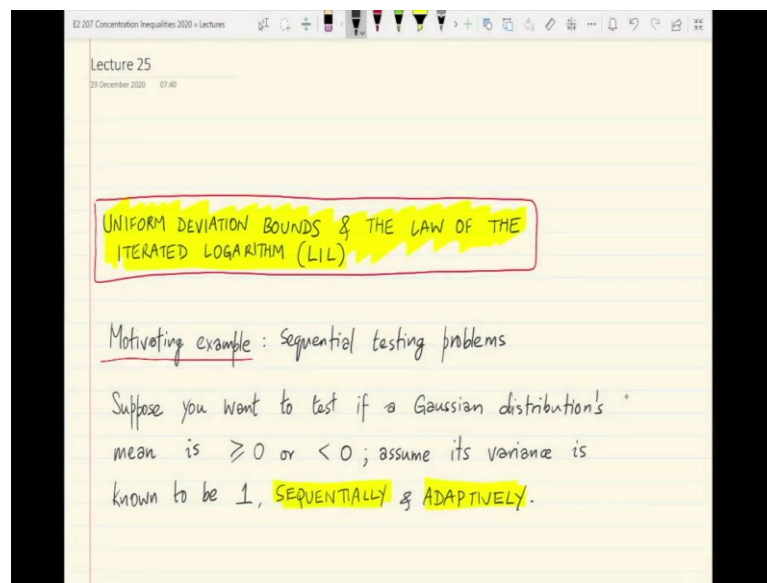


Concentration Inequalities
Prof. Aditya Gopalan
Prof. Himanshu Tyagi
Department of Electrical Communication Engineering
Indian Institute of Science, Bengaluru

Lecture - 25
Uniform Deviation Bounds and the Law of the Iterated Logarithm (LIL)

(Refer Slide Time: 00:21)



Hi all. In the next few lectures we will be looking at Concentration Inequalities for sequential processes. Also you can think of these as means to control the fluctuations of entire parts of sequences of random variables. And in this context an important topic that we will present is called Uniform Deviation Bounds and we will also explore the connection of these uniform deviation bounds to what is called the Law of Iterated Logarithm which is a well known result in probability theory.

So, let us start with a simple motivating example this is the example of sequential testing ok. So, what is a typical sequential testing problem? Let us say we are thinking of a sequential testing problem about the mean of a Gaussian distribution. So, let us say you want to test if an unknown Gaussian distributions mean is either positive or negative to the right side or to the left side of 0 moreover let us assume that we know its variance.

So, let us say its variance is unit variance and we would like to test as soon as possible sequentially and adaptively by drawing iid samples from this unknown Gaussian

distribution whether each of these hypothesis is true that is the real the true mean is either positive or negative. So, more formally speaking here is what you as a learner or a testing algorithm would do in order to solve this problem?

(Refer Slide Time: 01:57)

At each time $t = 1, 2, 3, \dots$ *based on past samples*

- If you decide to stop:
 - Output $A \equiv \begin{cases} 1 & \text{if you think Mean} \geq 0 \\ 0 & \text{" " " " Mean} < 0 \end{cases}$
 - Your "answer" or "guess"* ←
 - Break
- Else
 - Get a new indep't sample $X_t \sim \mathcal{N}(\mu, 1)$.

* Define $\tau :=$ the (random) time @ which the procedure stops & outputs an answer A .

