
—

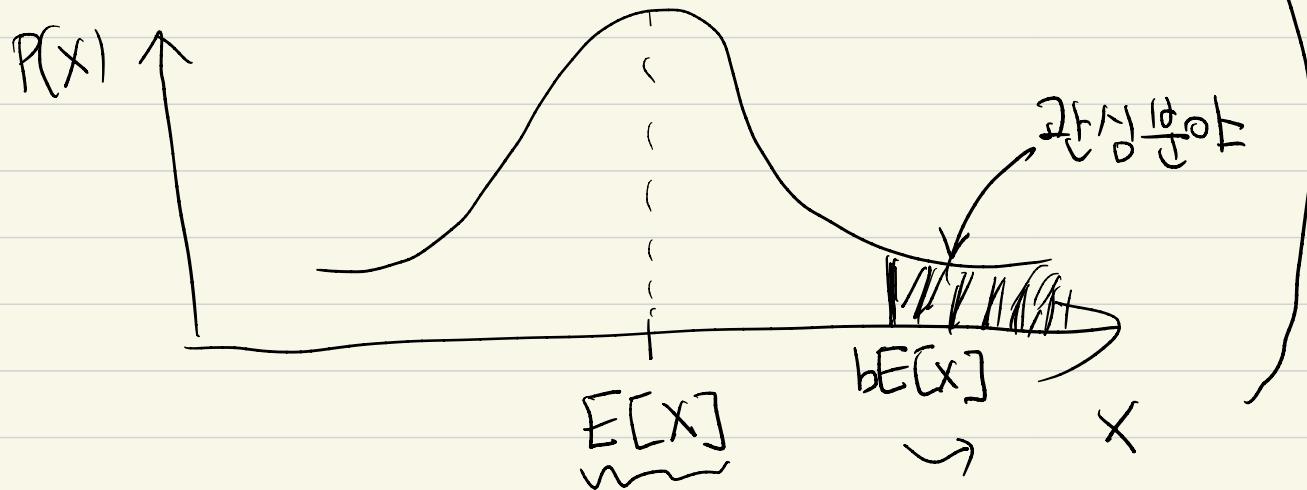


• Learning Theory 러닝 퍼리 오리

210311

- Concentration Inequality \leftarrow Statistics

• Tail Probability



• Concentration Property

How the r.v.s are concentrated around $E[X]$?

Theorem.

• Markov Inequality

nonnegative X , any $a > 0$

$$P[X \geq a] \leq \frac{E[X]}{a}$$

pf) $E[X] = E[X \cdot 1_{\{X \geq a\}}] \geq E[a \cdot 1_{\{X \geq a\}}] = aP(X \geq a)$

↳ 초기 조건, $X=0, a$ 만 ...

Corollary. $a = bE[X] \Rightarrow P(X \geq bE[X]) \leq \frac{1}{b}$ ($b > 0$)

Rmk

- 등호 조건
- nonnegative

Corollary. Chebyshev's Inequality

$$\begin{aligned} P(|X - E[X]| \geq a) &= P((X - E[X])^2 \geq a^2) \\ &\leq \frac{E((X - E[X])^2)}{a^2} = \frac{\text{Var}(X)}{a^2}. \end{aligned}$$

$$a \leftarrow a\sigma$$

$$P(|X - E[X]| \geq a\sigma) \leq \frac{\text{Var}(X)}{a^2\sigma^2} = \frac{1}{a^2}.$$

Theorem. Law of Large Numbers (Weak)

$$S_n = X_1 + \dots + X_n$$

$$P\left(\left|\frac{1}{n}S_n - E[X]\right| \geq \varepsilon\right) \leq \frac{1}{n} \frac{\text{Var}(X)}{\varepsilon^2}$$

✓

• Markov 사용법 ✓

$X: RV, t \in \mathbb{R}_+$ (Ω, \mathcal{F}, P) : probability space
event
 $E \in \mathcal{F}, E = \{X \geq t\} = \{\omega \in \Omega : X(\omega) \geq t\} = \{\phi(x) \geq \phi(t)\}$
under 특정 condition of ϕ .

