

Notes for High-Dimensional Probability*

Vishwak Srinivasan

Contents

1 Preliminaries	2
1.1 Example on approximate Caratheodory's Theorem	2
1.1.1 Auxiliary Lemmata	3
1.2 Quantities and Inequalities associated with RVs	4
1.3 Basic Limit Theorems	6
2 Concentration inequalities	8
2.1 Basic Gaussian Inequalities	8
2.1.1 Auxiliary Lemmata	9
2.2 Hoeffding's Inequality	9
2.2.1 Auxiliary Lemmata	13
2.3 Chernoff's Inequality	14
2.3.1 Auxiliary Lemmata	16
2.4 Sub-Gaussian Random Variables	17
2.4.1 Auxiliary Lemmata	24
2.5 Hoeffding's and Khintchine Inequalities for sub-Gaussian random variables, Centering	25
2.5.1 Auxiliary Lemmata	28
2.6 Sub-Exponential Random Variables, Orlicz spaces	28
2.6.1 Auxiliary Lemmata	32
2.7 Bernstein's Inequality for sub-Exponential random variables	32
2.7.1 Auxiliary Lemmata	34

*Vershynin [2019]

1 Preliminaries

1.1 Example on approximate Caratheodory's Theorem

First, we begin by discussing Caratheodory's Theorem:

Theorem 1.1.1 (Caratheodory's Theorem). *Consider a convex set $S \subseteq \mathbb{R}^p$. Any point $x \in S$ can be represented as a convex combination of at most $p + 1$ distinct points from S .*

Remark. This result is a popular result in convex analysis, and is tight. The tight lower bound is achieved by a simplex in p dimensions, which corresponds to $p + 1$ vertices.

Now, we seek an approximation of the above theorem like so: given k points $\{x_i\}_{i=1}^k \subset S$, is it possible to approximate a point $x \in S$? We answer this in the affirmative below:

Theorem 1.1.2 (Approx. Caratheodory's Theorem). *Given $x \in S \subseteq \mathbb{R}^p$, where S is convex, there exists a set of k points $\{x_i\}_{i=1}^k \in S$, such that the following holds:*

$$\left\| x - \frac{1}{k} \sum_{i=1}^k x_i \right\|_2 \leq \frac{\text{diam}(S)}{\sqrt{k}}$$

where $\text{diam}(S) = \sup_{s, t \in S} \|s - t\|_2$.

Proof. By the fact that $x \in S$, we know that we can write x as a convex combination of a subset $\{z_i\}_{i=1}^m$ that satisfy $\text{CONV}(\{z_i\}_{i=1}^m) = S$, where $m \leq p + 1$. Let the coefficients be $\{\lambda_i\}_{i=1}^m$ where $\sum_{i=1}^m \lambda_i = 1$ and $\lambda_i \geq 0$ for all $i \in [m]$.

Consider a random variable Z that takes m different values from the set $\{z_i\}_{i=1}^m$ with probability λ_i . Note that $\mathbb{E}[Z] = x$, since $\mathbb{E}[Z] = \sum_{i=1}^m \Pr(Z = z_i) z_i = \sum_{i=1}^m \lambda_i z_i = x$.

We know that for any $x \in \mathbb{R}^p$ and independent random variables $\{Z_i\}_{i=1}^k$ that satisfy $\mathbb{E}[Z_i] = x$ for all $i \in [k]$:

$$\begin{aligned} \mathbb{E} \left[\left\| x - \frac{1}{k} \sum_{i=1}^k Z_i \right\|_2^2 \right] &= \mathbb{E} \left[\left\| \frac{1}{k} \sum_{i=1}^k (x - Z_i) \right\|_2^2 \right] \\ &= \frac{1}{k^2} \mathbb{E} \left[\left\| \sum_{i=1}^k (x - Z_i) \right\|_2^2 \right] \\ &\stackrel{(i)}{=} \frac{1}{k^2} \sum_{i=1}^k \mathbb{E} [\|x - Z_i\|_2^2] \\ &\stackrel{(ii)}{\leq} \frac{1}{k^2} \sum_{i=1}^k \text{diam}(S)^2 = \frac{\text{diam}(S)^2}{k} \end{aligned}$$

Step (i) holds true due to Lemma 1.1.1. Step (ii) follows from the fact that $Z_i, x \in S$ which implies that $\|Z_i - x\|_2 \leq \text{diam}(S)$ followed by the fact that $\mathbb{E}[c] = c$ for constant c .

Therefore, there exists a realization of $\{Z_i\}_{i=1}^k$, that satisfies:

$$\left\| x - \frac{1}{k} \sum_{i=1}^k Z_i \right\|_2 \leq \frac{\text{diam}(S)}{\sqrt{k}}$$

□

Remark. First note the dimension independence in the result. Secondly, in the special case where S consists of elements with bounded norms i.e., $\|x\|_2 \leq B$ for all $x \in S$, the diameter of the set is bounded by $2B$ by an application of the triangle inequality. Finally, note that if we have $k \rightarrow \infty$ samples from the set, then our approximation is going to be perfect.

The method used to prove Theorem 1.1.1 is called Maurey's Empirical Method.

1.1.1 Auxiliary Lemmata

Lemma 1.1.1. *Let $\{X_i\}_{i=1}^k$ be a set of independent zero-mean random variables. The following holds true:*

$$\mathbb{E} \left[\left\| \sum_{i=1}^k X_i \right\|_2^2 \right] = \sum_{i=1}^k \mathbb{E} [\|X_i\|_2^2]$$

Proof. First note that:

$$\begin{aligned} \left\| \sum_{i=1}^k X_i \right\|_2^2 &= \left\langle \sum_{i=1}^k X_i, \sum_{j=1}^k X_j \right\rangle \\ &= \sum_{i=1}^k \sum_{j=1}^k X_i^T X_j \\ &= \sum_{i=1}^k \|X_i\|_2^2 + 2 \sum_{\substack{i,j=1 \\ i \neq j}}^k X_i^T X_j \end{aligned}$$

Taking expectations on both sides:

$$\begin{aligned} \mathbb{E} \left[\left\| \sum_{i=1}^k X_i \right\|_2^2 \right] &= \mathbb{E} \left[\sum_{i=1}^k \|X_i\|_2^2 \right] + 2 \mathbb{E} \left[\sum_{\substack{i,j=1 \\ i \neq j}}^k X_i^T X_j \right] \\ &= \sum_{i=1}^k \mathbb{E} [\|X_i\|_2^2] + 2 \sum_{\substack{i,j=1 \\ i \neq j}}^k \mathbb{E} [X_i^T X_j] \end{aligned}$$

Since X_i s are independent, $\mathbb{E} [X_i^T X_j] = \mathbb{E} [X_i]^T \mathbb{E} [X_j] = 0$, and this completes the proof. \square

Lemma 1.1.2. *For all integers $m \in [1, n]$, we have the following series of inequalities:*

$$\left(\frac{n}{m} \right)^m \leq \binom{n}{m} \leq \sum_{k=0}^m \binom{n}{k} \leq \left(\frac{en}{m} \right)^m$$

Proof. First inequality:

$$\binom{n}{m} m^m = \frac{n!}{(n-m)! \cdot m!} m^m \geq \frac{n!}{(n-m)!} \geq n^m \Rightarrow \binom{n}{m} \geq \left(\frac{n}{m} \right)^m$$

Second inequality:

$$\binom{n}{m} \leq \binom{n}{m} + \sum_{k=0}^{m-1} \binom{n}{k} = \sum_{k=0}^m \binom{n}{k}$$

Third inequality:

$$\left(\frac{m}{n}\right)^m \sum_{k=0}^m \binom{n}{k} \leq \sum_{k=0}^m \binom{n}{k} \left(\frac{m}{n}\right)^k \leq \sum_{k=0}^n \binom{n}{k} \left(\frac{m}{n}\right)^k = \left(1 + \frac{m}{n}\right)^n \leq e^m \Rightarrow \sum_{k=0}^m \binom{n}{k} \leq \left(\frac{en}{m}\right)^m$$

□

1.2 Quantities and Inequalities associated with RVs

- Expectation: $\mathbb{E}[X]$
- Variance: $\text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2]$
- MGF: $M_X(t) = \mathbb{E}[e^{tX}]$, $t \in \mathbb{R}$
- p^{th} moment: $\mathbb{E}[X^p]$ and p^{th} absolute moment: $\mathbb{E}[|X|^p]$
- L^p norm: $\|X\|_{L^p} = \sqrt[p]{\mathbb{E}[|X|^p]}$
- L^∞ norm: $\|X\|_{L^\infty} = \text{ess sup } |X|$, where $\text{ess sup } |X|$ denotes the supremum over all set with measure not 0. Also note that: $\text{ess sup } |X| \leq \sup |X|$.
- Covariance: $\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$
- CDF: $F_X(t) = \Pr(X \leq t)$, $t \in \mathbb{R}$

For a convex function f and any random variable X , we have by *Jensen's inequality* that:

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)]$$

Consequently, for a concave function f and any random variable X , we have:

$$f(\mathbb{E}[X]) \geq \mathbb{E}[f(X)]$$

As a special case, consider $f(x) : x^{q/p}$ where $q > p$. Note that f is convex. Therefore:

$$(\mathbb{E}[|X|^p])^{q/p} \leq \mathbb{E}[|X|^q] \Rightarrow \|X\|_{L^p} \leq \|X\|_{L^q}$$

Another inequality is the *Cauchy-Schwarz inequality*, which states that for any two RVs X and Y :

$$\mathbb{E}[|XY|] \leq \sqrt{\mathbb{E}[X^2]} \sqrt{\mathbb{E}[Y^2]} = \|X\|_{L^2} \|Y\|_{L^2}$$

We also have *Holder's inequality* which generalizes *Cauchy-Schwarz* to dual norms as:

$$\mathbb{E}[|XY|] \leq \|X\|_{L^p} \|Y\|_{L^q} \quad ; \quad \frac{1}{p} + \frac{1}{q} = 1$$

The following lemma characterizes the expectation as a quantity involving only tails:

Lemma 1.2.1. *Consider a non-negative random variable X . The expectation of this random variable can be written as:*

$$\mathbb{E}[X] = \int_0^\infty \Pr(X > t) dt$$

Proof. For any $x \geq 0$, we have that:

$$x = \int_0^\infty \mathbf{1}_{\{t < x\}} dt = \int_0^x 1 dt + \int_x^\infty 0 dt$$

Therefore:

$$\begin{aligned} X = \int_0^\infty \mathbf{1}_{\{t < X\}} dt &\Rightarrow \mathbb{E}[X] = \mathbb{E} \left[\int_0^\infty \mathbf{1}_{\{t < X\}} dt \right] \\ &= \int_0^\infty \int_{-\infty}^\infty \mathbf{1}_{\{t < x\}} \Pr(X = x) dx dt \\ &= \int_0^\infty \int_t^\infty \Pr(X = x) dx dt \\ &= \int_0^\infty \Pr(X > t) dt \end{aligned}$$

□

A simple generalization for real-valued random variables from the proof of Lemma 1.2.1 is as follows:

Corollary 1.2.1. *Consider a real valued random variable X . The expectation of this random variable can be written as:*

$$\mathbb{E}[X] = \int_0^\infty \Pr(X > t) dt - \int_{-\infty}^0 \Pr(X < t) dt$$

An application of Lemma 1.2.1 is to use it to bound the p^{th} absolute moments via tails:

Corollary 1.2.2. *For any random variable X :*

$$\mathbb{E}[|X|^p] = \int_0^\infty p t^{p-1} \Pr(|X| > t) dt$$

Classical inequalities: Markov and Chebyshev's:

Lemma 1.2.2 (Markov's Inequality). *Consider a non-negative random variable X . Then the tails of X can be bounded as:*

$$\Pr(X > t) \leq \frac{\mathbb{E}[X]}{t}$$

Proof. Note that:

$$\mathbb{E}[X] = \mathbb{E}[X \cdot \mathbf{1}_{\{X > t\}}] + \mathbb{E}[X \cdot \mathbf{1}_{\{X \leq t\}}] \geq \mathbb{E}[X \cdot \mathbf{1}_{\{X > t\}}] \geq t \mathbb{E}[\mathbf{1}_{\{X > t\}}] = t \Pr(X > t) \Rightarrow \Pr(X > t) \leq \frac{\mathbb{E}[X]}{t}$$

□

Corollary 1.2.3 (Chebyshev's Inequality). *Consider a random variable X . Then the probability of deviation from the expectation of X can be bounded as:*

$$\Pr(|X - \mathbb{E}[X]| > t) \leq \frac{\text{Var}(X)}{t^2}$$

Proof. Take $Y = |X - \mathbb{E}[X]|$ as the random variable as apply Markov's inequality (Lemma 1.2.2):

$$\Pr(Y > t) = \Pr(Y^2 > t^2) \leq \frac{\mathbb{E}[(X - \mathbb{E}[X])^2]}{t^2}$$

□

Remark. Note that one can achieve better dependence on t by using higher moments - provided they exist:

$$\Pr(Y > t) = \Pr(Y^{2k} > t^{2k}) \leq \frac{\mathbb{E}[(X - \mathbb{E}[X])^{2k}]}{t^{2k}}$$

1.3 Basic Limit Theorems

Theorem 1.3.1 (Strong Law of Large Numbers). *Let $\{X_i\}_{i=1}^n$ be a sequence of identically and independently distributed random variables with mean μ . The quantity $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ satisfies:*

$$\bar{X}_n \xrightarrow{a.s.} \mu$$

as $n \rightarrow \infty$.