We'd like the misclassification probability to be $\leq \delta$, i.e.,

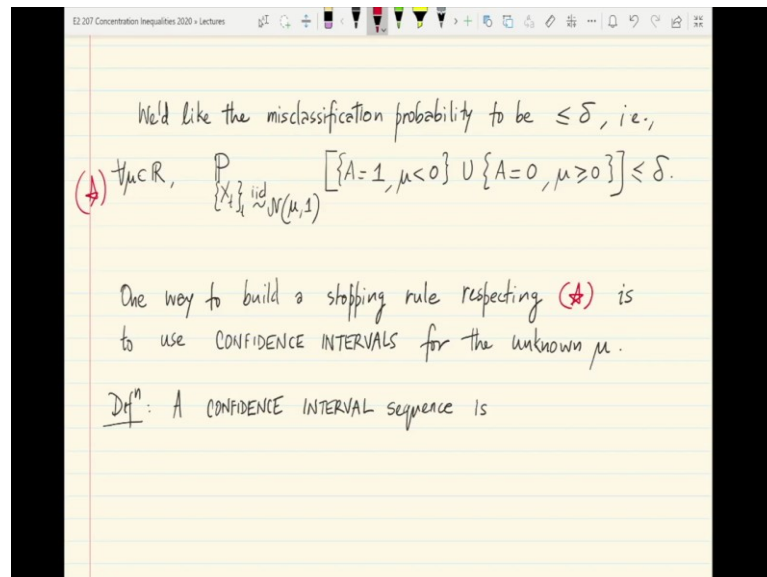
So, this is the protocol at each time index by t . So, you as a learner or tester you can decide to stop at this time t . So, if you decide to stop you can either stop or take one more sample. If you decide to stop and this is based on past samples based on everything you have seen in the past you can decide to stop past samples then you are also required to guess or decide which of these situations is actually active whether the mean is positive or negative.

So, you output random variable A , which is your guess and let us say that by convention you set A to be 1, if you believe or if you think the mean of the Gaussian is non negative and the output 0 if you think the mean is negative ok. And after your output this you break you break this loop else you request one more iid sample a fresh iid sample from the unknown distribution. So, you get to observe a new independent sample.

Let us see we call that X_t drawn from normal distribution with unknown mean μ and known variance 1 ok. So, this is the decision making loop. A essentially represents your answer or guess. So, typically a good sequential procedure is some is a procedure that takes as few samples X_t as possible and outputs A which is correct answer with high probability ok.

So, towards this let us formalize the requirement from such a sequential hypothesis testing procedure. Let us define tau to be the time or the random time at which the procedure stops and also outputs an answer A. So, it is the first time that the stop condition is invoked by the algorithm.

(Refer Slide Time: 05:00)



And a basic performance requirement from any algorithm is this correctness requirement or in other words we would like the misclassification probability of the answer written by the algorithm A to be always less than equal to δ . So, here is a mathematical way of writing it. For any value of the true mean of the distribution that is generating samples if you run this algorithm then all the samples provided to it are drawn from a normal distribution with mean μ and variance 1.

So, you run this algorithm and wait for it to stop and you look at the answer A the random answer A that it outputs. So, if A equal to 1; that means, the algorithm stops and thinks that the mean is larger than 0, but the actual mean is less than 0. So, this event or the event where A is returned to be 0, but the actual mean is greater than equal to 0. So, both these are error events. Only one of these can occur depending on the actual value of μ .

So, in any case the probability of such an event is required to be bounded by at most a given number δ say δ equal to 0.1 or 0.01 ok. So, the decision the design choice here is when to decide to stop and if the algorithm has decided to stop or break this loop, what is

the hypothesis to output based on all the data that has been seen so far or all the samples that have been collected so far right. So, a common way to build correct stopping rules which essentially satisfy this requirement.

So, let us; so, let us call this correctness requirement star. So, one way to build a stopping rule that respects star given a target parameter δ is to use what is called confidence intervals for the unknown parameter that governs the distribution of the samples. So, what do we mean by confidence intervals?

(Refer Slide Time: 07:18)

to use CONFIDENCE INTERVALS for the unknown μ .

Defⁿ: A CONFIDENCE INTERVAL sequence is a family of functions, one for each $n \geq 1$, that maps $X_1^n \mapsto I_n = [L_n, U_n] \subseteq \mathbb{R}$, such that, for any stopping time τ (depending on $\{X_i\}_i$) & $\forall \mu \in \mathbb{R}$,

$$\mathbb{P}_\mu [I_\tau \ni \mu] \geq 1 - \delta.$$

If we have a $(1-\delta)$ confidence interval construction, then a good sequential hypothesis testing algorithm is:

* Stop at $\tau \equiv$ first time t when $0 \notin I_t$

So, here is a definition of here is one possible definition of confidence interval sequence. So, a confidence interval sequence is nothing but a family of functions, one for each value of sample index one for each let us say n sample number of samples n that maps. So, what does it take as input and what does it give as output? It maps n samples seen so far we denote it by X_1 through X_n and use X superscript n subscript 1 and outputs an interval I_n ok.

