

• Youtube 02/09/2016

• Lecture incomplete + cheat sheet

• Looks like lemma it has proved: go over what is missed.

•  $P(|\bar{Y}_n| > \epsilon)$



• Bound  $P(\bar{Y}_n > \epsilon)$

•  $P(\bar{Y}_n > \epsilon) = P(\sum Y_i \geq n\epsilon)$  3:02  $\rightarrow e^{(\cdot)}$  is monotone

$= P(e^{\sum Y_i} \geq e^{n\epsilon})$

$= P(e^{t \sum Y_i} \geq e^{t n \epsilon})$  - variational trick  
- extra parameter

- via Markov inequality

$\leq e^{-t n \epsilon} E[e^{t \sum_{i=1}^n Y_i}]$

$= e^{-t n \epsilon} \prod_i E[e^{t Y_i}]$

via independence

$= e^{-t n \epsilon} (E[e^{t Y_i}])^n$

$\leq e^{-t n \epsilon} e^{\frac{n t^2 (b-a)^2}{8}}$

(\*) - notes

$\rightarrow$  we use a bound on this

• Advantage of including  $t$  - variational trick

• Choose:  $t = \frac{4\epsilon}{(b-a)^2}$  (can be found with derivatives)

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

• Can apply to Bernoulli r.v.s. (coin flips) :-

$a=0, b=1$

$\rightarrow$  goes to 0 very quickly

(tight bound)

• Gives a very tight bound:- (Hoeffding inequality gives a very, very

• Intuitively: probability fraction of heads is going to depart from  $\theta$  (parameter) is very small.

W: If this doesn't go to 0 exponentially quickly, then 95% of stats and ML would not work.

Remarks

- for the case  $X_i$  with mean  $\mu$ ; then define  $Y_i = (X_i - \mu)$
- prove whole theorem in terms of  $Y_i$ ; then substitute  $\rightarrow$  (8:21)
- choice of  $t = \frac{4\epsilon}{(b-a)^2}$ ; (\*) holds for any  $t > 0$ , we want to make it as small as possible (16): check this (41)
- use deriv, solve; plot on wolfram.

### Section 3 (optional)

- extends Hoeffding to functions other than sums/averages of r.v.s.
- imagine arbitrary function of  $n$  r.v.s.
- same result (Hoeffding/variant) holds as long as function is smooth; in the sense that changing one of the function arguments by a small amount does not change the overall value of function much
- there for interest
- previously, we found bounds on probabilities (Markov, Chebyshev, Hoeffding)
- Bounds on expectations are slightly different
- theorem 1 - Cauchy-Schwarz inequality
- will be used a lot; also

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

W: no details (16) (17): (42) fill in

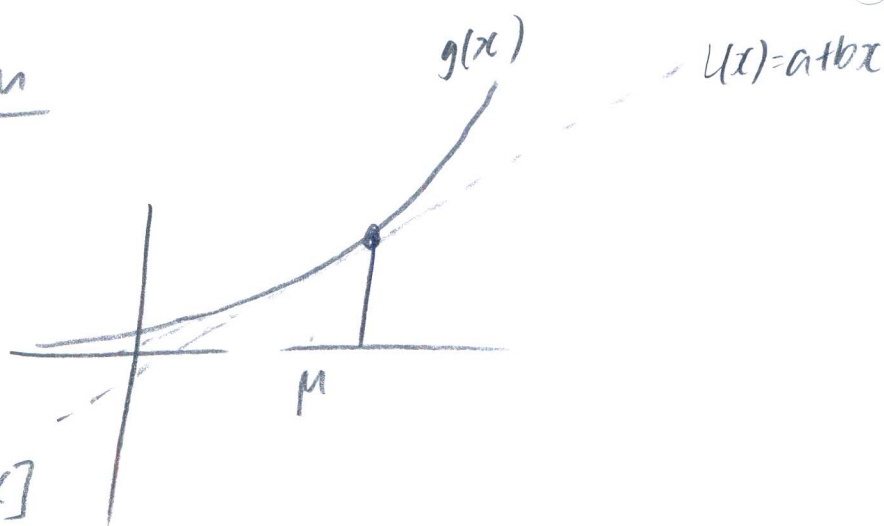
### Jensen's inequality

If  $g$  is convex then:

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

useful due to presence of  $\mathbb{E}[X]$

Proof: Okay -





## useful facts

- Require ideas of 'distance' between distributions (there are many)

## Kullback-Leibler

- this arises very 'naturally' from information theory

lw. not really a 'distance': not symmetric; but has distance-like prop.