Here $\xrightarrow{a.s.}$ denotes *almost sure convergence*, which is:

$$\Pr \left(\lim_{n \rightarrow \infty} \bar{X}_n = \mu \right) = 1$$

There is a *weak law of large numbers*, which can be derived from Chebyshev's Inequality, for distributions with bounded variance. It is stated below:

Corollary 1.3.1 (Weak Law of Large Numbers). *Let $\{X_i\}_{i=1}^n$ be a sequence of identically and independently distributed random variables with mean μ and variance $\sigma^2 < \infty$. The quantity $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ satisfies:*

$$\bar{X}_n \xrightarrow{p} \mu$$

where \xrightarrow{p} denotes convergence in probability, which is;

$$\forall \epsilon > 0, \quad \lim_{n \rightarrow \infty} \Pr(|\bar{X}_n - \mu| > \epsilon) = 0$$

Proof. First note $\mathbb{E}[\bar{X}_n] = \mu$, and hence $\text{Var}(\bar{X}_n) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{\sigma^2}{n}$.

By Chebyshev's inequality, for any $\epsilon > 0$:

$$\Pr(|\bar{X}_n - \mu| > \epsilon) \leq \frac{\sigma^2}{n\epsilon} \Rightarrow \lim_{n \rightarrow \infty} \Pr(|\bar{X}_n - \mu| > \epsilon) = 0 \quad (\because \text{Sandwich theorem})$$

□

Remark. This weak result is *weak* because $\xrightarrow{a.s.}$ implies \xrightarrow{p} .

Next, we state a result that gives the asymptotic distribution of \bar{X}_n .

Theorem 1.3.2 (Central Limit Theorem). *Let $\{X_i\}_{i=1}^n$ be a sequence of identically and independently distributed random variables with mean μ and variance $\sigma^2 < \infty$. Define $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. Then:*

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{d} \mathcal{N}(0, 1) \quad \text{as } n \rightarrow \infty$$

While this result states that the deviation between the sample mean and population mean is 0 in the limit, we can give some non asymptotic guarantees on the deviation as follows:

Lemma 1.3.1. *Let $\{X_i\}_{i=1}^n$ be a sequence of identically and independently distributed random variables with mean μ and variance $\sigma^2 < \infty$. We have that:*

$$\mathbb{E} \left[\left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| \right] = O \left(\frac{1}{\sqrt{n}} \right)$$

Proof. By Jensen's inequality:

$$\mathbb{E}[|Z|] \leq \sqrt{\mathbb{E}[Z^2]}$$

(Note that this also follows from the fact that $\|Z\|_{L^1} \leq \|Z\|_{L^2}$)

Therefore:

$$\begin{aligned} \mathbb{E}\left[\left|\frac{1}{n}\sum_{i=1}^n X_i - \mu\right|\right] &\leq \sqrt{\mathbb{E}\left[\left(\frac{1}{n}\sum_{i=1}^n X_i - \mu\right)^2\right]} \\ &= \sqrt{\mathbb{E}\left[\left(\frac{1}{n}\sum_{i=1}^n (X_i - \mu)\right)^2\right]} \\ &= \sqrt{\frac{1}{n^2}\mathbb{E}\left[\left(\sum_{i=1}^n (X_i - \mu)\right)^2\right]} \\ &\stackrel{(i)}{=} \sqrt{\frac{1}{n^2}\sum_{i=1}^n \mathbb{E}[(X_i - \mu)^2]} \\ &= \sqrt{\frac{\sigma^2}{n}} = O\left(\frac{1}{\sqrt{n}}\right) \end{aligned}$$

where Step (i) follows from Lemma 1.1.1 for 1-D random variables. □

A special case of the Central Limit Theorem is to provide approximate distributions for binomial distributions. Recall that the binomial distribution $\text{Bin}(n, p)$ is the sum of n independent Bernoulli distribution with parameter p . Therefore, we get that:

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} = \frac{n\bar{X}_n - n\mu}{\sigma\sqrt{n}} = \frac{B_{n,p} - np}{\sqrt{n}\sqrt{p(1-p)}} \xrightarrow{d} \mathcal{N}(0, 1) \text{ as } n \rightarrow \infty$$

where $X_i \sim \text{Ber}(p)$, $i \in [n]$ and $B_{n,p} \sim \text{Bin}(n, p)$. This means that $B_{n,p} \xrightarrow{d} \mathcal{N}(np, np(1-p))$ as $n \rightarrow \infty$.

However, there is a better limit theorem in the regime where $p \rightarrow 0$, $n \rightarrow \infty$ and $np = \lambda > 0$. This is the Poisson Limit Theorem:

Theorem 1.3.3 (Poisson Limit Theorem). *Consider $\{X_i\}_{i=1}^n$ to be n independent Bernoulli variables with parameters p_i . Then, for $n \rightarrow \infty$, $\max_{i \in [n]} p_i \rightarrow 0$ and $\sum_{i=1}^n p_i = \lambda > 0$, we have that:*

$$\sum_{i=1}^n X_i \xrightarrow{d} \text{Poi}(\lambda)$$

Remark. In the special case when all p_i s are equal, we obtain the same result with $n \rightarrow \infty$, $p \rightarrow 0$ and $np = \lambda > 0$ as described informally earlier.

2 Concentration inequalities

2.1 Basic Gaussian Inequalities

Lemma 2.1.1 (Mill's inequalities). *Let $g \sim \mathcal{N}(0, 1)$. We have the following lower and upper bounds for the tail $\Pr(g > t)$, $t > 0$ as follows:*

$$\left(\frac{1}{t} - \frac{1}{t^3}\right) \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \leq \Pr(g > t) \leq \frac{1}{t} \cdot \frac{1}{\sqrt{2\pi}} e^{-t^2/2}$$

Proof. First, the upper bound:

$$\begin{aligned} \Pr(g > t) &= \frac{1}{\sqrt{2\pi}} \int_t^\infty e^{-x^2/2} dx \\ &= \frac{1}{t} \cdot \frac{1}{\sqrt{2\pi}} \int_t^\infty t e^{-x^2/2} dx \\ &\leq \frac{1}{t} \cdot \frac{1}{\sqrt{2\pi}} \int_t^\infty x e^{-x^2/2} dx \\ &= \frac{1}{t} \cdot \frac{1}{\sqrt{2\pi}} \int_{t^2/2}^\infty y e^{-y} dy \quad \left(\because y = \frac{x^2}{2}\right) \\ &= \frac{1}{t} \cdot \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \end{aligned}$$

Second, the lower bound:

$$\begin{aligned} \Pr(g > t) &= \frac{1}{\sqrt{2\pi}} \int_t^\infty e^{-x^2/2} dx \\ &\geq \frac{1}{\sqrt{2\pi}} \int_t^\infty \left(1 - \frac{3}{x^4}\right) e^{-x^2/2} dx \quad \left(\because 1 - \frac{3}{x^4} \leq 1 \ \forall \ x > 0\right) \\ &\geq \left(\frac{1}{t} - \frac{1}{t^3}\right) \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \end{aligned}$$

An alternative proof for the lower bound can be obtained as follows. First note that for $\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$, we have (see [Wainwright, 2019]):

$$\phi'(z) = \frac{1}{\sqrt{2\pi}} \cdot -ze^{-z^2/2} = -z\phi(z)$$

Therefore:

$$\begin{aligned} \Pr(g > t) &= \int_t^\infty \phi(x) dx \\ &= \int_t^\infty -\frac{\phi'(x)}{x} dx \\ &= -\frac{1}{x} \phi(x) \Big|_t^\infty - \int_t^\infty \frac{1}{x^2} \phi(x) dx \quad \left(\because \text{int. by parts with } f(x) = -\frac{1}{x}, g(x) = \phi'(x)\right) \\ &= \frac{\phi(t)}{t} + \int_t^\infty \frac{1}{x^3} \phi'(x) dx \\ &= \frac{\phi(t)}{t} + \frac{\phi(x)}{x^3} \Big|_t^\infty + \underbrace{\int_t^\infty \frac{3}{x^4} \phi(x) dx}_{\geq 0} \quad \left(\because \text{int. by parts with } f(x) = -\frac{1}{x^3}, g(x) = \phi'(x)\right) \\ &\geq \frac{\phi(t)}{t} - \frac{\phi(t)}{t^3} \end{aligned}$$

□

Remark. Note that one can get tail bounds for $\mathcal{N}(0, \sigma^2)$ by simply reparameterising the integrals as:

$$\left(\frac{\sigma}{t} - \frac{\sigma}{t^3}\right) \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \leq \Pr(g > t) \leq \frac{\sigma}{t} \cdot \frac{1}{\sqrt{2\pi}} e^{-t^2/2}$$

The Central Limit Theorem (Theorem 1.3.2) states that averages tend in distribution to a Gaussian. But what can be said about the distribution function itself? The following theorem gives this result:

Theorem 2.1.1 (Berry-Esseen CLT). *Let $\{X_i\}_{i=1}^n$ be a sequence of identically and independently distributed random variables with mean μ and variance $\sigma^2 < \infty$. Define $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ and $\bar{Z}_n = \sqrt{n} \frac{\bar{X}_n - \mu}{\sigma}$.*

Then:

$$|\Pr(\bar{Z}_n > t) - \Pr(g > t)| \leq \frac{\rho}{\sqrt{n}}$$

where $\rho = \frac{\mathbb{E}[|X_i - \mu|^3]}{\sigma^3}$, $i \in [n]$ and $g \sim \mathcal{N}(0, 1)$.

Remark. This theorem basically states that the error of approximation scales as $O\left(\frac{1}{\sqrt{n}}\right)$, which is bad, since we can't always leverage the normal approximation from the central limit theorem always.

2.1.1 Auxiliary Lemmata

Lemma 2.1.2. *Let $g \sim \mathcal{N}(0, 1)$. For $t \geq 1$, we have that:*

$$\mathbb{E}[g^2 \mathbf{1}_{\{g > t\}}] = \frac{t}{\sqrt{2\pi}} e^{-t^2/2} + \Pr(g > t) \leq \left(t + \frac{1}{t}\right) \frac{1}{\sqrt{2\pi}} e^{-t^2/2}$$

Proof.

$$\begin{aligned} \mathbb{E}[g^2 \mathbf{1}_{\{g > t\}}] &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x^2 \mathbf{1}_{\{x > t\}} e^{-x^2/2} dx \\ &= \frac{1}{\sqrt{2\pi}} \int_t^{\infty} x^2 e^{-x^2/2} dx \\ &= \frac{1}{\sqrt{2\pi}} \left((x \cdot e^{-x}) \Big|_t^{\infty} + \int_t^{\infty} e^{-x^2/2} dx \right) \quad \left(\cdot \text{ int. by parts with } f(x) = x, g(x) = x e^{-x^2/2} \right) \\ &= \frac{t}{\sqrt{2\pi}} e^{-t^2/2} + \Pr(g > t) \\ &\leq \frac{t}{\sqrt{2\pi}} e^{-t^2/2} + \frac{1}{t} \cdot \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \\ &= \left(t + \frac{1}{t}\right) \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \end{aligned}$$

□

2.2 Hoeffding's Inequality

First, we describe a simple bounded random variable namely the Rademacher random variable. X with support $\{-1, +1\}$ is said to be distributed w.r.t. a Rademacher distribution if:

$$\Pr(X = +1) = \Pr(X = -1) = \frac{1}{2}$$

Let's discuss about the basic quantities attributed to a random variable. Note that $\mathbb{E}[X] = 0$ and $\text{Var}[X] = \frac{1}{2}$. Also the moment generating function is:

$$M_X(t) = \mathbb{E}[e^{tX}] = \frac{1}{2}e^t + \frac{1}{2}e^{-t} = \cosh(t)$$

We have the first concentration inequality for these random variables.

Theorem 2.2.1. *Let $\{X_i\}_{i=1}^n$ be a sequence of n i.i.d. Rademacher random variables. Then, for any $a \in \mathbb{R}^n$ we have:*

$$\Pr\left(\sum_{i=1}^n a_i X_i > t\right) \leq \exp\left(-\frac{t^2}{2\|a\|_2^2}\right)$$

Proof. Let's begin with Markov's inequality (Lemma 1.2.2):

$$\begin{aligned} \Pr\left(\sum_{i=1}^n a_i X_i > t\right) &= \Pr\left(\exp\left(s \cdot \sum_{i=1}^n a_i X_i\right) > \exp(s \cdot t)\right) \quad (s > 0) \\ &\stackrel{(i)}{\leq} \frac{\mathbb{E}\left[\exp\left(s \cdot \sum_{i=1}^n a_i X_i\right)\right]}{\exp(st)} \\ &= \frac{1}{\exp(st)} \mathbb{E}\left[\prod_{i=1}^n \exp(s \cdot a_i X_i)\right] \\ &\stackrel{(ii)}{=} \frac{1}{\exp(st)} \prod_{i=1}^n \mathbb{E}[\exp(s \cdot a_i X_i)] \\ &\stackrel{(iii)}{=} \frac{1}{\exp(st)} \prod_{i=1}^n \cosh(s \cdot a_i) \\ &\stackrel{(iv)}{\leq} \frac{1}{\exp(st)} \prod_{i=1}^n \exp\left(\frac{s^2 a_i^2}{2}\right) \\ &= \frac{1}{\exp(st)} \exp\left(\frac{s^2 \|a\|_2^2}{2}\right) \end{aligned}$$

Step (i) uses Markov's inequality over $Y = \exp\left(s \cdot \sum_{i=1}^n a_i X_i\right)$. Step (ii) uses the fact that X_i s are independent, and Step (iii) uses the MGF derived earlier. Step (iv) can be attributed to Lemma 2.2.1.