So, I_n is a real interval of real numbers with lower index called L_n and upper index called u_n . So, we will assume in general that it is a closed interval of this form subset of \mathbb{R} such that the following property is satisfied, for any stopping time or any stopping rule. So, we already defined τ as a stopping time τ . So, recall that τ can depend on the data on the X_i sequence.

And for all values of the distribution parameter that governs the probability distribution of the samples we must have that the probability ok. So, P_μ is shorthand notation for the same thing that is probability where all the samples each successive sample is generated iid from $N(\mu, 1)$ the probability that the interval I at time τ actually contains the real parameter μ is at least $1 - \delta$ ok.

So, this is basically a coverage property. It says that no matter when you decide to stop or no matter under what kind of stopping rule when the algorithm decides to stop, I_τ of τ the interval I_τ which is $[L_\tau, U_\tau]$ the set of all numbers which are at least L_τ and at most U_τ must contain μ with significant probability ok. So, note that a confidence interval sequence is an entire family of intervals.

So, for any given n , it is desired that there is a interval I_n such that whenever you happen to stop at time n , I_n must contain μ with high probability. So, how do confidence interval constructions help design good sequential tests? Here is the connection.

(Refer Slide Time: 10:39)

τ (depending on $\{X_i\}_i$) $\delta \forall \mu \in \mathbb{R}$,

$$P_\mu[I_\tau \ni \mu] \geq 1 - \delta$$

If we have a $(1-\delta)$ confidence interval construction, then a good sequential hypothesis testing algorithm is: $[L_t, U_t]$

- * Stop at $\tau \equiv$ first time t when $0 \notin I_t$
- * Output $A = \begin{cases} 1, & \text{if Lower endpoint of } I_\tau > 0 \\ 0, & \text{"Upper" " " } < 0. \end{cases}$

Why?

So, if we have a valid confidence interval construction with probability of failure no more than δ as above then here is a good sequential testing algorithm for the Gaussian mean testing problem ok. So, here is what such an algorithm would do. The algorithm would stop at the first time t when 0 falls out of the interval I_t that is when the interval I

it does not contain 0. So, recall that the interval I_t is at time t comprised of a lower endpoint L_t and an upper endpoint u_t .

So, if $L_t \leq 0$ or $u_t \geq 0$; that means, if either L_t is larger than 0 or u_t is less than 0 then you stop. And you output the natural decision; that means, if the lower endpoint of I_τ when you stop at time τ L_τ if its larger than 0 then you guess that the mean of the Gaussian distribution is larger than 0. If it is if the upper endpoint u_τ happens to be less than 0 then you would naturally guess that the mean of the Gaussian you are dealing with is less than 0 ok.

Let us call this a rule an algorithm this gives you an entire algorithm for performing sequential testing and stopping when required with confidence ok. So, why is why does this work? So, what is the proof that this kind of procedure actually has a misclassification probability which is bounded at most by δ ?

(Refer Slide Time: 12:21)

Why? If $\mu < 0$, $P_\mu[A=1] = P_\mu[L_\tau > 0]$

$$= P_\mu[L_\tau > 0, \mu \in I_\tau] + P_\mu[L_\tau > 0, \mu \notin I_\tau]$$

$$\leq P_\mu[\mu \notin I_\tau] \leq \delta, \text{ so } (*) \text{ is met.}$$

* How to construct confidence intervals?

Well, it is easy to see this. Let us assume for the moment that if the real mean of the Gaussian generating the samples is actually negative let us evaluate the probability that the answer returned at time τ according to the rule above when you stop is equal to 1. So, this is the error event when μ is actually negative and the answer output is actually that the mean is positive.

