Represent X using a codebook of N binary vectors, denoted $\{z^1, \dots, z^N\}$, where $z^j = \llbracket z_i^j \rrbracket \in \mathbb{R}^n$ is a vector, for all $j \in [N]$. The rescaled Hamming distortion is defined as

$$\delta := \mathbb{E} \left[\min_{j \in [N]} \rho_H(X, z^j) \right] = \mathbb{E} \left[\min_{j \in [N]} \frac{1}{n} \sum_{i=1}^n \mathbb{1} [X_i \neq z_i^j] \right].$$

The goal of lossy data compression is to find a representation of X with rescaled Hamming distortion as small as possible. When we use a codebook with $N = 2^n$ codewords, we always achieve the 0 distortion. However, obtaining the distortion with $N = 2^n$ codewords is computationally expensive. Our goal is to find a proper rate $R < 1$ such that using a codebook with $N = 2^{Rn}$ codewords achieves a good compression performance.

Lemma 19. Let $X = (X_1, \dots, X_n)$ be a sequence of i.i.d. Bernoulli variables with parameter $1/2$. Let $\delta \in (0, 1)$. Suppose that the rate R is upper bounded as

$$R < D_2(\delta || 1/2) = \delta \log_2 \frac{\delta}{1/2} + (1 - \delta) \log_2 \frac{1 - \delta}{1/2}.$$

Let $\Delta_R := R - D_2(\delta || 1/2)$. If $\Delta_R < 0$, then the probability of achieving distortion δ goes to 0 as n goes to the infinity, for all codebook $\{z^1, \dots, z^N\}$. If $\Delta_R > 0$, the probability of achieving distortion δ goes to 1 as n goes to the infinity, for all codebook $\{z^1, \dots, z^N\}$.

Proof. First, we prove the probability of achieving the distortion under $\Delta_R < 0$.

Consider the random variables $V^j = \mathbb{1} [\rho_H(X, z^j) \leq \delta]$, for all $j \in [N]$. Given the codeword z^j , the variable $\mathbb{1} [X_i - z_i^j]$ is still a Bernoulli random variable. Recalling the tail bound (41) in Exercise 2.9, we have

$$\mathbb{P} [V^j = 1] = \mathbb{P} \left[\sum_{i=1}^n \mathbb{1} [X_i - z_i^j] \leq n\delta \right] \leq e^{-nD(\delta || 1/2)}, \quad \text{for all } j \in [N].$$

Since the sum of V^j , denoted $V = \sum_{j=1}^N V^j$, follows the Binomial distribution $\text{Bin}(N, \mathbb{P}[V^j = 1])$, the probability of not achieving the distortion δ is

$$\mathbb{P}[V = 0] \geq \left(1 - e^{-nD(\delta||1/2)}\right)^N. \quad (84)$$

Let $R_0 = D_2(\delta||1/2)$ and $N_0 = 2^{R_0 n} = (2\delta)^{n\delta}(2(1-\delta))^{n(1-\delta)}$. By the definition of KL divergence, we have $e^{-nD(\delta||1/2)} = N_0^{-1}$. Note that the ratio $\frac{N}{N_0} = 2^{(\Delta_R)n}$ goes to 0 as n goes to the infinity. Hence, the limit of the probability (84) is

$$\lim_{n \rightarrow +\infty} \mathbb{P}[V = 0] \geq \lim_{n \rightarrow +\infty} \left(\left(1 - \frac{1}{N_0}\right)^{N_0} \right)^{\frac{N}{N_0}} = 1.$$

Therefore, when $\Delta_R < 0$, the probability of achieving distortion δ goes to 0 as n goes to the infinity.

Next, we prove the probability of achieving the distortion under $\Delta_R > 0$.

Recall the tail bound (44) in Exercise 2.10. We have the upper bound of $\mathbb{P}[V^j = 1]$, where

$$\mathbb{P}[V^j = 1] = \mathbb{P}\left[\sum_{i=1}^n \mathbb{1}[X_i - z_i^j] \leq n\delta\right] \geq \frac{1}{n+1} e^{-nD(\delta||1/2)}, \quad \text{for all } j \in [N]. \quad (85)$$

Then the probability of not achieving the distortion δ is

$$\mathbb{P}[V = 0] \leq \left(1 - \frac{1}{n+1} e^{-nD(\delta||1/2)}\right)^N.$$

Since the difference $\Delta_R > 0$, then the ratio $\frac{N}{N_0(n+1)} = \frac{2^{(\Delta_R)n}}{n+1}$ goes to infinity as n goes to infinity. Hence, the limit of the probability of not achieving the distortion δ is

$$\lim_{n \rightarrow +\infty} \mathbb{P}[V = 0] \leq \lim_{n \rightarrow +\infty} \left(1 - \frac{1}{n+1} e^{-nD(\delta||1/2)}\right)^N = \lim_{n \rightarrow +\infty} \left(\left(1 - \frac{1}{(n+1)N_0}\right)^{(n+1)N_0} \right)^{\frac{N}{(n+1)N_0}} = 0.$$

Therefore, when $\Delta_R > 0$, the probability of achieving distortion δ goes to 1 as n goes to the infinity.

An alternative version of the proof under $\Delta_R > 0$ is below.

Let $p = \mathbb{P}[V^j = 1]$. By the Bernoulli's inequality, the probability of not achieving the distortion δ satisfies

$$\mathbb{P}[V = 0] = (1 - p)^N \leq \frac{1}{1 - Np}$$

Therefore, by inequality (85), the limit of the probability of achieving the distortion δ is

$$\lim_{n \rightarrow +\infty} \mathbb{P}[V \geq 1] \geq \lim_{n \rightarrow +\infty} 1 - \frac{1}{1 - Np} \geq \lim_{n \rightarrow +\infty} 1 - \frac{N_0(n+1)}{N_0(n+1) - N} = 1.$$

□

Remark 21. To compress a random vector with Bernoulli entries and achieve a desired rescaled Hamming distortion δ , the number of necessary codewords, N , should be larger than $2^{D(\delta||1/2)n}$, where $D(\delta||1/2)$ is the KL divergence between the two Bernoulli distributions with parameter δ and $1/2$ and n is the size of the random vector.