Therefore, we have shown for any $s > 0$ that:

$$\Pr\left(\sum_{i=1}^n a_i X_i > t\right) \leq \exp\left(s^2 \frac{\|a\|_2^2}{2} - st\right)$$

To obtain the tightest / smallest upper bound, we minimize the RHS. The minimizer of the quadratic: $s^2 \frac{\|a\|_2^2}{2} - st$ w.r.t s is $\frac{t}{\|a\|_2^2}$, resulting in:

$$\Pr\left(\sum_{i=1}^n a_i X_i > t\right) \leq \exp\left(-\frac{t^2}{2\|a\|_2^2}\right)$$

□

Remark. Note that the first step generally holds for any invertible function f . We used such a trick in proving Chebyshev's inequality. Also note that since $\mathbb{E}[X_i] = 0$ for all $i \in [n]$, we didn't have an explicit mean term.

Also, in the case where $a_i = \frac{1}{n}$ i.e., we care about the concentration of the sample average to the mean, we get:

$$\Pr(\bar{X}_n > t) \leq \exp\left(-\frac{nt^2}{2}\right)$$

Generally, high probability guarantees are given by setting the RHS to a confidence parameter δ , like so:

$$\bar{X}_n \leq \sqrt{\frac{2}{n} \log\left(\frac{1}{\delta}\right)} \quad \text{w. p.} \geq 1 - \delta$$

Fix the error to ϵ . Smaller the δ , the more confident the estimate, but requires more samples.

There is a two sided version to this inequality as well, and the stated below:

Corollary 2.2.1. *Let $\{X_i\}_{i=1}^n$ be a sequence of n i.i.d. Rademacher random variables. Then, for any $a \in \mathbb{R}^n$ we have:*

$$\Pr\left(\left|\sum_{i=1}^n a_i X_i\right| > t\right) \leq 2 \exp\left(-\frac{t^2}{2\|a\|_2^2}\right)$$

Proof. Simple probability states that if for two events A and B such that $A \Rightarrow B$, then $\Pr(A) \leq \Pr(B)$. Therefore:

$$\begin{aligned} \left|\sum_{i=1}^n a_i X_i\right| > t &\Rightarrow \sum_{i=1}^n a_i X_i > t \vee \sum_{i=1}^n a_i X_i < -t \\ \Rightarrow \Pr\left(\left|\sum_{i=1}^n a_i X_i\right| > t\right) &\leq \Pr\left(\sum_{i=1}^n a_i X_i > t \vee \sum_{i=1}^n a_i X_i < -t\right) \leq \Pr\left(\sum_{i=1}^n a_i X_i > t\right) + \Pr\left(\sum_{i=1}^n a_i X_i < -t\right) \\ &\Rightarrow \Pr\left(\left|\sum_{i=1}^n a_i X_i\right| > t\right) \leq 2 \exp\left(-\frac{t^2}{2\|a\|_2^2}\right) \end{aligned}$$

□

An interesting application of this inequality is as follows: say that you have a set of objects of two kinds, equally probably occurring. If you want to provide high probability guarantees on the number of objects of one kind, then you could use it as follows. For each object, assign $Z_i = 1$ if the object is of the kind you want and 0 otherwise. Now, $S_n = \sum_{i=1}^n Z_i$ is the number of objects of the kind you desire. However, Z_i is supported on $\{0, 1\}$, which means you would have to reparameterize. Note that:

$$\Pr(S_n > t) = \Pr(2S - n > 2t - n) = \Pr\left(\sum_{i=1}^n (2Z_i - 1) > 2t - n\right) \leq \exp\left(-\frac{(2t - n)^2}{2n}\right)$$

So, if you anticipated having at least $\frac{3}{4}n$ objects of the kind you want, then this would happen with probability at most $e^{-n/8}$.

Now the derivation above has given an idea of how to obtain bounds for *bounded random variables* - of which the Rademacher random variable is an example. Below, we have a generalized theorem:

Theorem 2.2.2. Let $\{X_i\}_{i=1}^n$ be a sequence of n independent bounded random variables i.e., every X_i satisfies $X_i \in [l_i, u_i]$. Then, we have that:

$$\Pr \left(\sum_{i=1}^n (X_i - \mathbb{E}[X_i]) > t \right) \leq \exp \left(-\frac{2t^2}{\|u-l\|_2^2} \right)$$

Proof. The standard series of steps used in the proof of Theorem 2.2.1 yields:

$$\begin{aligned} \Pr \left(\sum_{i=1}^n (X_i - \mathbb{E}[X_i]) > t \right) &= \Pr \left(\exp \left(s \cdot \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \right) > \exp(st) \right) \quad (s > 0) \\ &\leq \frac{1}{\exp(st)} \mathbb{E} \left[\exp \left(s \cdot \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \right) \right] \\ &\leq \frac{1}{\exp(st)} \prod_{i=1}^n \mathbb{E} [\exp(s \cdot (X_i - \mathbb{E}[X_i]))] \end{aligned}$$

To bound $\mathbb{E} [\exp(s \cdot (X_i - \mathbb{E}[X_i]))]$: let $Z_i = X_i - \mathbb{E}[X_i]$ for $i \in [n]$. Therefore, we focus on the quantity $\mathbb{E} [\exp(s \cdot Z_i)]$. Adapting the proof from [Wainwright, 2019, Example 2.3], consider Z'_i to be an independent copy of Z_i .

$$\mathbb{E}_{Z_i} [\exp(s \cdot Z_i)] = \mathbb{E}_{Z_i} [\exp(s \cdot (Z_i - \mathbb{E}_{Z'_i}[Z'_i]))] = \mathbb{E}_{Z_i} [\exp(s \cdot (\mathbb{E}_{Z'_i}[Z_i - Z'_i]))] \leq \mathbb{E}_{Z_i, Z'_i} [\exp(s \cdot (Z_i - Z'_i))]$$

Now, we have by symmetricity:

$$\mathbb{E}_{Z_i, Z'_i} [\exp(s \cdot (Z_i - Z'_i))] = \mathbb{E}_{\epsilon \sim \text{Rad}} [\mathbb{E}_{Z_i, Z'_i} [\exp(s\epsilon(Z_i - Z'_i))]] = \mathbb{E}_{Z_i, Z'_i} [\mathbb{E}_{\epsilon \sim \text{Rad}} [\exp(\epsilon \cdot s(Z_i - Z'_i))]]$$

We showed via Lemma 2.2.1 that:

$$\mathbb{E}_{\epsilon \sim \text{Rad}} [\exp(\epsilon \cdot s(Z_i - Z'_i))] \leq \exp \left(\frac{s^2(Z_i - Z'_i)^2}{2} \right) \leq \exp \left(\frac{s^2(u_i - l_i)^2}{2} \right)$$

and therefore:

$$\mathbb{E}_{Z_i} [\exp(s \cdot Z_i)] \leq \exp \left(\frac{s^2(u_i - l_i)^2}{2} \right)$$

Therefore, we obtain:

$$\Pr \left(\sum_{i=1}^n (X_i - \mathbb{E}[X_i]) > t \right) \leq \frac{1}{\exp(st)} \exp \left(\frac{s^2 \|u-l\|_2^2}{2} \right)$$

From the proof of Theorem 2.2.1, with $a = u - l$, we get:

$$\Pr \left(\sum_{i=1}^n (X_i - \mathbb{E}[X_i]) > t \right) \leq \exp \left(-\frac{t^2}{2\|u-l\|_2^2} \right)$$

Note the sub-optimal constant. Hoeffding showed that one can obtain a better bound on the moments as:

$$\mathbb{E}_{Z_i} [\exp(s \cdot Z_i)] \leq \exp \left(\frac{s^2(u_i - l_i)^2}{8} \right)$$

which leads to:

$$\Pr \left(\sum_{i=1}^n (X_i - \mathbb{E}[X_i]) > t \right) \leq \exp \left(-\frac{2t^2}{\|u-l\|_2^2} \right)$$

□

Remark. The introduction and use of the Rademacher random variable is known as the *symmetrization trick*. The result of Hoeffding that gives the improved bound on the MGF is classically known as *Hoeffding's Lemma*. Setting $n = 1$, results in:

$$\Pr(|X_1 - \mathbb{E}[X_1]| > t) \leq 2 \exp\left(-\frac{2t^2}{(u-l)^2}\right)$$

The basic idea is to obtain bounds on the MGF. We will see later about classes of distributions that have a bounded MGF, and therefore allow for the application of this theorem despite being unbounded.

With these results, let's discuss some guarantees for robust mean estimation. Given n samples from a distribution with bounded variance σ^2 and mean μ , the mean \bar{X}_n satisfies:

$$\Pr(|\bar{X}_n - \mu| > \epsilon) \leq \frac{\text{Var}(\bar{X}_n)}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2}$$

Therefore, we require $n = O\left(\frac{\sigma^2}{\epsilon^2}\right)$ samples to obtain an estimate with high probability, say at least 0.75.

Now, given k sets of n samples, can we do any better w.r.t. dependence on success probability? This is tantamount to being given $\tilde{n} = n \cdot k$ samples, and splitting them into k subsets of size n each. The answer is yes, and can be obtained by considering the median of these k means.

To see this, consider k Bernoulli random variables $\{Z_i\}_{i=1}^k$, which take value 1 if $|\bar{X}_{i,n} - \mu| > \epsilon$ and 0 otherwise. Since this is a bounded variable, we can apply the result from Theorem 2.2.2 and the fact that $\mathbb{E}[Z_i] = \Pr(Z_i = 1) \leq \frac{1}{4}$ to get:

$$\Pr\left(\sum_{i=1}^k Z_i > \frac{k}{2}\right) \leq \Pr\left(\sum_{i=1}^k (Z_i - \mathbb{E}[Z_i]) > \frac{k}{4}\right) \leq \exp\left(-\frac{k}{8}\right)$$

Therefore, if $k = \log\left(\frac{8}{\delta}\right)$, we get a mean estimate which with probability at least $1 - \delta$, provides an ϵ -error estimate.

2.2.1 Auxiliary Lemmata

Lemma 2.2.1. *For any x , we have that $\cosh(x) \leq \exp\left(\frac{x^2}{2}\right)$.*

Proof. Note that $\cosh(x) = \frac{1}{2}e^x + \frac{1}{2}e^{-x}$

The Maclaurin series for e^x thus gives:

$$\begin{aligned} \cosh(x) &= \frac{1}{2} \sum_{k=0}^{\infty} \frac{x^{2k}}{(2k)!} + \frac{1}{2} \sum_{k=0}^{\infty} \frac{x^{2k+1}}{(2k+1)!} + \frac{1}{2} \sum_{k=0}^{\infty} \frac{x^{2k}}{(2k)!} - \frac{1}{2} \sum_{k=0}^{\infty} \frac{x^{2k+1}}{(2k+1)!} \\ &= \sum_{k=0}^{\infty} \frac{x^{2k}}{(2k)!} \\ &\stackrel{(i)}{\leq} \sum_{k=0}^{\infty} \frac{x^{2k}}{2^k k!} = \sum_{k=0}^{\infty} \frac{(x^2/2)^k}{k!} = \exp\left(\frac{x^2}{2}\right) \end{aligned}$$

where in Step (i) we have used $(2k)! \geq 2^k \cdot k!$. □

Lemma 2.2.2. *Let X is a non-negative random variable whose density function is uniformly bounded by 1. Then:*

$$M_X(-t) \leq \frac{1}{t} \quad t > 0$$

Proof. By definition:

$$M_X(-t) = \mathbb{E}[e^{-tX}] = \int_0^\infty p(x)e^{-tx} dx \leq \int_0^\infty e^{-tx} dx = \frac{1}{t}$$

□

2.3 Chernoff's Inequality

We discussed earlier about the Poisson Limit Theorem 1.3.3. While Hoeffding's inequality gives us good concentration bounds for bounded random variables such as Bernoulli, we can hope for better bounds by leveraging the fact that for small parameters, the distribution of sum of Bernoulli's could be closer to a Poisson as compared to a Gaussian.

Theorem 2.3.1. *Let $\{X_i\}_{i=1}^n$ be a sequence of n independent Bernoulli random variables with parameters $\{p_i\}_{i=1}^n$. Define the sum $S_n = \sum_{i=1}^n X_i$ and mean $\mu = \mathbb{E}[S_n]$. For any $t > \mu$, we have:*

$$\Pr(S_n > t) \leq e^{-\mu} \left(\frac{e\mu}{t} \right)^t$$

Proof. Again, following the usual steps give us:

$$\Pr(S_n > t) \leq \frac{1}{\exp(st)} \prod_{i=1}^n \mathbb{E}[\exp(s \cdot X_i)] \quad (s > 0)$$

Now, since X_i is a Bernoulli random variable, we get:

$$\mathbb{E}[\exp(s \cdot X_i)] = \exp(s)p_i + (1 - p_i) = 1 + (\exp(s) - 1)p_i \leq \exp((\exp(s) - 1)p_i)$$

and hence:

$$\Pr(S_n > t) \leq \frac{1}{\exp(st)} \exp(\exp(s) - 1)\mu$$

To obtain the tightest upper bound, we minimize for $s > 0$. The minimizer is $s = \log\left(\frac{t}{\mu}\right)$, and the minimum is:

$$\Pr(S_n > t) \leq \left(\frac{\mu}{t} \right)^t e^t e^{-\mu}$$

□

We can also obtain a bound on the other side of the tail like so.