- $D(p, p) = 0$ ;  $D(p, q) \geq 0$  (some desirable properties of distance)

- can use Jensen to prove KL-divergence:

- Any time you see  $\int p \dots$  or  $\sum p \dots \rightarrow$  think expected value

- note we only require specification of  $X \sim p$ ; the rest is in

$\mathbb{E} \left[ \log \frac{p(x)}{q(x)} \right]$  is just a function/trans.  $r(x)$  of  $X$ .

$$D(p, q) = \int p(x) \log \frac{p(x)}{q(x)} dx = \mathbb{E} \left[ \log \frac{p(x)}{q(x)} \right]$$

$$-D(p, q) = \mathbb{E} \left[ \log \frac{q(x)}{p(x)} \right] \leq \log \left[ \mathbb{E} \frac{q(x)}{p(x)} \right] = \log \int \frac{q(x)}{p(x)} p(x) dx = \log \int q(x) dx$$

via Jensen  
+ concavity of  $\log$

= 0  
(as  $\int q(x) dx = 1$  (density))

- lw: Switching  $\mathbb{E}$  and  $\log$ : common use of Jensen

- so  $D(p, q) \leq 0$ ; hence  $D(p, q) \geq 0$

Q13: lw does not cover ex. 13

- lw: skip Theorem 15 and proof

useful: Bounding expected value of max. of i.v. (recurs a lot)

- $X_1, \dots, X_n$  IID r.v.s.

- can impose an ordering  $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ ;  $X_{(n)} = \max \{X_1, \dots, X_n\}$   
by magnitude

- How to compute distri. of maximum? (conventionally find CDF)

lw: If we directly compute CDF  $\Pr(\max \{X_1, \dots, X_n\} < t) \Leftrightarrow \Pr(X_{(1)} < t, X_{(2)} < t, \dots, X_{(n)} < t)$

- use independence; 'multiply together'  $\rightarrow$  (2) (4)
- Get expression for CDF of max in terms of original CDF, get PDF, get  $\mathbb{E}$
- in principle this is feasible; but not in practice sometimes; especially if we don't know exact distri; but some properties on shape
- $\mathbb{E}[\max\{X_1, \dots, X_n\}]$  ("growing like  $n$ ") (28:10)  
 $\uparrow$   
 agmt.
- If we know distri is thin-tailed  $\rightarrow$  either from normality
- suppose we have a bound on MGF  $\leftarrow$  and/or Chernoff

$$\mathbb{E}[e^{tx_i}] \leq e^{\frac{t^2 \sigma^2}{2}} \quad \text{or } (2) (15) (*)$$

- $\sigma$  - not necessarily variance  $\downarrow$  (mean grows " $\log n$ ")

Theorem 16

• see notes  $\mathbb{E}[\max_{1 \leq i \leq n} X_i] \leq \sigma \sqrt{2 \log n}$

Proof:-

• start with  $\mathbb{E}[\max_{1 \leq i \leq n} X_i]$ ; apply transf:-

$$\begin{aligned} \exp\left\{t \mathbb{E}[\max_{1 \leq i \leq n} X_i]\right\} &\leq \mathbb{E}\left(\exp\left\{t \max_{1 \leq i \leq n} X_i\right\}\right) \quad \text{via Jensen (2) (16) - check} \\ &= \mathbb{E}\left(\max_{1 \leq i \leq n} \exp\{t X_i\}\right) \quad \text{(property of max(-); think)} \\ &\leq \sum_{i=1}^n \mathbb{E}[\exp\{t X_i\}] \\ &\leq n e^{t^2 \sigma^2 / 2} \end{aligned}$$

Apply logs:

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

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

• This is again a variational situation; set  $t$  to minimise

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

Giving:  $\mathbb{E} \left( \max_{1 \leq i \leq n} X_i \right) \leq \sigma \sqrt{2 \log n}$

WW: Recurring theme:

- if there exists a thin tail

- If we have an exponential bound on the MGF, we use that and Hoeffding's inequality style argument to say / make statements about

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

- why thin-tailed distribution?

⑦ WW: only thin-tailed distributions have properties where MGF have exponential bounds  $\rightarrow$  formally a Sub-Gaussian r.v.; informally thin tails

⑦ WW

⑦ Hoeffding's inequality is most important; and related  $\left\{ \begin{array}{l} \text{Bound MGF} \\ \text{Chernoff bound} \\ \text{'tacks'} \end{array} \right.$

- next lecture: Asymptotic notation