ϕ 를 써워준 품은 Markov Ineq. 주해주는 Bound가 된다!
non-negative, non-decreasing

$$P(X \geq t) = P(\phi(X) \geq \phi(t))$$

$$\leq \frac{E[\phi(X)]}{\phi(t)} = \phi(t)^{-1} E[\phi(X)]$$

$$(f) \phi(x) = x^q \quad q > 0$$

$$E[(x - E[x])^q]$$

$$P(|X - E[X]| \geq \varepsilon) \leq \frac{E[(x - E[x])^q]}{\varepsilon^q}$$

• MGF, Chernoff Bounds

$\phi(x) = e^{\theta x}$ 로 두기 ($\theta > 0$).

$$\rightarrow P(X \geq t) \leq e^{-\theta t} \underbrace{E[e^{\theta X}]}_{}, \quad (\text{Chernoff Bound})$$

θ 값에 무관하므로, $P(X \geq t) \leq \inf_{\theta > 0} e^{-\theta t} \underbrace{E[e^{\theta X}]}_{\text{MGF}}$

- Recall MGF. ✓

$$\text{MGF } M_X(\theta) := \underbrace{E[e^{\theta X}]}_{\substack{\text{터일려전개시} \\ \text{증명 쌍 가능}}} \rightarrow E[X^n] = \frac{d^n}{d\theta^n} M_X(\theta) \Big|_{\theta=0}$$

$$(M_X(\theta))' = E[(e^{\theta X})'] = E[X e^{\theta X}]$$

- Cumulant Function. ✓

$$\psi_X(\theta) := \log E[e^{\theta X}] \quad (\text{로그 오른 MGF})$$

비(부)무리화하기

$$\psi_X'(\theta) = \frac{(M_X(\theta))'}{\underbrace{E[e^{\theta X}]}} = \frac{\underbrace{E[X e^{\theta X}]}_{}}{\underbrace{E[e^{\theta X}]}}$$

- Cumulant function is convex. $\psi_x(\theta) = \log \mathbb{E}[e^{\theta x}]$

Hölder Inequality

$$f, g \quad p, q > 0, \quad \frac{1}{p} + \frac{1}{q} = 1$$

$$\int fg \leq \left(\int |f|^p \right)^{1/p} \left(\int |g|^q \right)^{1/q}$$

$p=q=2$ Cauchy-Schwarz, +

$$\frac{1}{p} + \frac{1}{q} = 1$$

$$f(\lambda \theta_1 + (1-\lambda) \theta_2) \leq \lambda f(\theta_1) + (1-\lambda) f(\theta_2)$$

$$f \leftarrow \mathbb{E}[e^{\theta X}] \quad p = \frac{1}{\lambda} \quad q = \frac{1}{1-\lambda}$$

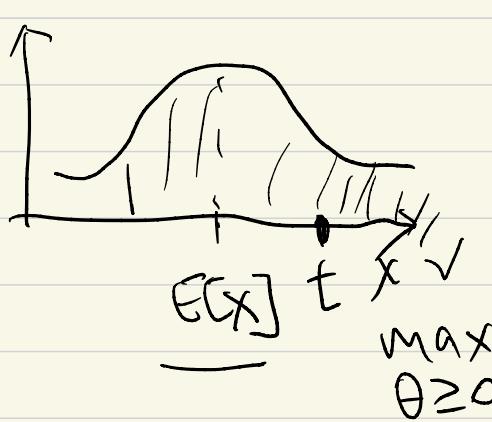
$$\mathbb{E}[e^{(\lambda \theta_1 + (1-\lambda) \theta_2)}]$$

From Markov... inequality

$$\phi(x) = e^{\theta x}$$

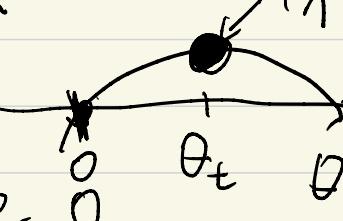
$$P(X \geq t) \leq e^{-\theta t} \cdot E[e^{\theta X}] = e^{-\theta t} \cdot e^{\log E[e^{\theta X}]}$$

$$= e^{-\theta t + \psi_x(\theta)}, \theta \geq 0$$



$$(E\theta - \psi_x(\theta))$$

concave



$$\psi_x^*(t) := \sup_{\theta \geq 0} (t\theta - \psi_x(\theta)) \quad \text{Cramer-transform}$$

$$\theta_t: (t - \psi_x(\theta)) = 0$$

$$\psi'_x(\theta) = t$$

$$\psi'_x(\theta) = t \quad \text{if } \theta \leq \theta_t$$

$$P(X \geq t) \leq e^{-\psi_x^*(t)}$$

$$\theta=0: t - \psi'_x(0) = t - E[X] > 0 \rightarrow \theta_t > 0$$

$$\psi_x^*(t) = \sup_{\theta \in \mathbb{R}} (t\theta - \psi_x(\theta))$$

$\psi_x(\theta)$ Lf Fenchel-Legendre Dn al.
Conjugate

정리하면...