Corollary 2.3.1. *Let $\{X_i\}_{i=1}^n$ be a sequence of n independent Bernoulli random variables with parameters $\{p_i\}_{i=1}^n$. Define the sum $S_n = \sum_{i=1}^n X_i$ and mean $\mu = \mathbb{E}[S_n]$. For any $t > \mu$, we have:*

$$\Pr(S_n \leq t) \leq e^{-\mu} \left(\frac{e\mu}{t} \right)^t$$

Proof. We repeat the proof of Theorem 2.3.1 using the fact that $\Pr(S_n \leq t) = \Pr(-S_n \geq -t)$.

$$\Pr(-S_n \geq -t) \leq \frac{1}{\exp(-st)} \prod_{i=1}^n \mathbb{E}[\exp(-s \cdot X_i)] \quad (s > 0)$$

Now, since X_i is a Bernoulli random variable, we get:

$$\begin{aligned}\mathbb{E}[\exp(-s \cdot X_i)] &= \exp(-s)p_i + (1 - p_i) = 1 + (\exp(-s) - 1)p_i \leq \exp((\exp(-s) - 1)p_i) \\ \Rightarrow \Pr(-S_n \geq -t) &\leq \exp(st) \exp((\exp(-s) - 1)\mu)\end{aligned}$$

To obtain the tightest upper bound, we minimize for $s > 0$. The minimizer is $s = \log\left(\frac{\mu}{t}\right)$, and the minimum is:

$$\Pr(S_n \leq t) \leq \left(\frac{\mu}{t}\right)^t e^t e^{-\mu}$$

□

An interesting consequence of the Chernoff bound when $\mu = \sum_{i=1}^n p_i = \lambda$, but when $n \rightarrow \infty$ is that, we retrieve tail bounds for the Poisson random variable:

Corollary 2.3.2. *Consider the setting of Theorem 2.3.1. Then, for $n \rightarrow \infty$ and $\mu = \lambda$, we have:*

$$\Pr(X \geq t) \leq \left(\frac{e\lambda}{t}\right)^t e^{-\lambda}$$

where $X \sim \text{Poi}(\lambda)$.

Proof. Note by the Poisson Limit Theorem 1.3.3, we have that:

$$\lim_{n \rightarrow \infty} \Pr(S_n \geq t) = \Pr(X \geq t)$$

and this completes the proof. □

A popular version of Chernoff's inequality is to give bounds on the event of deviation from the mean μ . We state it as follows:

Corollary 2.3.3. *Let $\{X_i\}_{i=1}^n$ be a sequence of n independent Bernoulli random variables with parameters $\{p_i\}_{i=1}^n$. Define the sum $S_n = \sum_{i=1}^n X_i$ and mean $\mu = \mathbb{E}[S_n]$. For any $\delta \in (0, 1]$, we have:*

$$\Pr(|S_n - \mu| \leq \delta\mu) \leq 2 \exp(-c\mu\delta^2)$$

where $c > 0$ is a universal constant.

Proof. To prove this bound, we will transform it to the form we say in Theorem 2.3.1.

$$|S_n - \mu'| \leq \delta\mu' \Rightarrow S_n - \mu' > \delta\mu' \vee S_n - \mu' < -\delta\mu'$$

Now, we will bound the probability of the event $S_n - \mu' > \delta\mu'$. For any $\mu' \geq \mu$:

$$\begin{aligned}\Pr(S_n - \mu' > \delta\mu') &= \Pr(S_n > \mu'(\delta + 1)) \\ &\leq \exp(-s\mu'(1 + \delta)) \exp((\exp(s) - 1)\mu) \quad (s > 0) \\ &\leq \exp((\exp(s) - 1)\mu - s\mu'(1 + \delta))\end{aligned}$$

Set $s = \log(1 + \delta)$ to get:

$$\Pr(S_n - \mu' > \delta\mu') \leq \exp(\mu\delta - \mu'(1 + \delta)\log(1 + \delta)) \leq \exp(\mu'\delta - \mu'(1 + \delta)\log(1 + \delta)) \leq \exp\left(-\frac{\mu'\delta^2}{3}\right)$$

where the last step uses Lemma 2.3.1.

Therefore, for any $\mu' \geq \mu$, we have:

$$\Pr(S_n > \mu'(1 + \delta)) \leq \exp\left(-\frac{\mu'\delta^2}{3}\right)$$

Now let's bound the probability of the event $S_n - \mu' < -\delta\mu'$. For any $\mu' \leq \mu$:

$$\begin{aligned} \Pr(S_n - \mu' \leq -\delta\mu') &= \Pr(-S_n \geq \mu'(\delta - 1)) \\ &\leq \exp(s\mu'(1 - \delta)) \exp((\exp(-s) - 1)\mu) \quad (s > 0) \\ &\leq \exp((\exp(-s) - 1)\mu + s\mu'(1 - \delta)) \end{aligned}$$

Set $s = -\log(1 - \delta)$ to get:

$$\Pr(S_n - \mu' \leq -\delta\mu') \leq \exp(-\delta\mu - \mu'(1 - \delta)\log(1 - \delta)) \leq \exp(-\delta\mu' - \mu'(1 - \delta)\log(1 - \delta)) \leq \exp\left(-\frac{\mu'\delta^2}{3}\right)$$

Therefore, for any $\mu' \leq \mu$, we have:

$$\Pr(S_n - \mu' \leq -\delta\mu') \leq \exp\left(-\frac{\mu'\delta^2}{3}\right)$$

Now, one can apply a union bound for $\mu' = \mu$, to get:

$$\Pr(|S_n - \mu| > \delta\mu) \leq 2 \exp\left(-\frac{\mu\delta^2}{3}\right)$$

□

2.3.1 Auxiliary Lemmata

Lemma 2.3.1. *For $x \in [0, 1]$, we have the following inequality:*

$$\log(1 + x) \geq \frac{2x}{2 + x}$$

and consequently:

$$x - (1 + x)\log(1 + x) \leq \frac{-x^2}{2 + x} \leq \frac{-x^2}{3}$$

Proof. We know that:

$$\log(1 + x) = \sum_{k=1}^{\infty} \frac{x^k(-1)^{k-1}}{k} \quad \frac{2x}{2 + x} = \sum_{k=1}^{\infty} \frac{x^k(-1)^{k-1}}{2^{k-1}}$$

Therefore:

$$\log(1 + x) - \frac{2x}{2 + x} = \sum_{k=1}^{\infty} x^k(-1)^{k-1} \left(\frac{1}{k} - \frac{2}{2^k} \right)$$

We know that since $x \in [0, 1]$:

$$\frac{x^k}{k} - \frac{x^{k+1}}{k+1} \geq \frac{x^k}{k(k+1)}$$

and

$$\frac{x^k}{2^{k-1}} - \frac{x^{k+1}}{2^k} \geq \frac{x^k}{2^{k+1}}$$

Due to these facts, and using $k(k+1) \leq 2^{k+1}$, we get $\log(1+x) - \frac{2x}{2+x} \geq 0$. As a consequence:

$$x - (1+x)\log(1+x) \leq x \left(\frac{-x}{2+x} \right) = \frac{-x^2}{2+x} \leq -\frac{x^2}{3}$$

□

2.4 Sub-Gaussian Random Variables

Earlier, we saw Hoeffding's inequality applied to Rademacher random variables. In general, can we say anything about the random variables that satisfy:

$$\Pr \left(\left| \sum_{i=1}^n a_i X_i \right| > t \right) \leq 2 \exp \left(-c \frac{t^2}{\|a\|_2^2} \right)$$

In the special case that $n = 1$ and $a_i = 1$, we have:

$$\Pr(|X_i| > t) \leq 2 \exp(-ct^2)$$

Such random variables are called *sub-Gaussian random variables*. From Theorem 2.2.2, we know that bounded random variables are sub-Gaussian with an appropriate constant $c > 0$. Incidentally, Gaussian random variables also fall in this class (hence the name). The following lemma proves this:

Lemma 2.4.1. *Let $X \sim \mathcal{N}(0, 1)$. Then for any $t > 0$:*

$$\Pr(|X| > t) \leq 2 \exp \left(-\frac{t^2}{2} \right)$$

Proof. First, let's consider $\Pr(X > t)$.

$$\Pr(X > t) = \Pr(e^{sX} > e^{st}) \leq \frac{\mathbb{E}[e^{sX}]}{e^{st}} = e^{\frac{s^2}{2} - st}$$

Minimizing the upper bound, gives $s = t$ and consequently:

$$\Pr(X > t) \leq \exp \left(-\frac{t^2}{2} \right)$$

Next, let's consider $\Pr(X < -t)$. Note that $\Pr(X < -t) = \Pr(-X > t) = \Pr(X > t)$ since X is symmetric about 0, and we have bounded this earlier. Combining them we get:

$$\Pr(|X| > t) \leq 2 \Pr(X > t) \leq 2e^{-\frac{t^2}{2}}$$

□

Now, let's also compute moments of a standard normal random variable:

Lemma 2.4.2. *Let $X \sim \mathcal{N}(0, 1)$. Then for any $p \geq 1$, we have:*

$$\|X\|_{L^p} = \sqrt{2} \left(\frac{\Gamma((1+p)/2)}{\Gamma(1/2)} \right)^{1/p}$$

Proof.

$$\mathbb{E}[|X|^p] = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} |x|^p e^{-x^2/2} dx = \frac{2}{\sqrt{2\pi}} \int_0^{\infty} x^p e^{-x^2/2} dx \stackrel{(i)}{=} \frac{\sqrt{2}^p}{\sqrt{\pi}} \int_0^{\infty} t^{\frac{p+1}{2}-1} e^{-t} dt = \sqrt{2}^p \frac{\Gamma((p+1)/2)}{\sqrt{\pi}}$$

where Step (i) uses the change of variable $t = \frac{x^2}{2}$

By the definition of the L^p norm, we have:

$$\|X\|_{L^p} = (\mathbb{E}[|X|^p])^{1/p} = \sqrt{2} \left(\frac{\Gamma((1+p)/2)}{\Gamma(1/2)} \right)^{1/p}$$

Note that:

$$\Gamma((1+p)/2) = \underbrace{\frac{p-1}{2} \frac{p-3}{2} \dots \frac{1}{2}}_{\frac{p}{2} \text{ terms}} \Gamma(1/2) = O\left(\left(\frac{p}{2}\right)^{p/2}\right) \Rightarrow \left(\frac{\Gamma((1+p)/2)}{\Gamma(1/2)}\right)^{1/p} = O(\sqrt{p})$$

which means that $\|X\|_{L^p} = O(\sqrt{p})$. □

Now, let's look at some properties of sub-Gaussian random variables. Recall that a sub-Gaussian random variable X is one that satisfies:

$$\Pr(|X| > t) \leq 2 \exp\left(-\frac{t^2}{K_1^2}\right) \quad (1)$$

for some $K_1 > 0$, and $t \geq 0$. We will see that the following properties are equivalent:

Lemma 2.4.3. *Let X be a sub-Gaussian random variable as defined in Eqn 1. The following statements are equivalent:*

- (i) X is sub-Gaussian
- (ii) $\|X\|_{L^p} \leq K_2 \sqrt{p}$ for all $p \geq 1$
- (iii) $\mathbb{E}[\exp(\lambda^2 X^2)] \leq \exp(K_3^2 \lambda^2)$ for all $|\lambda| \leq \frac{1}{K_3}$
- (iv) $\mathbb{E}\left[\exp\left(\frac{X^2}{K_4^2}\right)\right] \leq 2$
- (v) $\mathbb{E}[\exp(\lambda X)] \leq \exp(K_5^2 \lambda^2)$ for all $\lambda \in \mathbb{R}$ (if X is zero mean sub-Gaussian)

Proof. (i) \Rightarrow (ii) First, consider $\mathbb{E}[|X|^p]$. From Corollary 1.2.2, we have:

$$\begin{aligned} \mathbb{E}[|X|^p] &= \int_0^{\infty} p t^{p-1} \Pr(|X| > t) dt \\ &\leq \int_0^{\infty} 2 p t^{p-1} \exp\left(-\frac{t^2}{K_1^2}\right) dt \\ &= p \int_0^{\infty} K_1^p u^{\frac{p}{2}-1} \exp(-u) du \\ &= p K_1^p \Gamma(p/2) \\ &\leq 3 K_1^p p \left(\frac{p}{2}\right)^{p/2} \end{aligned}$$

This implies: $\|X\|_{L^p} \leq 3 K_1 \sqrt{p}$, and therefore for $K_2 \leq 3 K_1$, we get obtain the proof.

(ii) \Rightarrow (iii) By the series expansion:

$$\begin{aligned}
\mathbb{E} [\exp(\lambda^2 X^2)] &= \mathbb{E} \left[\sum_{k=0}^{\infty} \frac{(\lambda^2 X^2)^k}{k!} \right] \\
&= 1 + \sum_{k=1}^{\infty} \frac{\lambda^{2k} \mathbb{E}[X^{2k}]}{k!} \\
&\leq 1 + \sum_{k=1}^{\infty} \frac{\lambda^{2k} (2K_1^{2k} k \Gamma(k))}{k!} \\
&= 1 + \sum_{k=1}^{\infty} \lambda^{2k} 2K_1^{2k} \\
&\leq \sum_{k=0}^{\infty} 2^{2k} \lambda^{2k} K_1^{2k}
\end{aligned}$$

The above series converges if $2^2 \lambda^2 K_1^2 < 1 \Rightarrow |\lambda| < \frac{1}{2K_1}$. If this condition is satisfied:

$$\mathbb{E} [\exp(\lambda^2 X^2)] \leq \frac{1}{1 - 4\lambda^2 K_1^2} \stackrel{(i)}{\leq} e^{8\lambda^2 K_1^2}$$

where Step (i) is due to Lemma 2.4.6 where $4\lambda^2 K_1^2 \leq \frac{1}{2} \Rightarrow |\lambda| \leq \frac{1}{2\sqrt{2}K_1}$. This yields (iii) with $K_3 = 2\sqrt{2}K_1$.