So, in other words this is just $P(\mu \text{ that the lower endpoint of the interval } I_\tau \text{ that is } L_\tau \text{ is larger than } 0)$. So, we can break this sum into two parts by the law of total probability. So, this is $P(\mu \text{ } L_\tau \text{ greater than } 0 \text{ and } \mu \text{ belongs to } I_\tau) - \text{probability of the same event intersected with } \mu \text{ naught belonging to } I_\tau$.

Now, the probability of the first event is actually 0 because if this event were to occur then it means that μ is at least L_τ , but L_τ is at least 0. So, μ is greater than 0 which is not possible because μ was assumed to be less than 0 and so, we are left with only the second term for which an upper bound can be obtained by just dropping the first event inside.

So, this is just upper bounded by the probability that μ does not belong to I_τ and we know by the confidence interval property of the sequence of intervals I_1, I_2, I_3 and so on that this event occurs with probability at most τ ok. So, the condition star which is the correctness condition. So, we call the star condition, the misclassification probability condition is met ok.

So, what we have seen is that good confidence interval constructions or valid confidence interval constructions when coupled with the natural stopping rules yield non trivial hypothesis sequential hypothesis tests ok. So, how do you go about constructing confidence intervals? So, this is what we will expose the connection between confidence intervals and uniform deviations of stochastic processes.

(Refer Slide Time: 14:59)

* How to construct confidence intervals?

1) Fix $n \in \mathbb{N}$. Since $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, 1)$, by Chernoff, $P_\mu \left[\left| \underbrace{\frac{1}{n} \sum_{i=1}^n X_i}_{\bar{X}_n} - \mu \right| \geq \varepsilon \right] \leq 2e^{-\frac{n\varepsilon^2}{2}}$

$\Rightarrow P_\mu \left[\mu \in \left[\underbrace{\bar{X}_n}_{L_n} - \underbrace{\sqrt{\frac{2 \log(2/\delta)}{n}}}_{U_n}, \bar{X}_n + \sqrt{\frac{2 \log(2/\delta)}{n}} \right] \right] \geq 1 - \delta.$

So, let us perform the following thought experiment. So, let us take a set of these random variables. So, let us take the first n successive rewards generated by the distribution θ_k or the unknown Gaussian distribution. So, since, so, fix. So, let us say let us first fix an integer N , θ_k a non negative integer N , since the first N samples θ_k . So, just think of the first N samples returned by the of the Gaussian distribution out of the infinitely many samples that are possibly generated.

This has no relation to what samples were actually consumed by an algorithm. So, we are just thinking of the first N samples for a fixed small n . So, since X_1 through X_n are drawn assumed to be independent and identically distributed as a normal distribution with mean μ and variance 1 by Chernoff by the Chernoff bound we already have that $P(\mu \text{ of let us say the sample mean } \frac{1}{n} \sum_{i=1}^n X_i - \mu \text{ exceeding } \epsilon \text{ a level } \epsilon \text{ is bounded by } 2e^{-\frac{n\epsilon^2}{2}} \text{ by 2.}$

So, this is a standard Chernoff bound for a Gaussian θ_k . Let us for shorthand denote the sample mean of the first n samples as \bar{X}_n θ_k . So, this means that; so, this is the say this in fact, means that if I set a particular value of ϵ to make the right hand side equal to δ . So, this is what we get the probability that μ belongs to this particular interval $\bar{X}_n \pm \sqrt{2 \log 2 \text{ by } \delta \text{ divided by } n}$ \bar{X}_n - the same deviation $\sqrt{2 \log 2 \text{ by } \delta \text{ over } n}$ this probability is at least $1 - \delta$ θ_k .

So, this holds for any fixed n θ_k . So, if there was a stopping rule that decided always to stop at time small n θ_k . So, never stop before that or after that. So, let us say there was a stopping rule that always stopped trivially at time small n then this interval I_n with left endpoint as this number L_n and right endpoint as u_n θ_k . So, this interval would be a valid confidence interval at that specific time small n θ_k . So, this is a $1 - \delta$ confidence interval.