$$P(X \geq t) \leq e^{-t\theta + \mathcal{V}_X(\theta)}$$

$$e^{-t\theta + \frac{1}{2}\theta^2\sigma^2}$$

① Gaussian $X \sim N(0, \sigma^2)$

$$\mathcal{V}_X(\theta) = \log \mathbb{E}[e^{\theta X}] = \log(e^{\frac{1}{2}\theta^2\sigma^2})$$

$$P(X \geq t) \leq e^{-\frac{t^2}{2\sigma^2}}$$

근데 정규분포는 수워서 사실 계산가능

② 안 가우시안 시리즈

- Poisson $F(x=k) = \frac{\lambda^k e^{-\lambda}}{k!} \quad X \text{ nonneg.}$

$$\mathbb{E}[e^{\theta X}] = e^{\lambda e^\theta - \lambda - \lambda \theta}$$

$$\mathcal{V}_X^*(t) = \lambda h\left(\frac{t}{\lambda}\right) \quad \text{where } h(x) = ((1+x)\log(1+x) - x\log x).$$

* Bernoulli (p) $\sim Y$ $Y = 1 \quad P(\text{prob})$

$$\cdot \mathbb{E}[Y] = p.$$

$$\cdot X = Y - P.$$

$$X = 1-p \quad P(\text{prob})$$

$$-p \quad \text{otherwise}$$

$$\mathbb{E}[e^{\theta X}] = p \cdot e^{\theta(1-p)} + (1-p) e^{\theta(-p)} = e^{-p\theta} (pe^\theta + 1-p)$$

$$\log \mathbb{E}[e^{\theta X}] = -p\theta + \log(pe^\theta + 1-p) = \mathcal{V}_X(\theta)$$

• $\exists \theta - \psi_x(\theta)$ 미분.

$$\psi_x(\theta) = \log(p e^\theta + 1-p) - p\theta$$

$$\frac{d}{d\theta} (\theta + p\theta - \log(p e^\theta + 1-p)) = 0$$

$$\begin{aligned}\theta_E &= \log \frac{(1-p)(p+t)}{p(1-p-t)} & g &:= p+t \\ &= \log \frac{g \cdot (1-p)}{p \cdot (1-g)}\end{aligned}$$

$$\Rightarrow \psi_x^*(t) = \underbrace{KL(P_g || P_p)}_{\text{구현가능.}} \quad \text{설명: } \dots$$

$$= g \log \frac{g}{p} + (1-g) \log \frac{1-g}{1-p}$$

- \Rightarrow 더 깔끔하게 ...

$$f(t) = (p+t) \log \frac{p+t}{p} + (1-p-t) \log \frac{1-p-t}{1-p}$$

$$f'(t) = \log \frac{p+t}{p} - \log \frac{1-p-t}{1-p}$$

$$f''(t) = \frac{1}{(p+t)(1-p-t)} = \frac{1}{a(1-a)} \geq \frac{4}{(a+(1-a))^2}$$

$$\Rightarrow f(t) = f(0) + f'(0)t + \frac{1}{2} f''(\xi) t^2 \geq 2t^2 \geq \frac{4}{2} t^2$$

$$P(X \geq t) \leq e^{-2t}$$

- Independent sum.

