- lecture Notes 2 - supplementary
- go over details you've highlighted

Theorem 1 - Gaussian Tail Inequality

1. Probability Inequal.
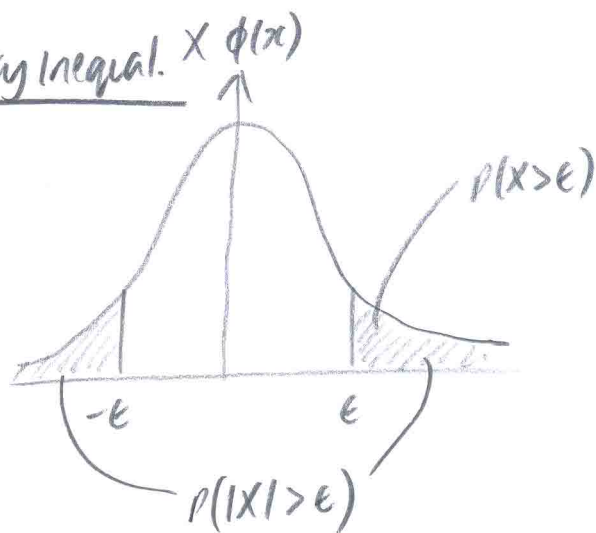


A (single r.v.)

- $X \sim N(0,1)$
- $P(|X| > \epsilon) \leq \dfrac{2e^{-\epsilon^2/2}}{\epsilon}$

- B Multiple r.vs.

- $X_1, \ldots, X_n \sim N(0,1)$
- $P(|\bar{X}_n| > \epsilon) \leq \dfrac{2}{\sqrt{n}\epsilon} e^{-n^2 \epsilon/2}$  large n $\leq e^{-n\epsilon^2/2}$

Proof (A)

- 2 parts :- i) Prove $P(X > \epsilon) \leq \dfrac{e^{-\epsilon^2/2}}{\epsilon}$    · Density of X: $\phi(x) = \dfrac{1}{\sqrt{2\pi}} e^{-x^2/2}$
          ii) use symmetry argument   (std Normal)

i)

$P(X > \epsilon) = \displaystyle\int_\epsilon^\infty \phi(s)\, ds = \underbrace{\int_\epsilon^\infty \frac{s}{s} \phi(s)\, ds}_{(i)} \leq \underbrace{\frac{1}{\epsilon} \int_\epsilon^\infty s\phi(s)\, ds}_{(ii)}$

And $\underbrace{\dfrac{1}{\epsilon} \int_\epsilon^\infty s\phi(s)\, ds = -\dfrac{1}{\epsilon}\int_\epsilon^\infty \phi'(s)\, ds}_{(iii)} = \underbrace{\dfrac{\phi(\epsilon)}{\epsilon} \leq \dfrac{e^{-\epsilon^2/2}}{\epsilon}}_{(iv)}$

(i) multiply/divide by s
(ii) over interval $[\epsilon, \infty)$, $s \geq \epsilon$
(iii) $\phi'(s) = -s\phi(s) \implies s\phi(s) = -\phi'(s)$
(iv) note:- $\dfrac{1}{\epsilon\sqrt{2\pi}} e^{-\epsilon^2/2} \leq \dfrac{1}{\epsilon} e^{-\epsilon^2/2}$

$P(|X| > \epsilon) = P(X > \epsilon) + P(-X > \epsilon)$
$= P(X > \epsilon) + P(X < -\epsilon)$
$= 2P(X > \epsilon)$  via appeal to symmetry

· So we have shown (i) : $P(X > \epsilon) \leq \dfrac{e^{-\epsilon^2/2}}{\epsilon}$   ↙ (???)