(Refer Slide Time: 18:10)

So, a confidence interval with level $(1-\delta)$ for μ after having observed exactly n samples is

$$\left[\bar{X}_n - \sqrt{\frac{2 \log(2/\delta)}{n}}, \bar{X}_n + \sqrt{\frac{2 \log(2/\delta)}{n}} \right]$$

BUT: At a random, possibly DATA-DEPENDENT stopping time τ , X_1, \dots, X_τ need not be iid (in general) !!

[Chernoff requires τ to be NON-RANDOM]

2) Chernoff + Union bound

So, let us record this by saying that a confidence interval with level $1 - \delta$ for μ after having observed exactly n samples where n is some fixed non random integer n samples is this interval to $\bar{X}_n - \sqrt{2 \log 2 \text{ by } \delta \text{ divided by } n}$ $\bar{X}_n + \sqrt{2 \log 2 \pi \delta \text{ divided by } n}$ ok.

So, the key thing here is that the n samples are fixed ok. So, this is set sort of a confidence interval for a fixed or a batch. So, this is like a batch of n iid samples ok. So however, if you have a stopping rule that actually stops at a random time depending on the data, so, at a random possibly data dependent that means, depending on the previous history of samples observed so far.

So, let us say someone has designed a stopping time τ that is sort of non trivial ok. So, its maybe complicated. Maybe you stop at the first time when the certain property of the sequence of samples seen so far has been satisfied. So, that could lead to various distribution of the stopping time τ not necessarily always at a fixed number of samples n ok.

Its not hard to argue that the set of samples that you have seen until stopping X_1, X_2 and so on up to the last time a sample was taken; that means, X_1 through X_τ the τ samples. In fact, these are no longer iid in general ok. So, these need not be iid in general ok. So, this is a probably a surprising fact the first time you encounter it, but on the other hand its not hard to come up with some counter examples to X_1 through X_τ being iid.

So, basically the upshot is that in the previous setting we considered n to be a fixed batch of samples.

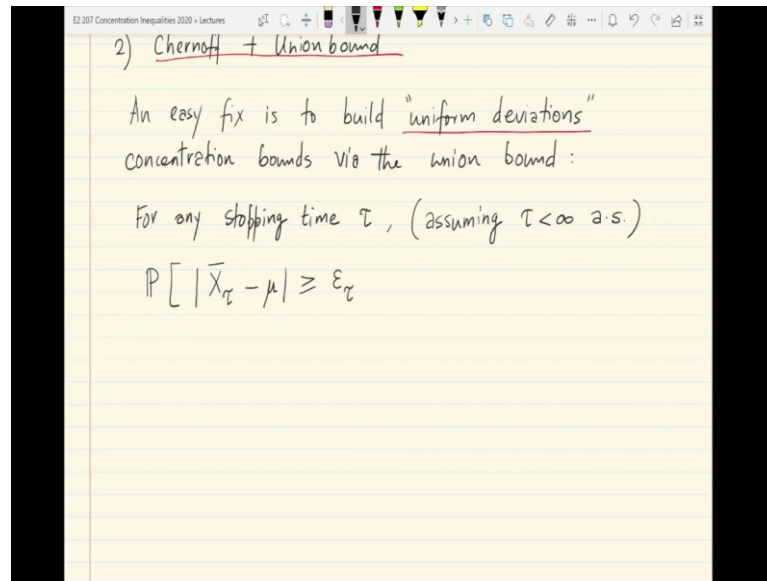
So, this is equivalent to saying that there is a stopping rule for which τ is exactly equal to n always; it always stops after taking n samples no less no more. And so, in that case naturally X_1 to X_n by definition are going to be iid, but if the stopping rule actually depends on the data seen in the past then the iid property is can be completely violated ok.

So, just as an example you can think of a rule which is quite strange in the sense that it says you know let us wait, let us stop when you see a sample larger than some value let us say 10 for the first time ok. So, it obviously means that the sequence of samples you have seen so far has the property that the last sample X_τ is going to be greater than equal to 10 and all the previous samples are by definition going to be less than 10 ok.

So, this is not an iid distribution over samples X_1 through X_τ ok, but. So, this sort of violates the ability to apply Chernoff because as we know Chernoff requires τ to be non random. So, you cannot ok. So, you could not have applied Chernoff in the setting earlier when n was when n was random. So, think of a random number of samples being taken which can actually depend on the samples themselves.