$$\mathbb{E}[X_1]\mathbb{E}[X_2] = \mathbb{E}[X_1 X_2]$$

$$\mathbb{E}[X_i] = 0, \quad X_1, X_2, \dots (i.i.d), \quad S_n = X_1 + \dots + X_n$$

$$\mathbb{E}[e^{\theta S_n}] = \mathbb{E}\left[e^{\theta(X_1 + X_2 + \dots + X_n)}\right] = \prod_{i=1}^n \mathbb{E}[e^{\theta X_i}]$$

$$\psi_{S_n}(\theta) = \sum_{i=1}^n \psi_{X_i}(\theta) = n \psi_{X_1}(\theta)$$

$$\psi_{S_n}^*(t) = \sup_{\theta} \{ t\theta - \psi_{X_1}(\theta) \}$$

$$= n \cdot \sup \left\{ \left(\frac{t}{n} \right) \theta - \psi_{X_1}(\theta) \right\}$$

$$= n \psi_{X_1}^*\left(\frac{t}{n}\right)$$

$$P(S_n \geq t) \leq e^{-n \psi_{X_1}^*\left(\frac{t}{n}\right)}$$

- Bernoulli 일 경우.

$$X_i \sim \text{Bernoulli}(p) \Rightarrow S_n \sim \text{Binomial}(n, p)$$

$$P(S_n - np \geq t) \leq e^{-n \cdot \psi_{X_1}^*\left(\frac{t}{n}\right)}$$

$$\leq e^{-2n\left(\frac{t}{n}\right)^2} = e^{-2\frac{t^2}{n}} \quad \checkmark$$

$$t = \underbrace{n \cdot l}_{\text{!}}$$

$$= e^{-2nl^2}$$

• 다른 부분

- Sub Gaussian RV. ($sG(\sigma)$)

Theorem 2.1.1. Let X be a centered random variable on \mathbb{R} , each statement below implies the next (we take $\sigma^2 > 0$ in the first definition as a variance proxy).

- Laplace transform: for any $s \in \mathbb{R}$, $\mathbb{E}[\exp(sX)] \leq \exp\left(\frac{\sigma^2 s^2}{2}\right)$.
- Concentration: for any $t > 0$, $\max\{\mathbb{P}(X \geq t), \mathbb{P}(X \leq -t)\} \leq \exp\left(\frac{-t^2}{2\sigma^2}\right)$.
- Moment condition: for any $q \in \mathbb{N}^*$, $\mathbb{E}[X^{2q}] \leq q!(4\sigma^2)^q$.
- Orlicz condition: $\mathbb{E}\left[\exp\left(\frac{X^2}{8\sigma^2}\right)\right] \leq 2$.
- Laplace transform: for any $t \in \mathbb{R}$, $\mathbb{E}[\exp(tX)] \leq \exp\left(\frac{24\sigma^2 t^2}{2}\right)$.

$$\mathbb{E}[e^{\theta X}] \leq e^{\frac{1}{2}\theta^2 \sigma^2}, \quad \psi_X(\theta) \leq \frac{1}{2}\theta^2 \sigma^2 \text{ 인 경우}$$

sub-gaussian!

$$\psi_X^A(t) = \sup_{\theta} \{t\theta - \psi_X(\theta)\} \geq \sup_{\theta} \{t\theta - \frac{1}{2}\theta^2 \sigma^2\} = \frac{t^2}{2\sigma^2}$$

$$\text{기본적으로 } \mathbb{P}(X \geq t) \leq e^{-\psi_X^A(t)} \leq e^{-\frac{t^2}{2\sigma^2}} \quad [\text{Thm 2.1}]$$

← Sum of i.i.d vars

$$\mathbb{E}[e^{\theta S_n}] \leq e^{\frac{1}{2}\theta^2 \sum_{i=1}^n \sigma_i^2}, \quad \mathbb{P}(S_n \geq t) \leq e^{-\frac{t^2}{2\sum_{i=1}^n \sigma_i^2}}$$

- Bounded RVs. $X_i \in [a, b]$

(Claim) Bounded RVs are sub-Gaussian. RV

(Hoeffding Lemma)

$$Y \in [a, b], E[Y] = \theta \quad \frac{1}{2}\theta^2\sigma^2$$

$$\mathcal{U}_Y(\theta) \leq \frac{1}{8}(b-a)^2\theta^2, \text{ s.t. } Y \sim SG\left(\frac{1}{2}(b-a)\right)$$

pf) $\left| Y - \frac{a+b}{2} \right| \leq \frac{b-a}{2}$

$$\text{Var}(Y) = \text{Var}\left(Y - \frac{b+a}{2}\right) \leq E\left[\left(Y - \frac{a+b}{2}\right)^2\right] \leq \left(\frac{b-a}{2}\right)^2$$