(iii) \Rightarrow (iv) Choose $\lambda = \frac{1}{\sqrt{2}K_3}$ in (iii) to get:

$$\mathbb{E} \left[\exp \left(\frac{X^2}{2K_3^2} \right) \right] \leq \exp(1/2) \leq 2$$

(iv) \Rightarrow (i)

$$\Pr \left(\frac{|X|}{K_4} > t \right) = \Pr \left(\frac{X^2}{K_4^2} > t^2 \right) = \mathbb{E} \left[\exp \left(\frac{X^2}{K_4^2} \right) \right] \exp(-t^2) \leq 2 \exp(-t^2)$$

With $t' = K_4 t$, we get:

$$\Pr(|X| > t') \leq 2 \exp \left(-\frac{t'^2}{K_4^2} \right)$$

(ii) \Rightarrow (v) From the series expansion:

$$\begin{aligned}
\mathbb{E}[\exp(\lambda X)] &= \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} \mathbb{E}[X^k] \\
&\leq 1 + \sum_{k=2}^{\infty} \frac{\lambda^k}{k!} \mathbb{E}[|X|^k] \\
&\stackrel{(i)}{\leq} 1 + \sum_{k=1}^{\infty} \frac{2^k \lambda^{2k}}{(2k)!} \mathbb{E}[|X|^{2k}] \\
&\leq 1 + \sum_{k=1}^{\infty} \frac{2^k \lambda^{2k}}{(2k)!} K_1^{2k} \underbrace{k \Gamma(k)}_{k!} \\
&\leq 1 + \sum_{k=1}^{\infty} \lambda^{2k} K_1^{2k} \frac{2^k k!}{(2k)!} \\
&\stackrel{(ii)}{\leq} \sum_{k=0}^{\infty} \frac{\lambda^{2k} (K_1)^{2k}}{k!} = \exp(\lambda^2 K_1^2)
\end{aligned}$$

where in Step (i) we have bounded odd moments using even moments (Lemma 2.4.7), and in Step (ii) we have used $\frac{2^k k!}{(2k)!} \leq \frac{1}{k!}$ (which follows from $\left(\frac{2k}{k}\right)^k \leq \binom{2k}{k}$).

(v) \Rightarrow (i)

$$\Pr(X > t) = \Pr(\exp(\lambda X) > \exp(\lambda t)) \leq \frac{\mathbb{E}[\exp(\lambda X)]}{e^{\lambda t}} \leq \exp(\lambda^2 K_1^2 - \lambda t)$$

Minimizing the upper bound over $\lambda > 0$ gives:

$$\Pr(X > t) \leq \exp\left(-\frac{t^2}{2K_1^2}\right)$$

Analogously, we can obtain:

$$\Pr(X < -t) \leq \exp\left(-\frac{t^2}{2K_1^2}\right)$$

Combining these via the union bound gives:

$$\Pr(|X| > t) \leq 2 \exp\left(-\frac{t^2}{2K_1^2}\right)$$

□

Remark. Note that (ii) states that sub-Gaussian random variables have bounded moments. It is also worth noting that when $\mathbb{E}[X] \neq 0$, we can still bound the moment generating function of X as follows:

Corollary 2.4.1. Let X be a non-zero mean sub-Gaussian random variable satisfying

$\Pr(|X| > t) \leq 2 \exp\left(-\frac{t^2}{K_1^2}\right)$. Then we have that:

$$\mathbb{E}[\exp(\lambda X)] \leq 2 \exp(\lambda^2 K_1^2)$$

Proof. The proof mostly follows from the proof of (ii) \Rightarrow (v), except that:

$$\mathbb{E}[\exp(\lambda X)] \leq 1 + \mathbb{E}[\lambda X] + \sum_{k=1}^{\infty} \frac{2^k \lambda^{2k}}{(2k)!} \mathbb{E}[|X|^{2k}] \leq 1 + \frac{1}{2} + \frac{\lambda^2}{2} \mathbb{E}[|X|^2] + \sum_{k=1}^{\infty} \frac{2^k \lambda^{2k}}{(2k)!} \mathbb{E}[|X|^{2k}] \leq 2 \sum_{k=0}^{\infty} \frac{2^k \lambda^{2k}}{(2k)!} \mathbb{E}[|X|^{2k}]$$

We know that $\sum_{k=0}^{\infty} \frac{2^k \lambda^{2k}}{(2k)!} \mathbb{E}[|X|^{2k}] \leq \exp(\lambda^2 K_1^2)$, and this completes the proof. □

Therefore, except for a constant factor, the bound on the MGF is the same for non-centered sub-Gaussian random variables.

Additionally, an interesting observation is that (iii) is equivalent to claiming that $\mathbb{E}[\exp(\lambda^2 X^2)]$ is finite in a bounded neighbourhood of 0, whose radius is $\frac{1}{K_3}$. We can verify this for the standard normal distribution.

$$\begin{aligned} \mathbb{E}[\exp(\lambda^2 X^2)] &= \int_{-\infty}^{\infty} e^{\lambda^2 x^2} p_X(x) dx = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{\frac{x^2}{2}(2\lambda^2 - 1)} dx = \sqrt{\frac{2}{\pi}} \int_0^{\infty} e^{\frac{x^2}{2}(2\lambda^2 - 1)} dx \\ &= \begin{cases} \frac{1}{\sqrt{1-2\lambda^2}} & \text{if } \lambda \leq \frac{1}{\sqrt{2}} \\ \infty & \text{o.w.} \end{cases} \end{aligned}$$

With these equivalences, let's bring some concreteness to the constants used above. We define the *sub-Gaussian norm* of a sub-Gaussian random variable X as:

$$\|X\|_{\psi_2} = \inf \left\{ t > 0 : \mathbb{E} \left[\exp \left(\frac{X^2}{t^2} \right) \right] \leq 2 \right\}$$

Recall that **(iv)** in Lemma 2.4.3, we note that a sub-Gaussian random variable satisfied the condition defining the set with some constant. The sub-Gaussian norm is merely the tightest constant. By terminology, this is a norm, which means:

- For a random variable whose support is 0, the norm $\|\mathbf{0}\|_{\psi_2} = 0$.
- For $\alpha > 0$, $\|\alpha X\|_{\psi_2} = \alpha \|X\|_{\psi_2}$.
- For two sub-Gaussian random variables X and Y , we have $\|X + Y\|_{\psi_2} \leq \|X\|_{\psi_2} + \|Y\|_{\psi_2}$.

Now, since $\|X\|_{\psi_2}$ is the smallest value t such that $\mathbb{E} \left[\exp \left(\frac{X^2}{t^2} \right) \right] \leq 2$ and by the equivalences in Lemma 2.4.3, we can restate the equivalent statements as follows:

- (i) $\Pr(|X| > t) \leq 2 \exp \left(-\frac{C_1 t^2}{\|X\|_{\psi_2}^2} \right)$ for all $t \geq 0$
- (ii) $\|X\|_{L^p} \leq C_2 \|X\|_{\psi_2} \sqrt{p}$ for all $p \geq 1$
- (iii) $\mathbb{E} [\exp(\lambda^2 X^2)] \leq \exp(C_3 \lambda^2 \|X\|_{\psi_2}^2)$ for λ satisfying $|\lambda| \leq \frac{1}{\sqrt{C_3} \|X\|_{\psi_2}}$
- (iv) If $\mathbb{E}[X] = 0$, then $\mathbb{E}[\exp(\lambda X)] \leq \exp(C_4 \lambda^2 \|X\|_{\psi_2}^2)$ for $\lambda \in \mathbb{R}$

where C_1, C_2, C_3, C_4 are universal constants.

Now, let's try computing the *sub-gaussian norms* for common sub-Gaussian distributions.

1. Gaussian distribution. Let $X \sim \mathcal{N}(0, \sigma^2)$. Then:

$$\mathbb{E} \left[\exp \left(\frac{X^2}{t^2} \right) \right] = \begin{cases} \frac{1}{\sigma} \left(\frac{1}{\sigma^2} - \frac{1}{t^2} \right)^{-1} & \text{if } t > \sigma \\ \infty & \text{o.w.} \end{cases}$$

For the LHS to be ≤ 2 , $t \geq \frac{4}{3}\sigma \Rightarrow \|X\|_{\psi_2} = C\sigma$. **Why is it supposed to be $\|X\|_{\psi_2} \leq C\sigma$ instead?**

2. Rademacher distribution. Let $X \sim \text{Rad}$. Then:

$$\mathbb{E} \left[\exp \left(\frac{X^2}{t^2} \right) \right] = e^{\frac{1}{t^2}} \leq 2 \Rightarrow \|X\|_{\psi_2} = \frac{1}{\sqrt{\log(2)}}$$

Remark. For random variables with support $\{-a, a\}$, the sub-Gaussian norm is $\|X\|_{\psi_2} = \frac{a}{\sqrt{\log(2)}}$.

This follows from $\|\cdot\|_{\psi_2}$ being a norm.

An important result concerning the maximum of a sequence of sub-Gaussian random variables is next.

Theorem 2.4.1. *Let $\{X_i\}_{i=1}^\infty$ be a sequence of sub-Gaussian random variables, that are not necessarily independent with sub-Gaussian norms $\{\|X_i\|_{\psi_2}\}_{i=1}^\infty$. Then:*

$$\mathbb{E} \left[\max_i \frac{|X_i|}{\sqrt{1 + \log(i)}} \right] \leq CK$$

where $C > 0$ is a universal constant and $K = \max_i \|X_i\|_{\psi_2}$.

Proof. First note that $\max_i \frac{|X_i|}{\sqrt{1+\log(i)}}$ is a positive random variable. This allows us to write:

$$\mathbb{E} \left[\max_i \frac{|X_i|}{\sqrt{1+\log(i)}} \right] = \int_0^\infty \Pr \left(\max_i \frac{|X_i|}{\sqrt{1+\log(i)}} > t \right) dt$$

Now:

$$\Pr \left(\max_i \frac{|X_i|}{\sqrt{1+\log(i)}} > t \right) \leq \Pr \left(\bigcup_i \left\{ \frac{|X_i|}{\sqrt{1+\log(i)}} > t \right\} \right) \leq \sum_{i=1}^\infty \Pr \left(\frac{|X_i|}{\sqrt{1+\log(i)}} > t \right)$$

Back to the integral - now, we will partition in the integral into two parts, one where $\frac{ct^2}{K^2} \geq 2$, thereby, leading to convergence of the sum, and one where $\frac{ct^2}{K^2} \leq 2$, where we will simply bound the tail probability by 1.

$$\begin{aligned} \mathbb{E} \left[\max_i \frac{|X_i|}{\sqrt{1+\log(i)}} \right] &= \int_0^{\frac{2K}{\sqrt{c}}} \Pr \left(\max_i \frac{|X_i|}{\sqrt{1+\log(i)}} > t \right) dt + \int_{\frac{2K}{\sqrt{c}}}^\infty \sum_{i=1}^\infty \Pr \left(\frac{|X_i|}{\sqrt{1+\log(i)}} > t \right) dt \\ &\leq \frac{2K}{\sqrt{c}} + \int_{\frac{2K}{\sqrt{c}}}^\infty \sum_{i=1}^\infty 2 \exp \left(-\frac{ct^2(1+\log(i))}{\|X_i\|_{\psi_2}^2} \right) dt \\ &\leq \frac{2K}{\sqrt{c}} + \int_{\frac{2K}{\sqrt{c}}}^\infty \sum_{i=1}^\infty 2 \exp \left(-\frac{ct^2(1+\log(i))}{K^2} \right) dt \\ &= \frac{2K}{\sqrt{c}} + \int_{\frac{2K}{\sqrt{c}}}^\infty \sum_{i=1}^\infty 2 \left(\frac{1}{ei} \right)^{\frac{ct^2}{K^2}} dt \\ &\leq \frac{2K}{\sqrt{c}} + \int_{\frac{2K}{\sqrt{c}}}^\infty 2e^{-\frac{ct^2}{K^2}} \sum_{i=1}^\infty \left(\frac{1}{i} \right)^2 dt \\ &\leq \frac{2K}{\sqrt{c}} + \frac{\pi^2}{3} \int_0^\infty e^{-\frac{ct^2}{K^2}} dt \\ &= \frac{2K}{\sqrt{c}} + \frac{\pi^2}{3} \frac{K\sqrt{2\pi}}{\sqrt{c}} = CK \end{aligned}$$

where $C = \frac{2}{\sqrt{c}} \left(1 + \frac{\pi^{5/2}}{3\sqrt{2}} \right)$ □

This is in the asymptotic regime, so in we were given n sub-Gaussian random variables, we have the following corollary:

Corollary 2.4.2. *Let $\{X_i\}_{i=1}^n$ be a sequence of sub-Gaussian random variables, that are not necessarily independent with sub-Gaussian norms $\{\|X_i\|_{\psi_2}\}_{i=1}^n$. Then:*

$$\mathbb{E} \left[\max_i |X_i| \right] \leq CK \sqrt{\log(n)}$$

for $n \geq 2$.