Part (ii)

- By symmetry, i.e. $P(|X| > \epsilon) = P(X > \epsilon) + P(-X < -\epsilon) = 2P(X > \epsilon)$
of std. Normal

   - we have $P(|X| > \epsilon) \leq \dfrac{2e^{-\epsilon^2/2}}{\epsilon}$

## Proof (B)

- $X_1, \ldots, X_n \sim N(0,1)$
- Then $\bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$ and $\bar{X}_n \sim N(0, \frac{1}{n})$ $\left(\text{recall } \text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}, \sigma^2 = 1\right)$

$$\Rightarrow \bar{X}_n \overset{d}{=} \frac{1}{\sqrt{n}} Z = n^{-\frac{1}{2}} Z \quad \text{where } Z \sim N(0,1)$$

$$\Rightarrow P(|\bar{X}_n| > \epsilon) = P(n^{-\frac{1}{2}} |Z| > \epsilon) = P(|Z| > \sqrt{n}\epsilon) \leq \frac{2}{\sqrt{n}\epsilon} e^{-\frac{n\epsilon^2}{2}}$$

(i) - ? - clarify ⊛

## Intuition + plots - ⊛

### Theorem 2 - Markov's inequality

- Let $X$ be a non-negative random variable and suppose $\mathbb{E}[X]$ exists
- That is $X > 0$, $\mathbb{E}[|X|] < \infty$
- For any $t > 0$,

$$P(X > t) \leq \frac{\mathbb{E}[X]}{t} \tag{1}$$

### Proof

- As $X > 0$,

$$\mathbb{E}[X] = \int_0^\infty x p(x)\, dx \overset{(i)}{=} \int_0^t x p(x)\, dx + \int_t^\infty x p(x)\, dx$$

$$\geq \int_t^\infty x p(x)\, dx \geq t \int_t^\infty p(x)\, dx = t\, P(X > t)$$

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

## Proof steps

i) Partition expectation/integral over $(0,\infty)$ into sum of integrals over $(0,t)$ and $(t,\infty)$

ii) Over interval $[t,\infty)$, $x \geqslant t \Rightarrow \int_t^\infty x\, p(x)\, dx \geqslant t \int_t^\infty p(x)\, dt$

## Theorem 3 - Chebyshev's inequality

- Let $\mu = \mathbb{E}[X]$ and $\sigma^2 = \text{Var}(X)$; then :-

$$P(|X-\mu| \geqslant t) \leqslant \frac{\sigma^2}{t} \quad \text{and} \quad P(|Z| \geqslant k) \leqslant \frac{1}{k^2}$$

- where $Z = \frac{(X-\mu)}{\sigma}$

- Note that $P(|Z| > 2) \leqslant \frac{1}{4} \quad P(|Z| > 3) \leqslant \frac{1}{9}$

- Assuming variance exists :- i.e. $V(|X|) = \int |x-\mu|^2 \, p(x)\, dx$ exists ⑦

### proof:

- via Markov's inequality;

$$P(|X-\mu| \geqslant t) = P(|X-\mu|^2 \geqslant t^2) \leqslant \frac{\mathbb{E}[(X-\mu)^2]}{t^2} = \frac{\sigma^2}{t^2}$$

- setting $t = k\sigma$ :- (explicitly)

$$P(|X-\mu| \geqslant t) = P(|X-\mu| \geqslant k\sigma) = P\left(\frac{|X-\mu|}{\sigma} \geqslant k\right) = P(|Z| \geqslant k) \leqslant \frac{\sigma^2}{(k\sigma)^2}$$

$$\Rightarrow P(|Z| \geqslant k) \leqslant \frac{1}{k^2}$$

- Application to Bernoulli r.v.s (from Wasserman) with exp.

- If $X_1, \ldots, X_n \sim \text{Bernoulli}(p)$, and $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$

- Then $\text{Var}(\bar{X}_n) = \frac{\text{Var}(X_i)}{n} = \frac{p(1-p)}{n}$

- And via Chebyshev:

$$P(|\bar{X}_n - p| > \epsilon) \leqslant \frac{\text{Var}(\bar{X}_n)}{\epsilon^2} = \frac{p(1-p)}{n\epsilon^2} \leqslant \frac{1}{4n\epsilon^2}$$