$E[e^{\theta Y}]$ well-defined, $\mathcal{U}_Y(\theta)$ differentiable

$$\mathcal{U}'_Y(\theta) = \frac{E[Y e^{\theta Y}]}{E[e^{\theta Y}]} = \boxed{E[Y e^{\theta Y} - \mathcal{U}_Y(\theta)]} = e^{\log E[e^{\theta Y}]} = \mathcal{U}_Y(\theta)$$

• Brilliant idea: $Q(A) := E[1_A \cdot e^{\theta Y}]$

$\Rightarrow Q$ becomes probability measure !!!!!!!

$$P(A) = E[1_A]$$

$$Q(\Omega) = E[] = 1$$

Ω, \mathcal{F}, P

Rmk.

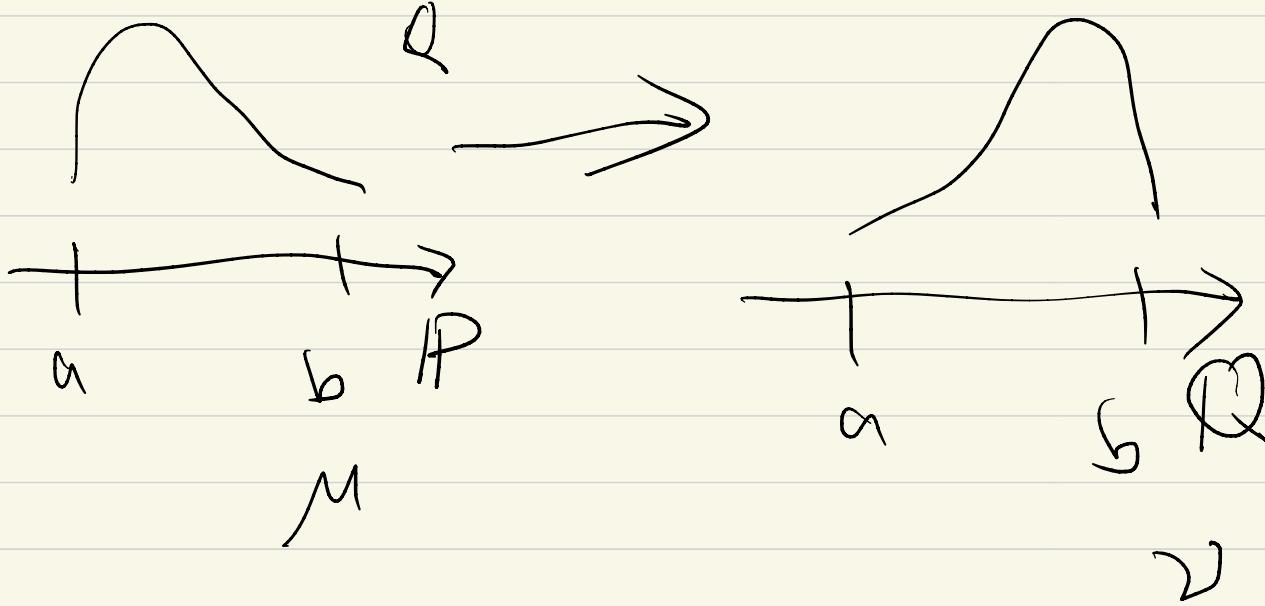
Radon-Nikodym Theorem.

Measurable Space (X, Σ) .

$\nu \ll \mu$.

($\mu(A) > 0$ implies $\nu(A) > 0$) absolutely continuous

$$\nu(A) = \int_A f d\mu \Rightarrow f := \left[\frac{d\nu}{d\mu} \right]_{\geq 0} = 0$$



$$\mu(A) = 0 \rightarrow \nu(A) = 0$$

$$\nu(A) = \int_A f d\mu = \int_A \left(\frac{d\nu}{d\mu} \right) d\mu$$

$P \Rightarrow Q$ measure change $dQ = e^{\theta Y - \psi_Y(\theta)} dP$

$$\psi'_Y(\theta) = E[Y e^{\theta Y - \psi_Y(\theta)}] = E_Q[Y]$$

$$\begin{aligned} \psi''_Y(\theta) &= \left(\frac{E[Y e^{\theta Y}]}{E[e^{\theta Y}]} \right)' = \frac{E[Y^2 e^{\theta Y}] E[e^{\theta Y}] - E[Y e^{\theta Y}]^2}{E[e^{\theta Y}]^2} \\ &= \frac{E[Y^2 e^{\theta Y}]}{E[e^{\theta Y}]} - \left(\frac{E[Y e^{\theta Y}]}{E[e^{\theta Y}]} \right)^2 \end{aligned}$$