Proof. For $n \geq 2$, we have that:

$$\frac{1}{\sqrt{1+\log(i)}} \geq \frac{1}{\sqrt{2\log(n)}} \quad \forall i \in [n] \Rightarrow \mathbb{E} \left[\max_i \frac{|X_i|}{\sqrt{1+\log(i)}} \right] \geq \mathbb{E} \left[\max_i \frac{|X_i|}{\sqrt{2\log(n)}} \right]$$

In the proof of Theorem 2.4.1, we used the fact that $n \rightarrow \infty$ to obtain the series $\sum_{i=1}^{\infty} \frac{1}{i^2}$. Since the summand is positive, the partial sum can be bounded by the series to get:

$$\mathbb{E} \left[\max_i \frac{|X_i|}{\sqrt{2 \log(n)}} \right] \leq \mathbb{E} \left[\max_i \frac{|X_i|}{\sqrt{1 + \log(i)}} \right] \leq CK \Rightarrow \mathbb{E} \left[\max_i |X_i| \right] \leq CK \sqrt{2 \log(n)}$$

and this completes the proof.

An alternative proof can be given using Jensen's inequality as follows (from [Devroye and Lugosi, 2012, Lemma 2.2]). Note that $\max_{i \in [n]} \max\{X_i, -X_i\} = \max_{i \in [n]} |X_i|$. Therefore, consider a sequence of $2n$

random variables $\{Y_i\}_{i=1}^{2n}$, where $Y_i = \begin{cases} X_i & \text{if } i \bmod 2 = 0 \\ -X_i & \text{o.w.} \end{cases}$.

$$\begin{aligned} \exp \left(\lambda \mathbb{E} \left[\max_i |X_i| \right] \right) &= \exp \left(\lambda \mathbb{E} \left[\max_i Y_i \right] \right) \\ &\leq \mathbb{E} \left[\exp \left(\lambda \max_i Y_i \right) \right] \\ &= \mathbb{E} \left[\max_i \exp(\lambda Y_i) \right] \\ &\leq \sum_{i=1}^{2n} \mathbb{E} [\exp(\lambda Y_i)] \\ &\leq 4n \exp(C \lambda^2 K^2) \\ \Rightarrow \mathbb{E} \left[\max_i |X_i| \right] &\leq \frac{\log(4n)}{\lambda} + C \lambda K^2 \quad (\lambda > 0) \end{aligned}$$

The tightest upper bound can be obtained by $\lambda = \sqrt{\frac{\log(4n)}{C K^2}}$ to get:

$$\mathbb{E} \left[\max_i |X_i| \right] \leq C' K \sqrt{\log(4n)}$$

□

Remark. The alternative proof provides a general strategy to deal with maxima of absolute random variables, and indicates that this yields a multiplicative constant factor in looseness. However, this cannot be used when the independence condition is required, as $|X_i| = \max\{X_i, -X_i\}$ means that the new set of $2n$ random variables are no longer independent.

This is tight upto constant factors, as shown in the lemma below:

Lemma 2.4.4. *Let $\{X_i\}_{i=1}^n$ be a sequence of independent standard normal variables. Then:*

$$\mathbb{E} \left[\max_i X_i \right] \geq c \sqrt{\log(n)}$$

where c is a universal constant.

Proof. **TODO**

□

One other interesting fact to know is about the linear combination of independent sub-Gaussian random variables. As a motivating example, let's look at the linear combination of independent normal variables. We are familiar with the following result:

$$\sum_{i=1}^n a_i X_i \sim \mathcal{N} \left(\sum_{i=1}^n a_i \mu_i, \sum_{i=1}^n a_i^2 \sigma_i^2 \right)$$

This follows from analyzing the MGF of random variable $Y = \sum_{i=1}^n X_i$, and using the independence condition to break to MGF into a product. One can expect such a result for sub-Gaussian random variables as well.

Lemma 2.4.5. *Let $\{X_i\}_{i=1}^n$ be a sequence of n independent sub-Gaussian random variables with mean 0, and sub-Gaussian norms $\{\|X_i\|_{\psi_2}\}_{i=1}^n$. Then, $Y = \sum_{i=1}^n a_i X_i$ is sub-Gaussian with sub-Gaussian norm satisfying:*

$$\|Y\|_{\psi_2} \leq \sqrt{C \sum_{i=1}^n a_i^2 \|X_i\|_{\psi_2}^2} \leq \sqrt{C \|a\|_2^2 K^2}$$

where $K = \max_i \|X_i\|_{\psi_2}$, and C is a universal constant.

Proof. As remarked earlier, we will look at the MGF of Y .

$$\begin{aligned} \mathbb{E}[e^{\lambda Y}] &= \mathbb{E}\left[\exp\left(\sum_{i=1}^n \lambda a_i X_i\right)\right] \\ &= \prod_{i=1}^n \mathbb{E}\left[e^{(\lambda a_i) X_i}\right] \\ &\leq \prod_{i=1}^n e^{C \lambda^2 a_i^2 \|X_i\|_{\psi_2}^2} = \exp\left(C \lambda^2 \sum_{i=1}^n a_i^2 \|X_i\|_{\psi_2}^2\right) \end{aligned}$$

The final step uses:

$$\sum_{i=1}^n a_i^2 \|X_i\|_{\psi_2}^2 \leq \sum_{i=1}^n a_i^2 K^2 = \|a\|_2^2 K^2$$

□

Remark. Now, since we have shown that Y is sub-Gaussian, by manner of equivalences in Lemma 2.4.3, we can derive a concentration inequality for general sub-Gaussian random variables.

2.4.1 Auxiliary Lemmata

Lemma 2.4.6. *For $x \in [0, 1/2]$, we have that $\frac{1}{1-x} \leq e^{2x}$.*

Proof. Both functions are monotonically increasing in $[0, 1/2]$ and satisfy: $f(0) = g(0)$ and $f(1/2) \leq g(1/2)$, with no crossings i.e., $\nexists x \in (0, 1/2]$ such that $f(x) = g(x)$. This means that $f(x) \leq g(x)$ for all $x \in [0, 1/2]$. □

Lemma 2.4.7. *For a random variable X , we have for $k \geq 0$:*

$$\mathbb{E}[|\lambda X|^{2k+1}] \leq \frac{1}{2} \mathbb{E}[|\lambda X|^{2k}] + \frac{1}{2} \mathbb{E}[|\lambda X|^{2k+2}]$$

As a consequence:

$$\mathbb{E}[\exp(\lambda X)] \leq 1 + \mathbb{E}[X] + \sum_{k=1}^{\infty} \frac{2^k \lambda^{2k} \mathbb{E}[|X|^{2k}]}{(2k)!}$$

Proof. By Cauchy-Schwarz and AM-GM inequality, we have:

$$\mathbb{E}[|\lambda X|^{2k+1}] \leq \sqrt{\mathbb{E}[|\lambda X|^{2k}] \mathbb{E}[|\lambda X|^{2k+2}]} \leq \frac{1}{2} \mathbb{E}[|\lambda X|^{2k}] + \frac{1}{2} \mathbb{E}[|\lambda X|^{2k+2}]$$

Now, by the series expansion:

$$\begin{aligned}
\mathbb{E}[\exp(\lambda X)] &\leq 1 + \mathbb{E}[\lambda X] + \sum_{k=2}^{\infty} \frac{\lambda^k \mathbb{E}[X^k]}{k!} \\
&\leq 1 + \mathbb{E}[X] + \sum_{k=2}^{\infty} \frac{\lambda^k \mathbb{E}[|X|^k]}{k!} \\
&\leq 1 + \mathbb{E}[X] + \sum_{k=1}^{\infty} \mathbb{E}[|X|^{2k}] \left(\frac{\lambda^{2k}}{(2k)!} + \frac{1}{2} \left(\frac{\lambda^{2k}}{(2k-1)!} + \frac{\lambda^{2k}}{(2k+1)!} \right) \right) \\
&\leq 1 + \mathbb{E}[X] + \sum_{k=1}^{\infty} \mathbb{E}[|X|^{2k}] \frac{\lambda^{2k} 2^k}{(2k)!}
\end{aligned}$$

□

2.5 Hoeffding's and Khintchine Inequalities for sub-Gaussian random variables, Centering

In Lemma 2.4.5, we saw that a sum of independent sub-Gaussian random variables is sub-Gaussian. Following the remark from earlier, we have the generalized Hoeffding's inequality:

Corollary 2.5.1. *Let $\{X_i\}_{i=1}^n$ be a sequence of independent, zero mean sub-Gaussian random variables with sub-Gaussian norms $\{\|X_i\|_{\psi_2}\}_{i=1}^n$, and $a \in \mathbb{R}^n$ be a vector. Then for $t \geq 0$:*

$$\Pr \left(\left| \sum_{i=1}^n a_i X_i \right| > t \right) \leq 2 \exp \left(- \frac{ct^2}{\sum_{i=1}^n a_i \|X_i\|_{\psi_2}^2} \right) \leq 2 \exp \left(- \frac{ct^2}{\|a\|_2^2 K^2} \right)$$

Remark. Set $a_i = \frac{1}{n}$ for all $i \in [n]$ gives a mean concentration inequality like so:

$$\Pr (|\bar{X}_n| > t) \leq 2 \exp \left(- \frac{cnt^2}{K^2} \right)$$

where \bar{X}_n is the mean of the n random variables $\{X_i\}_{i=1}^n$.

Khintchine inequalities are a collection of inequalities providing upper and lower bounds for L^p norms of linear combination of independent sub-Gaussian random variables. We state them formally in Lemmas 2.5.1 and 2.5.2.

Lemma 2.5.1. *Let $\{X_i\}_{i=1}^n$ be a sequence of independent, zero mean sub-Gaussian random variables with sub-Gaussian norms $\{\|X_i\|_{\psi_2}\}_{i=1}^n$. Then for $a \in \mathbb{R}^n$ and $p \geq 2$, we have that:*

$$\|a\|_2 \leq \left\| \sum_{i=1}^n a_i X_i \right\|_{L^p} \leq C \sqrt{p} K \|a\|_2$$

where C is a universal constant, and $K = \max_i \|X_i\|_{\psi_2}$.

Proof. We know from Corollary 2.5.1 that the linear combination of independent sub-Gaussian random variables are sub-Gaussian with sub-Gaussian norm $\leq K\|a\|_2$. By the equivalences (Lemma 2.4.3), we have that:

$$\left\| \sum_{i=1}^n a_i X_i \right\|_{L^p} \leq C \sqrt{p} K \|a\|_2 \quad \forall \quad p \geq 1$$

We also know that: $\|X\|_{L^2} \leq \|X\|_{L^p}$ for $p \geq 2$, which gives:

$$\left\| \sum_{i=1}^n a_i X_i \right\|_{L^p} \geq \left\| \sum_{i=1}^n a_i X_i \right\|_{L^2} = \sqrt{\mathbb{E} \left[\left| \sum_{i=1}^n a_i X_i \right|^2 \right]} = \sqrt{\sum_{i=1}^n a_i^2} = \|a\|_2$$

□

Lemma 2.5.2. *Let $\{X_i\}_{i=1}^n$ be a sequence of independent, zero mean sub-Gaussian random variables with sub-Gaussian norms $\{\|X_i\|_{\psi_2}\}_{i=1}^n$. Then for $a \in \mathbb{R}^n$, we have that:*

$$c(K)\|a\|_2 \leq \left\| \sum_{i=1}^n a_i X_i \right\|_{L^1} \leq \|a\|_2$$

where $c(K)$ is a constant depending on K where $K = \max_i \|X_i\|_{\psi_2}$.

Proof. The upper bound follows from the fact that $\|X\|_{L^1} \leq \|X\|_{L^2}$ for any random variable. The proof of 2.5.1 derives $\left\| \sum_{i=1}^n a_i X_i \right\|_{L^2}$.

By Cauchy-Schwarz:

$$\mathbb{E}[|X|^k] \leq \sqrt{\mathbb{E}[|X|^{k-1}] \mathbb{E}[|X|^{k+1}]}$$

For $k = 2$, we get: $\mathbb{E}[X^2] \leq \sqrt{\mathbb{E}[|X|]} \sqrt{\mathbb{E}[|X|^3]}$.

From Lemma 2.5.1, we know that for $Y = \sum_{i=1}^n a_i X_i$

- $\mathbb{E}[|Y|^3] \leq CK^3 \|a\|_2^3$
- $\mathbb{E}[|Y|^2] \geq \|a\|_2^2$

Assembling them together, we get:

$$\mathbb{E}[|Y|] \geq \frac{\|a\|_2^4}{CK^3 \|a\|_2^3} = c(K)\|a\|_2$$

□

Remark. For the proof to mean something, it is worth noting that $c(K) \leq 1$.

Lemma 2.5.3. *Let $\{X_i\}_{i=1}^n$ be a sequence of independent, zero mean sub-Gaussian random variables with sub-Gaussian norms $\{\|X_i\|_{\psi_2}\}_{i=1}^n$. Then for $a \in \mathbb{R}^n$ and $p \in (0, 2)$, we have that:*

- If $p \in (1, 2)$, then:

$$\frac{\|a\|_2}{C_1 K} \sqrt{\frac{p}{p-1}} \leq \left\| \sum_{i=1}^n a_i X_i \right\|_{L^p} \leq \|a\|_2$$

- If $p \in (0, 1)$, then:

$$\left(\frac{C_2}{K^5} \right)^{1/p} \|a\|_2 \leq \left\| \sum_{i=1}^n a_i X_i \right\|_{L^p} \leq c(K)\|a\|_2$$

where C_1, C_2 are universal constants, $c(K)$ is a constant depending of K where $K = \max_i \|X_i\|_{\psi_2}$.