· As $p(1-p) \leq \frac{1}{4}$ $\forall p$ ✓

(*)
∴ Note how the bound is used here

## 2. Hoeffding's Inequality

- note the proof strategy in notes

### Lemma 4 (Hoeffding's Lemma)

- Suppose $a \leq X \leq b$
- Then

$$\mathbb{E}[e^{tX}] \leq e^{t\mu} e^{\frac{t^2(b-a)^2}{8}}$$  where $\mu = \mathbb{E}[X]$

___

### convexity

- recall - A function is convex iff for each $x, y$ and each $\alpha \in [0,1]$

$$g(\alpha x + (1-\alpha)y) \leq \alpha g(x) + (1-\alpha)g(y)$$

### proof of lemma 4

- Assume $\mu = 0$

Since $a \leq X \leq b$, we write $X$ as a convex combination of $a$ and $b$ :-

$$X = \alpha b + (1-\alpha)a$$

where $\alpha = \frac{(X-a)}{(b-a)}$ and $(1-\alpha) = \frac{b-X}{b-a}$

note that the function $g(y): y \to e^{ty}$ is convex, meaning that

$$e^{tX} \leq \alpha e^{tb} + (1-\alpha)e^{ta}$$

more explicitly; note that we are setting $X = \alpha b + (1-\alpha)a$  $x=b$  $y=a$
in the general definition of a convex function

$$g(X) = g(\alpha b + (1-\alpha)a) \leq \alpha g(b) + (1-\alpha)g(a)$$

$$e^{tX} \leq \alpha e^{tb} + (1-\alpha)e^{ta} = \frac{X-a}{b-a}e^{tb} + \frac{b-X}{b-a}e^{ta}$$

$\mathbb{E}[\cdot]$ both sides; $\mathbb{E}[X] = 0 \Rightarrow$

$$\mathbb{E}[e^{tX}] \leq \frac{e^{tb}}{b-a}\mathbb{E}[X-a] \quad \frac{e^{ta}}{b-a}\mathbb{E}[b-X]$$

**Giving**

$$\mathbb{E}[e^{tx}] \leq \frac{-ae^{tb}}{b-a} + \frac{be^{ta}}{b-a}$$

- At this stage we will express RHS in a form $e^{g(u)}$, using properties of $g(u)$ and Taylor's theorem for 1st three terms (up to quadratic)

- Define:  $u = t(b-a)$

$$g(u) = -\gamma u + \log(1 - \gamma + \gamma e^u)$$

$$\gamma = \frac{-a}{b-a}$$

NOTE: $g(0) = g'(0) = 0$ , $g''(u) \leq \frac{1}{4} \; \forall u > 0$  (*)

- i. Taylor expansion: ①ⓦⓧ

$\exists \, \xi \in (0, u)$ such that

$$g(u) = \underbrace{g(0)}_{=0} + \underbrace{ug'(0)}_{=0} + \frac{u^2}{2}g''(\xi) = \frac{u^2}{2}\underbrace{g''(\xi)}_{\leq \frac{1}{4}} \leq \frac{u^2}{8} = \frac{t^2(b-a)^2}{8}$$

② - What is the 'point' about which Taylor expansion is being carried out about?

(*) check this:

$$g(0) = -\gamma(0) + \log(1 - \gamma + \gamma e^0) = 0$$

$$g'(u) = -\gamma + \frac{\gamma e^u}{1 - \gamma + \gamma e^u} \qquad g'(0) = -\gamma + \frac{\gamma e^0}{1 - \gamma + \gamma e^0} = -\gamma + \gamma = 0$$

$$g''(u) = \frac{\gamma e^u(1 - \gamma + \gamma e^u) - (\gamma e^u)^2}{(1 - \gamma + \gamma e^u)^2} = \frac{\gamma e^u}{1 - \gamma + \gamma e^u}\left(1 - \frac{\gamma e^u}{1 - \gamma + \gamma e^u}\right) = s(1-s) \leq \frac{1}{4}$$

$s > 0$ ;

Hence
$$\mathbb{E}[e^{tX}] \le e^{g(u)} \le e^{\frac{t^2(b-a)^2}{8}}$$

Remark: Lemma 4 is known as Hoeffding's lemma
· Uses Taylor's theorem and Jensen's inequality
· It is an inequality that bounds the moment generating function of any bounded random variable (above and below)

