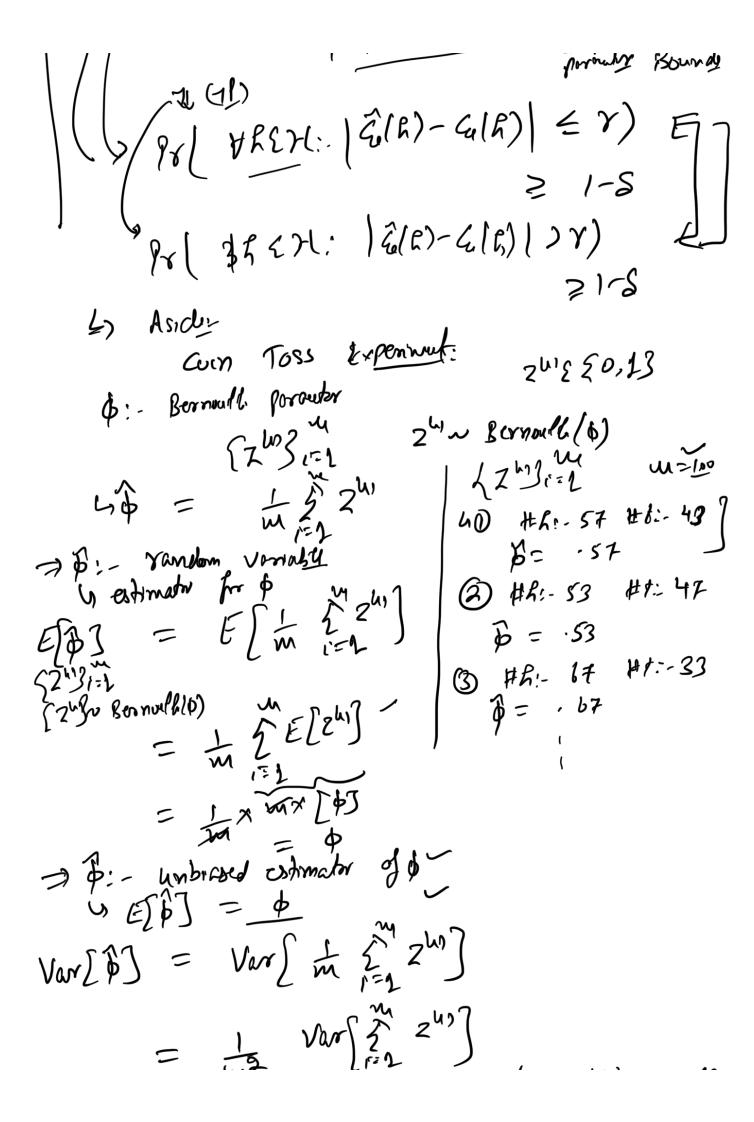
COLT+4 Nov 3rd, 2021

Lost Class:heaving Thing Shilying in yus so, 13 Stir Finite (classification) Du Hi initimite (classification) Du Hi initimite (classification) Du Hi initimite (vec-orm)

hood: h & ye hypothene space (vec-orm)

Hi f(nu)) = yer (numeroring traing con) $\frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle E(R) \rangle}{\langle E(R) \rangle} \times \frac{\langle E(R) \rangle}{\langle E(R) \rangle} = \frac{\langle$ Probably Approximately breat Bound Pr (38624: - |4/6) -4/6) > V) -0 (=) (Pr (3 RE) (! | G/R) - G/R)) V) - () - () - () - ())) - (8 = 1711.2c-272021 = 5



(: Z"/s one wid) = 1 2 Var [24] 7~ Bornou 1660) 4 Vov[2] = E[29-[6[2])² = 6(1-4) Van [6] = 1 x 24x \$(1-6) Ly Goes down as Hofexamples is involved Hoeffoling inequality [therrist sounds. $P_{Y}(|\phi-\hat{\phi}|>\gamma) \leq 2e^{-2\gamma^{2}m}$ tru prouter Li PAL bounds for superused hearing seding (nh, yh) ~ Dist 5 247, 3473 col 1 { flam) + yui} is orner of 2th example alt) = 1 2 12 tooks + yin 2 Zun Bernoullife (ny, yar) ~ Dist 25 x (my) + y y ? = 2 h) 1 { zln) + y3}. - Gonwaledmo = ETG (2, 1, 2) m) ~ DISZ

 $\frac{1}{2} \left(\frac{1}{2} \right) = \frac{1}{2} \left(\frac{1}{2} \right) = \frac{1}{2} \left(\frac{1}{2} \right) \left(\frac{1}{2} \right)$ > 4/h) = es timetor destribution = 1 Charnoff bound Elt) $P(\left(\begin{array}{c} \left(\frac{\pi}{6}\right) - 4(\pi) \\ \left(\frac{\pi}{6}\right) - 4(\pi) \\ \left(\frac{\pi}{6}\right) \end{array}\right) > \gamma) \leq 2e^{-2\gamma^2}$ 1/ () Elle) - 4/h) > r) = ze-2rm 1 2(he) - 4 (he) | > ~ AZUAZ -- AK) = ZP(Ae) Prol 3 her | Ga (he) - Ga (he) | >r)

notine commence sound

notine commence sound