Proof. We will break it down into two cases:

$p \in (1, 2)$: By Holder's inequality: if $p \in (1, 2)$, then $q \in (2, \infty)$. We therefore have:

$$\mathbb{E}[Y^2] \leq (\mathbb{E}[|Y|^p])^{\frac{1}{p}} (\mathbb{E}[|Y|^q])^{\frac{1}{q}}$$

We know that $\|Y\|_q \leq CK\sqrt{q}\|a\|_2$ and $\mathbb{E}[Y^2] = \|a\|_2^2$ for $Y = \sum_{i=1}^n a_i X_i$. Therefore:

$$\|Y\|_{L^p} \geq \frac{\|a\|_2}{CK} \sqrt{\frac{p}{p-1}}$$

The upper bound is trivial since $\|Y\|_{L^p} \leq \|Y\|_{L^2} = \|a\|_2$.

$p \in (0, 1)$: By Cauchy Schwarz:

$$\mathbb{E}[|Y|^{p+2}] \leq \sqrt{\mathbb{E}[|Y|^p]} \sqrt{\mathbb{E}[|Y|^{p+4}]}$$

We know that since $p \in (0, 1)$, $p+4 \in (4, 5)$, and hence $\mathbb{E}[|Y|^{p+2}] \leq \|Y\|_{L^5}^{p+2} \leq CK^{p+4}\|a\|_2^{p+4}$. Similarly, $p+2 \in (2, 3)$, and hence $\mathbb{E}[|Y|^{p+2}] \geq \|Y\|_{L^2}^{p+2} = \|a\|_2^{p+2}$.

Therefore:

$$\mathbb{E}[|Y|^p] \geq \frac{\|a\|_2^{2p+4}}{CK^{p+4}\|a\|_2^{p+4}} = \frac{C}{K^5} \|a\|_2^p \Rightarrow \|Y\|_{L^p} \geq \left(\frac{C}{K^5}\right)^{1/p} \|a\|_2$$

The upper bound is trivial since $\|Y\|_{L^p} \leq \|Y\|_{L^1} \leq c(K)\|a\|_2$.

□

Remark. Note that the proof for the case when $p \in (0, 1)$, the bounds on p where “boosted” by adding a suitably large constant. A larger “boost” would increase the effect of K in the lower bound.

While we are still in the topic of sub-Gaussian norms, we will now discuss an important concept called *centering*. This concerns the sub-Gaussianity of non-zero mean sub-Gaussian random variables once centered.

Lemma 2.5.4. *Let X be a sub-Gaussian random variable with non-zero mean and sub-Gaussian norm $\|X\|_{\psi_2}$. Then:*

$$\|X - \mathbb{E}[X]\|_{\psi_2} \leq C\|X\|_{\psi_2}$$

Proof. We first use the fact that $\|\cdot\|_{\psi_2}$ is a norm, and therefore triangle inequality holds.

$$\|X - \mathbb{E}[X]\|_{\psi_2} \leq \|X\|_{\psi_2} + \|\mathbb{E}[X]\|_{\psi_2} \stackrel{(i)}{\leq} \|X\|_{\psi_2} + \frac{|\mathbb{E}[X]|}{\sqrt{\log(2)}} \leq \|X\|_{\psi_2} + \frac{\mathbb{E}[|X|]}{\sqrt{\log(2)}} \leq C\|X\|_{\psi_2}$$

where in step (i) we have used the fact that $\mathbb{E}[X]$ is a constant random variable, and thus:

$$\mathbb{E} \left[\exp \left(\frac{\mathbb{E}[X]^2}{t^2} \right) \right] = \exp \left(\frac{\mathbb{E}[X]^2}{t^2} \right) \leq 2 \Rightarrow t \geq \frac{|\mathbb{E}[X]|}{\sqrt{\log(2)}}$$

□

Remark. Now, unlike Lemma 2.5.5, $C > 1$. This can be proven by showing that there exists a random variable such that $\|X - \mathbb{E}[X]\|_{\psi_2} > \|X\|_{\psi_2}$. Consider a random variable over $\{-1, +1\}$ which takes -1 with probability 0.25 and 1 with probability 0.75. A simple calculation reveals that $\mathbb{E}[X] = 0.5$. Now:

$$\mathbb{E} \left[\exp \left(\frac{(X - \mathbb{E}[X])^2}{t^2} \right) \right] = 0.75 \cdot e^{\frac{0.25}{t^2}} + 0.25 \cdot e^{\frac{2.25}{t^2}}$$

Now, set $e^{\frac{0.25}{t^2}} = x$ to get $0.75x + 0.25x^9 - 2 \leq 0$. Plotting reveals that $x \leq 1.18$, which implies that $\|X - \mathbb{E}[X]\|_{\psi_2} \approx \frac{0.5}{\sqrt{\log(1.18)}} \approx 1.229$.

However, we know that $\|X\|_{\psi_2} = \frac{1}{\sqrt{\log(2)}} = 1.201$, and therefore proves the assertion.

2.5.1 Auxiliary Lemmata

Lemma 2.5.5. *For any random variable X , we have:*

$$\|X - \mathbb{E}[X]\|_{L^2} \leq \|X\|_{L^2}$$

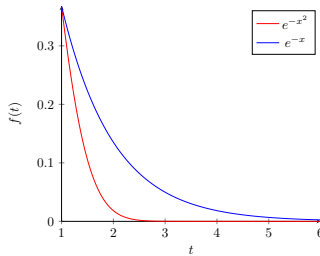
Proof.

$$\begin{aligned} \|X - \mathbb{E}[X]\|_{L^2} &= \sqrt{\mathbb{E}[|X - \mathbb{E}[X]|^2]} \\ &= \sqrt{\mathbb{E}[X^2] + \mathbb{E}[X]^2 - 2\mathbb{E}[X]\mathbb{E}[|X|]} \\ &\leq \sqrt{\mathbb{E}[X^2] - \mathbb{E}[X]^2} \\ &\leq \sqrt{\mathbb{E}[X^2]} = \|X\|_{L^2} \end{aligned}$$

□

2.6 Sub-Exponential Random Variables, Orlicz spaces

As seen earlier, the family of sub-Gaussian random variables encapsulate a vast range of distributions - bounded and Gaussian. To recall, a sub-Gaussian random variable is one where the tails decay at the rate of $\exp(-ct^2)$ (or faster). But what about distributions whose tails decay slower.



For instance, the exponential random variable $X \sim \text{Exp}(\lambda)$ satisfies:

$$\Pr(X > t) = e^{-\lambda t}$$

This doesn't decay as fast as the *sub-Gaussian rate*, as shown in the plot on the right, and therefore is said to have a *heavier* tail than a Gaussian.

Distributions such as these occur more commonly than just the exponential distribution in practice. For instance, let's consider the square of a sub-Gaussian random variable. Let $Y = X^2$ where X is sub-Gaussian with sub-Gaussian norm $\|X\|_{\psi_2} \leq K$. Then:

$$\Pr(Y > t) = \Pr(|X| > \sqrt{t}) \leq 2 \exp \left(-\frac{ct}{K^2} \right)$$

which shows that Y has a tail heavier than a sub-Gaussian. Such distributions are called *Sub-Exponential distributions*. Below is a lemma similar to Lemma 2.4.3 for sub-Exponential distributions that establishes equivalence properties.

Lemma 2.6.1. *Let X be a sub-Exponential random variable. The following statements are equivalent:*

- (i) $\Pr(|X| > t) \leq 2 \exp\left(-\frac{t}{K_1}\right)$ for all $t \geq 0$
- (ii) $\|X\|_{L^p} \leq K_2 p$ for all $p \geq 1$
- (iii) $\mathbb{E}[\exp(\lambda|X|)] \leq \exp(K_3 \lambda)$ for all $\lambda \in \left[0, \frac{1}{K_3}\right]$
- (iv) $\mathbb{E}\left[\exp\left(\frac{|X|}{K_4}\right)\right] \leq 2$
- (v) $\mathbb{E}[\exp(\lambda X)] \leq \exp(K_5^2 \lambda^2)$ for all λ such that $|\lambda| \leq \frac{1}{K_5}$ (if $\mathbb{E}[X] = 0$)

Proof. (i) \Rightarrow (ii) We use the corollary 1.2.2, we have that:

$$\begin{aligned}
\mathbb{E}[|X|^p] &= \int_0^\infty p t^{p-1} \Pr(|X| > t) dt \\
&\leq \int_0^\infty 2 p t^{p-1} \exp\left(-\frac{t}{K_1}\right) dt \\
&= 2 p K_1^p \int_0^\infty u^{p-1} \exp(-u) du \\
&= 2 p K_1^p \Gamma(p) = 2 K_1^p p! \leq (2 K_1)^p p^p
\end{aligned}$$

Therefore:

$$\|X\|_{L^p} = \sqrt[p]{\mathbb{E}[|X|^p]} \leq K_2 p$$

for $K_2 = 2 K_1$.

(ii) \Rightarrow (iii) The expansion of e^x yields:

$$\begin{aligned}
\mathbb{E}[\exp(\lambda|X|)] &= \sum_{k=0}^\infty \frac{\lambda^k \mathbb{E}[|X|^k]}{k!} \\
&\leq \sum_{k=0}^\infty \frac{\lambda^k (2 K_1)^k k!}{k!} \\
&\leq \sum_{k=0}^\infty (2 \lambda K_1)^k
\end{aligned}$$

This is geometric progression that converges if $|2 \lambda K_1| \leq 1 \Leftrightarrow |\lambda| \leq \frac{1}{2 K_1}$. Now we have:

$$\mathbb{E}[\exp(\lambda|X|)] \leq \frac{1}{1 - 2 \lambda K_1} \stackrel{(i)}{\leq} e^{K_3 \lambda}$$

for $K_3 = 4 K_1$ and Step (i) uses Lemma 2.4.6 with $\lambda \in \left[0, \frac{1}{2 K_1}\right]$.

(iii) \Rightarrow (iv) Choose $\lambda = \frac{1}{2 K_3}$ to get:

$$\mathbb{E}\left[\exp\left(\frac{|X|}{2 K_3}\right)\right] \leq \exp(1/2) \leq 2$$

for $K_4 = 2 K_3$.

(iv) \Rightarrow (i) Note that:

$$\Pr(|X| > t) = \Pr\left(\frac{|X|}{K_4} > \frac{t}{K_4}\right) = \Pr\left(\exp\left(\frac{|X|}{K_4}\right) > \exp\left(\frac{t}{K_4}\right)\right) \leq e^{\frac{-t}{K_4}} \mathbb{E}\left[\exp\left(\frac{X}{K_4}\right)\right] \leq 2 e^{-\frac{t}{K_4}}$$

(ii) \Rightarrow (v) Again from the expansion of e^x , we have:

$$\begin{aligned}
\mathbb{E}[\exp(\lambda X)] &= \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} \mathbb{E}[X^k] \\
&\leq 1 + \sum_{k=2}^{\infty} \frac{\lambda^k}{k!} \mathbb{E}[|X|^k] \\
&\stackrel{(i)}{\leq} 1 + \sum_{k=1}^{\infty} \frac{2^k \lambda^{2k}}{(2k)!} \mathbb{E}[|X|^{2k}] \\
&\leq 1 + \sum_{k=1}^{\infty} 2^{2k} \lambda^{2k} K_1^{2k} \\
&\leq \sum_{k=0}^{\infty} (2\lambda K_1)^{2k}
\end{aligned}$$

where Step (i) follows from Lemma 2.4.7.

The infinite geometric series converges if $|\lambda| \leq \frac{1}{2K_1}$. Under this condition:

$$\mathbb{E}[\exp(\lambda X)] \leq \frac{1}{1 - (2\lambda K_1)^2} \leq e^{8\lambda^2 K_1^2} = e^{K_5^2 \lambda^2}$$

where $K_5 = 2K_1$.

(v) \Rightarrow (ii) We will make use of this interesting inequality (refer Lemma 2.6.5 for a proof)

$$|x|^p \leq p^p (e^x + e^{-x})$$

$$\mathbb{E}[|X|^p] \leq p^p \mathbb{E}[e^X] + p^p \mathbb{E}[e^{-X}] \leq p^p \exp(K_5^2) + p^p \exp(K_5^2) \leq 2e^{K_5^2} p^p \Rightarrow \|X\|_{L^p} \leq K_6 p$$

for $K_6 = 2e^{K_5^2}$.

□

Remark. Again, note that when the variable is not necessarily mean zero, we have a weaker corollary like in the sub-Gaussian case. We state it without proof below.

Corollary 2.6.1. Let X be a non-zero mean sub-Exponential random variable satisfying

$\Pr(|X| > t) \leq 2 \exp\left(-\frac{t}{K_1}\right)$. Then we have that:

$$\mathbb{E}[\exp(\lambda X)] \leq 2 \exp(\lambda^2 K_5^2)$$

Analogous to the sub-Gaussian norm, we can define a *sub-Exponential norm* for sub-Exponential random variables. Formally, we define the sub-Exponential norm of a sub-Exponential random variable X as:

$$\|X\|_{psi_1} = \inf \left\{ t > 0 : \mathbb{E} \left[\exp \left(\frac{|X|}{t} \right) \right] \leq 2 \right\}$$

Just as noted earlier for sub-Gaussian random variables, we see that the sub-Exponential norm is merely the tightest constant in property (iii) in Lemma 2.6.1. Let us restate the equivalences in Lemma 2.6.1 in terms of this norm:

$$(i) \Pr(|X| > t) \leq 2 \exp\left(-\frac{C_1 t}{\|X\|_{psi_1}}\right)$$

(ii) $\|X\|_{L^p} \leq C_2 \|X\|_{\psi_1} p$ for all $p \geq 1$

(iii) $\mathbb{E}[\exp(\lambda|X|)] \leq \exp(C_3 \lambda \|X\|_{\psi_1})$ for all $\lambda \in \left[0, \frac{1}{C_3 \|X\|_{\psi_1}}\right]$

(iv) If $\mathbb{E}[X] = 0$, then $\mathbb{E}[\exp(\lambda X)] \leq \exp(C_4 \lambda^2 \|X\|_{\psi_1}^2)$ for all λ such that $|\lambda| \leq \frac{1}{C_4 \|X\|_{\psi_1}}$

where C_1, C_2, C_3, C_4 are universal constants.