## Lemma 5 - Chernoff's method

- Let X be a random variable. Then
$$P(X > \epsilon) \le \inf_{t > 0} e^{-t\epsilon} \mathbb{E}[e^{tX}]$$
where 'inf' can be understood as a 'min'

## Proof

for any t > 0
$$P(X > \epsilon) = P(e^X > e^\epsilon) = P(e^{tX} > e^{t\epsilon}) \le e^{-t\epsilon} \mathbb{E}[e^{tX}]$$
$$\underbrace{\qquad}_{(i)} \qquad \underbrace{\qquad}_{(ii)}$$

· since this is true for any t > 0, the result follows

(i) Raise/manipulate inequality within probability
(ii) introduce a variational parameter t

## Theorem 6 (Hoeffding's inequality)

- let $Y_1, \dots, Y_n$ be iid observations such that $\mathbb{E}[Y_i] = \mu$ and $a \le Y_i \le b$

- Then for any $\epsilon > 0$,
$$P(|\bar{X}_n - \mu| \ge \epsilon) \le 2e^{\frac{-2n\epsilon^2}{(b-a)^2}} \qquad (4)$$

## Corollary 7

If $X_1, \dots, X_n$ are independent with $P(a \le X_i \le b) = 1$ and common mean $\mu$ then with probability at least $1-\delta$

$$|\bar{X}_n - \mu| \le \sqrt{\frac{(b-a)^2}{2n} \log\left(\frac{2}{\delta}\right)} \qquad (5.)$$

# Proof of Hoeffding

- without loss of generality considerations:

  - Assume $\mu = 0$
  - And observe $P(|\bar{Y}_n| \geq \epsilon) = P(\bar{Y}_n \geq \epsilon) + P(\bar{Y}_n \leq -\epsilon)$

    $$= P(\bar{Y}_n \geq \epsilon) + P(-\bar{Y}_n \geq \epsilon)$$

---

- use Chernoff's method:-

$$P(\bar{Y}_n \geq \epsilon) = P\left(\frac{1}{n}\sum_{i=1}^{n} Y_i \geq \epsilon\right) = P\left(\sum_{i=1}^{n} Y_i \geq n\epsilon\right) \overset{(i)}{=} P\left(e^{\sum_{i=1}^{n} Y_i} \geq e^{n\epsilon}\right)$$

$$\overset{(ii)}{=} P\left(e^{t\sum_{i=1}^{n} Y_i} \geq e^{tn\epsilon}\right) \overset{(iii)}{\leq} e^{-tn\epsilon} \mathbb{E}\left(e^{t\sum_{i=1}^{n} Y_i}\right)$$

$$= e^{-tn\epsilon} \mathbb{E}\left[\prod_{i=1}^{n} e^{tY_i}\right] = e^{-tn\epsilon} \prod_{i=1}^{n} \mathbb{E}[e^{tY_i}]$$

$$= e^{-tn\epsilon} \left(\mathbb{E}[e^{tY_i}]\right)^n$$

(i) - $e(\cdot)$ monotone
(ii) - variational reparam.
(iii) - Markov's ineq.

- we bound $\mathbb{E}[e^{tY_i}]$ using lemma 4, Hoeffding's lemma :-

$$\mathbb{E}[e^{tY_i}] \leq e^{\frac{t^2(b-a)^2}{8}}$$

- So we have from above:-

$$P(\bar{Y}_n \geq \epsilon) \leq e^{-tn\epsilon} e^{\frac{-t^2 n(b-a)^2}{8}} \quad (*) \quad \text{- holds for any } t > 0$$

- Our variational tricks pays dividends as we can now minimise wrt $t$
- minimise RHS wrt $t$:
- select $t = \dfrac{4\epsilon}{(b-a)^2}$

- And we then have :-

$$P(|\bar{Y}_n| \geq \epsilon) \leq e^{\frac{-2n\epsilon^2}{(b-a)^2}}$$

- applying the same argument to $P(-\bar{Y}_n \geq \epsilon)$ yields the same result $\blacksquare$
- extending to case with $\mu$, define $Y_i = (X_i - \mu)$, prove in terms of $Y_i$, then s.b.