$$\begin{aligned} &= E_Q[Y^2] - E_Q[Y]^2 = \text{Var}_Q[Y] \\ &\leq \frac{(b-a)^2}{4} \end{aligned}$$

$$\psi_Y(\theta) = \psi_Y(0) + \psi'_Y(0)\theta + \frac{1}{2}\psi''_Y(\xi)\theta^2 \leq \frac{(b-a)^2}{8}\theta^2$$

so?

Rademacher distribution



+1 ↗

-1 ↘

Thm (Hoeffding)

$$Y_i \in [a_i, b_i] \quad Y_i - E[Y_i] \in [a_i - E[Y_i], b_i - E[Y_i]]$$

$$P\left(\sum_{i=1}^n (Y_i - E[Y_i]) \geq t\right) \leq e^{-\frac{t^2}{2 \sum_{i=1}^n (\frac{b_i - a_i}{2})^2}}$$

$$P\left(\frac{1}{n} \sum_{i=1}^n (Y_i - E[Y_i]) \geq t\right) \leq \exp\left(-\frac{2n^2 t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

- Bernoulli의 경우 앞에서 했던 대로 양하고 $a_i = 0, b_i = 1$ 하면

- 사실 variance of $Sgt = \sum_{i=1}^n (b_i - a_i)^2$ 보다 작은 경우
가 많아!

\Rightarrow Bernstein, Bennett (Boucheron, (2.10))

MGF에 더 tighter bound 주는 것이 가능함....

- Martingale 마팅게일, Hoeffding (1963)

$$\mathbb{E}[X_{n+1} | X_0, X_1, \dots, X_n] = X_n \quad \text{ex) } \frac{\sum x_i}{n} !$$

$$\sigma(X_0, \dots, X_n) = \mathcal{F}_n \quad (\text{information} \dots)$$

$$\mathbb{E}[X_{n+1} | \mathcal{F}_n] = X_n.$$

• Azuma-Hoeffding

$$X_0, X_1, \dots \text{ martingale} \quad a_i \leq X_i - X_{i-1} \leq b_i \quad \forall i \in \mathbb{N}.$$

$$Y_i := X_i - X_{i-1}$$

$$\mathbb{E}[Y_{k+1} | \mathcal{F}_k] = \mathbb{E}[X_{k+1} - X_k | \mathcal{F}_k] = \mathbb{E}[X_{k+1} | \mathcal{F}_k] - X_k$$

$$S_k := \sum_{i=1}^k Y_i = 0$$

$$\text{Hoeffding's Lemma} \quad \mathbb{E}[e^{\theta Y_{k+1}} | \mathcal{F}_k] \leq e^{\frac{\theta^2(b_k - a_k)^2}{8}}$$

$$\mathbb{E}[e^{\theta S_k}] = \mathbb{E}\left[\mathbb{E}[e^{\theta S_k} | \mathcal{F}_{k-1}]\right] \quad (\text{Tower Rule})$$

$$= \mathbb{E}\left[e^{\theta S_{k-1}} \cdot \mathbb{E}[e^{\theta Y_k} | \mathcal{F}_{k-1}]\right] \leq \mathbb{E}[e^{\theta S_{k-1}}] e^{\frac{\theta^2(b_k - a_k)^2}{8}}$$

$$\dots \leq \exp\left(\frac{1}{2}\theta^2 \sum_{i=1}^n \left(\frac{b_i - a_i}{2}\right)^2\right)$$

$$\therefore S_k = \sum_{i=1}^k Y_i = \sum_{i=1}^k (X_i - X_{i-1}) = X_k - X_0,$$

$$X_k - X_0 \sim \mathcal{N}\left(\sqrt{\frac{1}{4} \sum_{i=1}^{mk} (b_i - a_i)^2}\right)$$

나머진 Hoeffding의 증명 과정과 동일

- McDiarmid Inequality.

- Reference
 - Subgaussian random variables: An expository note
Omar Rivasplata
 - Boucheron, Concentration Inequalities, Oxford 2012
 - course A7603 KAIST : course notes
MAS480