As shown via an example earlier, we know that the square of a sub-Gaussian random variable is a sub-Exponential random variable. Remarkably, the sub-Gaussian and sub-Exponential norms of the variable and its square are directly related as shown in the lemma below.

Lemma 2.6.2. *Let X be a sub-Gaussian random variable. Then:*

$$\|X\|_{\psi_2}^2 = \|X^2\|_{\psi_1}$$

Proof. $\|X\|_{\psi_2}$ is the infimum of the values taken by $t > 0$ that satisfy $\mathbb{E}\left[\exp\left(\frac{X^2}{t^2}\right)\right] \leq 2$. Similarly, $\|X\|_{\psi_1}$ is the infimum of the values taken by $t > 0$ that satisfy $\mathbb{E}\left[\exp\left(\frac{|X|}{t}\right)\right] \leq 2$ which is equivalent to $\mathbb{E}\left[\exp\left(\frac{X^2}{t^2}\right)\right] \leq 2$, and the statement is proved. \square

We can generalize further to product of two different sub-Gaussian random variables.

Lemma 2.6.3. *Let X and Y be two sub-Gaussian random variables. Then XY is a sub-exponential random variable. Moreover:*

$$\|XY\|_{\psi_1} \leq \|X\|_{\psi_2} \|Y\|_{\psi_2}$$

Proof. **TODO**

\square

Now, we turn to a version of the centering lemma for the sub-Exponential norm.

Lemma 2.6.4. *Let X be a sub-Exponential random variable. Then $\|X - \mathbb{E}[X]\|_{\psi_1} \leq C \|X\|_{\psi_1}$.*

Proof. The proof is the same as the proof for the centering lemma for the sub-Gaussian norm. Since $\|\cdot\|_{\psi_1}$ is a norm:

$$\begin{aligned} \|X - \mathbb{E}[X]\|_{\psi_1} &\leq \|X\|_{\psi_1} + \|\mathbb{E}[X]\|_{\psi_1} \\ &\leq \|X\|_{\psi_1} + \frac{|\mathbb{E}[X]|}{\log(2)} \\ &\leq \|X\|_{\psi_1} + \frac{\mathbb{E}[|X|]}{\log(2)} \\ &\leq \|X\|_{\psi_1} + C_1 \|X\|_{\psi_1} = C \|X\|_{\psi_1} \end{aligned}$$

\square

Remark. Do we have that $\|X - \mathbb{E}[X]\|_{\psi_1} > \|X\|_{\psi_1}$ for some sub-Exponential random variable?

There is a lot of commonality in the properties of sub-Gaussian and sub-Exponential distributions, so naturally, can we generalize this any further? The answer is yes, and introduces us to *Orlicz spaces*. First, we define an *Orlicz function*, which is a convex, increasing function $\psi : [0, \infty) \rightarrow [0, \infty)$ satisfying $\psi(0) = 0$ and $\lim_{x \rightarrow \infty} \psi(x) = \infty$.

The *Orlicz norm* of a random variable X is:

$$\|X\|_\psi = \inf \left\{ t > 0 : \mathbb{E} \left[\psi \left(\frac{|X|}{t} \right) \right] \leq 1 \right\}$$

Note that for $\psi(x) = e^{x^2} - 1$, we retrieve the sub-Gaussian norm, and for $\psi(x) = e^x - 1$, we obtain the sub-Exponential norm. The *Orlicz space* is a collection of random variables that have finite *Orlicz norm*. That is: $L_\psi(\Omega, \Sigma, \mathbb{P}) = \{X \in (\Omega, \Sigma, \mathbb{P}) : \|X\|_\psi < \infty\}$.

The Orlicz norm is a valid norm over the Orlicz space, and therefore satisfies:

- $\|\mathbf{0}\|_\psi = 0$
- $\|\alpha X\|_\psi = \alpha \|X\|_\psi$
- $\|X + Y\|_\psi \leq \|X\|_\psi + \|Y\|_\psi$.

Note that there is a natural hierarchy between L_p -spaces as follows. $L_\infty \subset L_{\psi_2}$ because any distribution with bounded supremum is sub-Gaussian. $L_{\psi_2} \subset L_p$ because any sub-Gaussian distribution has bounded L^p norm.

2.6.1 Auxiliary Lemmata

Lemma 2.6.5. *For any $x \in \mathbb{R}$ and $p > 0$, we have:*

$$|x|^p \leq p^p (e^x + e^{-x})$$

Proof.

$$\begin{aligned} e^x + e^{-x} &= \left(e^{\frac{x}{p}}\right)^p + \left(e^{-\frac{x}{p}}\right)^p \\ &\geq 2 \left(\frac{e^{\frac{x}{p}} + e^{-\frac{x}{p}}}{2}\right)^p \\ &\geq \left(\frac{e^{\frac{x}{p}} + e^{-\frac{x}{p}}}{2}\right)^p \\ &= \cosh^p\left(\frac{x}{p}\right) \\ &\geq \left(\frac{|x|}{p}\right)^p \end{aligned}$$

where in the final step we have used $\cosh(x) \geq |x|$ for $x \in \mathbb{R}$. □

2.7 Bernstein's Inequality for sub-Exponential random variables

With the tails properties of a sub-Gaussian random variable, we could derive concentration inequalities for sums of sub-Gaussian random variables. Analogously, we will derive concentration inequalities for sub-Exponential random variables. The first is *Bernstein's inequality*, as stated below:

Theorem 2.7.1. *Let $\{X_i\}_{i=1}^n$ be a collection of n independent, zero-mean, sub-Exponential random variables. Then for $t \geq 0$, we have:*

$$\Pr \left(\left| \sum_{i=1}^n X_i \right| > t \right) \leq 2 \exp \left(-c \min \left\{ \frac{t^2}{\sum_{i=1}^n \|X_i\|_{\psi_1}^2}, \frac{t}{\max_i \|X_i\|_{\psi_1}} \right\} \right)$$

where $c > 0$ is a universal constant.

Proof. We will consider the event $\sum_{i=1}^n X_i > t$ for convenience.

By Markov's inequality (Lemma 1.2.2), and $\lambda > 0$,

$$\Pr\left(\sum_{i=1}^n X_i > t\right) = \Pr\left(\lambda \sum_{i=1}^n X_i > \lambda t\right) = \Pr\left(\exp\left(\lambda \sum_{i=1}^n X_i\right) > \exp(\lambda t)\right) \leq e^{-\lambda t} \mathbb{E}\left[\exp\left(\lambda \sum_{i=1}^n X_i\right)\right]$$

Since $\{X_i\}_{i=1}^n$ are independent:

$$\mathbb{E}\left[\exp\left(\lambda \sum_{i=1}^n X_i\right)\right] = \prod_{i=1}^n \mathbb{E}[e^{\lambda X_i}]$$

Using (v) from Lemma 2.6.1, we have that:

$$\mathbb{E}[\exp(\lambda X_i)] \leq \exp(C\lambda^2 \|X_i\|_{\psi_1}^2)$$

for $|\lambda| \leq \frac{1}{C\|X_i\|_{\psi_1}}$. And therefore:

$$e^{-\lambda t} \mathbb{E}\left[\exp\left(\lambda \sum_{i=1}^n X_i\right)\right] \leq \exp\left(-\lambda t + C\lambda^2 \sum_{i=1}^n \|X_i\|_{\psi_1}^2\right)$$

for $\lambda \in \left(0, \frac{1}{C\max_i \|X_i\|_{\psi_1}}\right)$.

λ that minimizes the inequality is $-\frac{t}{2C\sum_{i=1}^n \|X_i\|_{\psi_1}^2}$. However, $\lambda < \frac{1}{C\max_i \|X_i\|_{\psi_1}}$, and this gives rise to

$\lambda = \min\left\{\frac{1}{C\max_i \|X_i\|_{\psi_1}}, -\frac{t}{2\sum_{i=1}^n \|X_i\|_{\psi_1}^2}\right\}$, we gives the result:

$$\Pr\left(\sum_{i=1}^n X_i > t\right) \leq \exp\left(-c \min\left\{\frac{t^2}{\sum_{i=1}^n \|X_i\|_{\psi_1}^2}, \frac{t}{\max_i \|X_i\|_{\psi_1}}\right\}\right)$$

By symmetry, we have:

$$\Pr\left(\left|\sum_{i=1}^n X_i\right| > t\right) \leq 2 \exp\left(-c \min\left\{\frac{t^2}{\sum_{i=1}^n \|X_i\|_{\psi_1}^2}, \frac{t}{\max_i \|X_i\|_{\psi_1}}\right\}\right)$$

□

Similar to the sub-Gaussian case, we have an immediate corollary for linear combination of sub-Exponential random variables:

Corollary 2.7.1. *Let $\{X_i\}_{i=1}^n$ be a collection of n independent, zero-mean, sub-Exponential random variables and let $a \in \mathbb{R}^n$ be a vector. Then for $t \geq 0$, we have:*

$$\begin{aligned} \Pr\left(\left|\sum_{i=1}^n a_i X_i\right| > t\right) &\leq 2 \exp\left(-c \min\left\{\frac{t^2}{\sum_{i=1}^n a_i^2 \|X_i\|_{\psi_1}^2}, \frac{t}{\max_i a_i \|X_i\|_{\psi_1}}\right\}\right) \\ &\leq 2 \exp\left(-c \min\left\{\frac{t^2}{\|a\|_2^2 K^2}, \frac{t}{K\|a\|_\infty}\right\}\right) \end{aligned}$$

where $K = \max_i \|X_i\|_{\psi_1}$.

Remark. Set $a_i = \frac{1}{n}$ for all $i \in [n]$ gives a mean concentration inequality like so:

$$\Pr(|\bar{X}_n| > t) \leq 2 \exp\left(-c \cdot n \min\left\{\frac{t^2}{K^2}, \frac{t}{K}\right\}\right)$$

where \bar{X}_n is the mean of the n random variables $\{X_i\}_{i=1}^n$.

Looking at this inequality carefully reveals two phases: a sub-Exponential phase and a sub-Gaussian phase. When $t \leq K$, we have the sub-Gaussian phase - where the tail decays $\propto e^{-t^2}$, and when $t > K$, we have a sub-Exponential phase - where the tail decays $\propto e^{-t}$.

While bounded random variables are sub-Gaussian, and hence sub-Exponential, we can use the same result from above to provide concentration inequalities for sums of these random variables. However, a tighter analysis is possible, which is performed in the following theorem.

Theorem 2.7.2. *Let $\{X_i\}_{i=1}^n$ be a set of independent, zero-mean random variables such that $|X_i| \leq K$ for all $i \in [n]$. Then for $t \geq 0$, we have:*

$$\Pr\left(\left|\sum_{i=1}^n X_i\right| > t\right) \leq 2 \exp\left(-\frac{t^2/2}{\sigma^2 + Kt/3}\right)$$

where $\sigma^2 = \sum_{i=1}^n \text{Var}(X_i)$.

Proof. Just as done before, we will use Markov's inequality (Lemma 1.2.2) and $\lambda > 0$ to get:

$$\Pr\left(\sum_{i=1}^n X_i > t\right) = \Pr\left(\lambda \sum_{i=1}^n X_i > \lambda t\right) = \Pr\left(\exp\left(\lambda \sum_{i=1}^n X_i\right) > \exp(\lambda t)\right) \leq e^{-\lambda t} \mathbb{E}\left[\exp\left(\lambda \sum_{i=1}^n X_i\right)\right]$$

Now, by Lemma 2.7.1 we have that:

$$\mathbb{E}[\exp(\lambda X_i)] \leq \exp(g(\lambda)\mathbb{E}[X_i^2])$$

where $g(\lambda) = \frac{\lambda^2/2}{1-|\lambda|K/3}$ for $|\lambda| \leq \frac{3}{K}$.

This gives:

$$\Pr\left(\sum_{i=1}^n X_i > t\right) \leq \exp(-\lambda t + g(\lambda)\sigma^2)$$

We intend on finding the λ that minimizes the RHS for the tightest bound. Now:

$$\inf_{\lambda > 0} g(\lambda)\sigma^2 - \lambda t = -\sup_{\lambda > 0} \lambda t - g(\lambda)\sigma^2 \stackrel{(i)}{\leq} -\frac{t^2/2}{\sigma^2 + Kt/3}$$

where in Step (i) we have used Lemma 2.7.2 and this gives:

$$\Pr\left(\sum_{i=1}^n X_i > t\right) \leq \exp\left(-\frac{t^2/2}{\sigma^2 + Kt/3}\right)$$

By symmetry, we obtain the final result. □

2.7.1 Auxiliary Lemmata

Lemma 2.7.1.

Lemma 2.7.2.

References

- Luc Devroye and Gábor Lugosi. *Combinatorial methods in density estimation*. Springer Science & Business Media, 2012.
- Roman Vershynin. *High-dimensional probability*. 2019.
- Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.