## 3. Bounded difference Inequality

- Hoeffding's inequality can be extended
- McDiarmid's inequality extends the general insight to more general functions $g(x_1, \ldots, x_n)$ of Hoeffding.
- supplementary

## 4. Bounds on expected values

### Theorem 11 - Cauchy-Schwartz Inequal.

- If $X$ and $Y$ have finite variance i.e. $\text{var}(|X|)$ and $\text{var}(|Y|) < \infty$ then

$$\mathbb{E}[|XY|] \leq \sqrt{\mathbb{E}[X^2]\mathbb{E}[Y^2]} \qquad (9.)$$

- Some additional exposition on convex functions (wasserman)
- $g(\cdot)$ is convex if for each $x$ and $y$ and each $x \in [0,1]$

$$g(\alpha x + (1-\alpha)y) \leq \alpha g(x) + (1-\alpha)g(y)$$

- If $g$ is twice differentiable and $g''(x) \geq 0 \ \forall x$ then $g$ is convex (calculus def.)
- If $g$ is convex then $g$ lies above any line that $g$ touches at that point (tangent line) (geometric)
- A function $g$ is concave if $-g$ is convex
- convex examples:- $g(x) = x^2$, $g(x) = e^x$
- concave examples: $g(x) = -x^2$  $g(x) = \log x$

(w): Is inverse of a convex function concave?

- CS inequality can be given a more explicit statistical context:-

$$\text{cov}^2(X,Y) \leq \sigma_X^2 \sigma_Y^2$$

(A) (W: why is this the case?)

### Theorem 12 - Jensen's Inequality

- If $g$ is convex; then

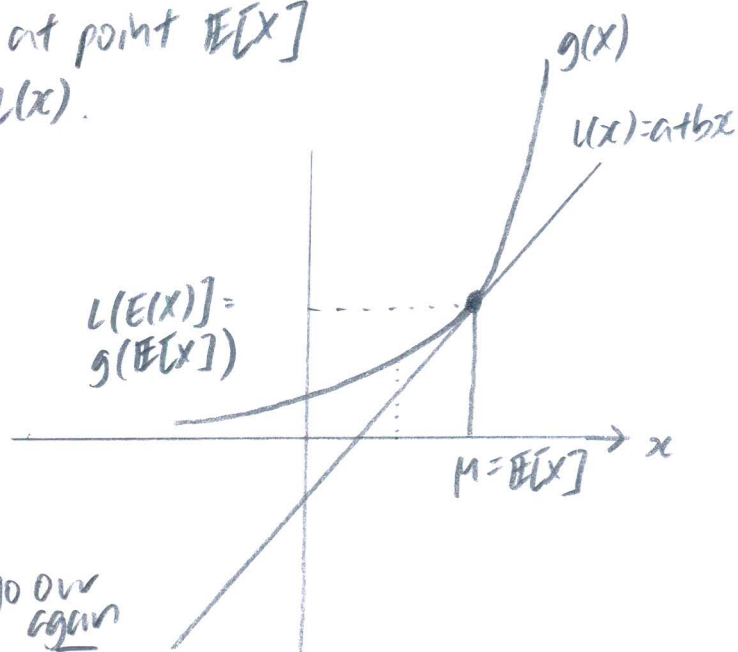$$\mathbb{E}[g(X)] \geq g(\mathbb{E}[X]) \qquad (10)$$

- If $g$ is concave, then

$$\mathbb{E}[g(X)] \leq g(\mathbb{E}[X]) \times \qquad (11)$$

**proof:**

- Let $L(x) = a + bx$ be a line tangent to $g(x)$ at point $\mathbb{E}[X]$
- Since $g$ is convex, it lies above the line $L(x)$.

$$\mathbb{E}[g(X)] \geq \mathbb{E}[L(X)] = \mathbb{E}[a + bX] = a + b\mathbb{E}[X]$$
$$= L(\mathbb{E}[X])$$
$$= g(\mathbb{E}[X])$$

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

(?) · For some reason the proof has become opaque to me again · ⓌⒶ · go over again



$g(x)$

$L(x) = a + bx$

$L(\mathbb{E}[X]) = g(\mathbb{E}[X])$

$\mu = \mathbb{E}[X]$

---

**Example 13**

- As $g(x) = x^2$ is convex, we have $\mathbb{E}[X^2] \geq (\mathbb{E}[X])^2$

i.e the 2nd moment is greater than the square of the 1st moment (mean)

- constitutes a proof of $\text{var}(x)$ being non-negative

**Example 14**

- Notes from lecture cover this fairly comprehensively

---

- Theorem 15 - skipped (spare time)

---

- consider bounding the maximum of a set of random variables

---

**Theorem 16**

- Let $X_1, ..., X_n$ be random variables
- Suppose there exists $\sigma > 0$ such that -

$$\mathbb{E}[e^{tX_i}] \leq e^{\frac{t^2 \sigma^2}{2}} \qquad \forall t$$

sub-Gaussian

⊛ Occurs for Normal, bounded r.v.s (thin-tailed r.v.s)

- Then $\mathbb{E}\left[\max_{1 \leq i \leq n} X_i\right] \leq \sigma \sqrt{2 \log n}$

## Proof

- Start with statement; apply trans.

$$\mathbb{E}\left[\max_{1 \leq i \leq n} X_i\right] \qquad (i)$$

- (W): I couldn't initially see how proof invoked Jensen's inequality
- Define the convex function $g(y) = e^{ty}$, apply to (i) and invoke J.I.

$$\exp\left\{t\,\mathbb{E}\left[\max_{1 \leq i \leq n} X_i\right]\right\} \leq \mathbb{E}\left[\exp\left\{t \max_{1 \leq i \leq n} X_i\right\}\right] \qquad \Bigg\downarrow (ii)$$

$$= \mathbb{E}\left[\max_{1 \leq i \leq n} \exp\{tX_i\}\right] \qquad \Bigg\downarrow (iii)$$

$$\leq \sum_{i=1}^{n} \mathbb{E}\left[\exp\{tX_i\}\right] \qquad \Bigg\downarrow (iv)$$

$$\leq n\, e^{\frac{t^2 \sigma^2}{2}}$$

- Apply logs :-

$$t\,\mathbb{E}\left[\max_{1 \leq i \leq n} X_i\right] \leq \log n + \frac{t^2 \sigma^2}{2}$$

$$\Rightarrow \mathbb{E}\left[\max_{1 \leq i \leq n} X_i\right] \leq \frac{\log n}{t} + \frac{t \sigma^2}{2}$$

- variational situation; set $t$ to minimise RHS.

i.e. set $t = \dfrac{\sqrt{2 \log n}}{\sigma}$

- Yielding $\mathbb{E}\left[\max_{1 \leq i \leq n} X_i\right] \leq \sigma \sqrt{2 \log n}$

2

---

(ii) Properties of max(·) function

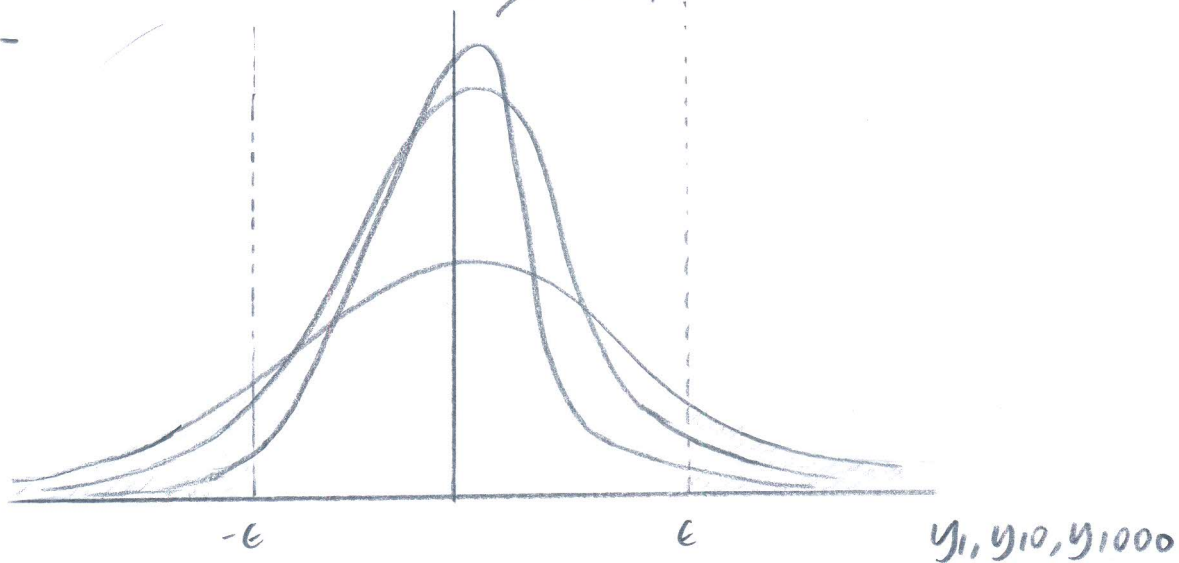(iii) $\max(x_1, \ldots, x_n) \leq \sum_{i=1}^{n} x_i$

(iv) by assumption

- Supplementary - review (Asymptotics)

- $Y_n = o_P(1) :-$ (convergence in prob.)          (sequence of r.v.s.)

$$P(|Y_n| > \epsilon) \xrightarrow{n \to \infty} 0 \quad \text{or} \quad \lim_{n \to \infty} P(|Y_n| > \epsilon) \to 0 \quad \forall \epsilon > 0^{(*)}$$

$$\to f_{Y_1}(y_1); \; F_{Y_{10}}(y_{10}); \; F_{Y_{1000}}(y_{1000})$$

- Diagram :-



$-\epsilon \qquad\qquad \epsilon \qquad\qquad y_1, y_{10}, y_{1000}$

· Note $P(|Y_n| > \epsilon)$ refers to <u>tail probabilities</u>

· A consequence of $Y_n = o_P(1)$ is that as $n \to \infty$; these tail probabilities (when you hold $\epsilon$ fixed at some arbitrary positive value) will approach 0.

· Note that for a fixed interval $[-\epsilon, \epsilon]$; the area under the PDF outside that interval, corresponding to tail probabilities, get smaller $[-\infty, -\epsilon]; [-\epsilon, \infty)$     and smaller, and approach 0.

$Y_n = O_P(1)$ (stochastic boundedness)

- Notes already cover this intuitively very well
- But I want to add a little more to capture some essential insight

- $Y_n = O_P(1)$
- If $\forall \epsilon > 0 \; \exists C_\epsilon : P(|Y_n| > C_\epsilon) \leq \epsilon \qquad \forall n > n_0$ (for finite $n_0, C_\epsilon$)

· Can we improve on this intuitively to better understand the differences in definition?

- $(*)$ - A subtlety in definition $\longrightarrow$ PTO

For $O_p(1)$ (convergence in probability); we require the statement to hold not only for one; but for any arbitrarily small $\epsilon$.

For $O_p(1)$ (stochastic boundedness); it suffices that there exists one arbitrarily large $C_\epsilon$ to satisfy the inequality; and $C_\epsilon$ is dependant on $\epsilon$.

· This yields the analysis/adversarial way of thinking for $O_p(1)$ as a pedagogical tool for proofs.

· (*) If you give me an $\epsilon > 0$; can I find an arbitrarily large $C_\epsilon$ (finite) such that statement holds for large $n$ greater than finite $n_0$?

(*) LW: If you give me $(1-\epsilon) = 0.9$; can I find an interval $[-C_\epsilon, C_\epsilon]$ to match that $\epsilon$ such that the interval traps 90% of the probability as $n \to \infty$?

(*) There are still some questions about this → place in overspill

- see stackexchange for proof examples; apply these insights to learn from them.

(*) Other interpretations (and formal def.); which can help in dusterdy formalism → see add. notes