So, if n was actually a random variable you would not be able to apply Chernoff for obvious reasons because Chernoff requires number 1 all the random variables to be independent. And number 2 there is this n that appears on the right hand side. So, you cannot have n as a random variable ok. So, this is a sort of subtle, but important issue that one has to worry about when defining concentration events for sequential procedures or sequentially sampled processes.

(Refer Slide Time: 23:14)



So, the obvious way one obvious way to get around this problem is to use a union bounding argument along with Chernoff ok. So, the idea here is that you can build uniform deviations concentration bounds by using the basic union bound idea from probability ok. So, what do we mean by this?

So, suppose we want to construct confidence you want to construct a family of confidence intervals one for each time such that no matter what the value of tau is the interval at time at the stop time I tau contains mu with high probability must hold. So, here is what you do.

At any for any stopping time induced by a stopping rule and we define mathematically the stopping time to be tau ok. We will make the small technical assumption that tau is almost surely finite, assuming tau is less than infinity almost surely. So, as long as tau is less than infinity almost surely one can do the following one can write the following. So, the probability.

So, I am omitting the subscript nu which is because that is understood the probability that the absolute value difference between X bar tau. So, imagine that we have stopped at time tau according to some data dependent rule and X bar tau is the mean of the first tau samples that you have taken before stopping. So, the probability that X bar tau - real mean mu exceeds let us say a number ϵ_τ that we will design later.

So, let us think about the probability of this event ok. So, we know that tau is less than infinity almost surely.

(Refer Slide Time: 25:14)

$$\begin{aligned}
 & \mathbb{P} \left[|\bar{X}_\tau - \mu| \geq \epsilon_\tau \right] \\
 & \leq \mathbb{P} \left[\exists i \in \mathbb{N} : |\bar{X}_i - \mu| \geq \epsilon_i \right] \\
 & \leq \sum_{i=1}^{\infty} \mathbb{P} \left[|\bar{X}_i - \mu| \geq \epsilon_i \right].
 \end{aligned}$$

So, if the i^{th} term above is $\leq \frac{\delta}{i^2}$,
then the sum is $\leq c \cdot \delta$.

So, what this event implies is that there must exist a value of i where i is any positive is some positive integer such that the absolute value difference between \bar{X}_i the first i samples average - μ must have exceeded ϵ_i ok, this is just by basic inclusion. So, there is no way that the first event would have happened if the if not for the second event ok.

So, at some point in time tau if your sample mean of tau samples exceeds is further than further from μ then by an amount ϵ_τ that it means that there must exist some integer i at which $\bar{X}_i - \mu$ is greater than equal to ϵ_i . And by the union bound you can bound this probability by summing over all possible values of this time tau ok of the probability that $\bar{X}_i - \mu$ exceeds, ϵ_i ok ϵ_i something that we are yet to derive.

So, let us try to arrange things. So, that this final sum that we have here, the sum from i equal to 1 to infinity is no more than δ . Let us say we are given a target value of δ for the for coming up with a confidence interval construction. So, if the i^{th} term in the sum above is let us say less than equal to let us say let us just say δ over i square.

So, δ is the number the probability the violation probability number we have been given and let us probably try to find ϵ_i such that each successive term is bounded as δ over i square then by summing all these 1 over i square is a summable sequence sum is for

instance less than some constant times δ ok. The constant is whatever you get when you sum 1 over i square ok. So, what is this?

(Refer Slide Time: 27:43)

then the sum is $\leq c \cdot \delta$.

$$\mathbb{P}[|\bar{X}_i - \mu| \geq \epsilon_i] \leq \frac{\delta}{i^2}$$

\Uparrow (by Chernoff)

$$2e^{-\frac{i\epsilon_i^2}{2}} \leq \frac{\delta}{i^2}$$

\Downarrow

$$\epsilon_i = \sqrt{\frac{2}{i} \log\left(\frac{2i^2}{\delta}\right)}$$

$I_\tau = \bar{X}_\tau \pm \sqrt{\frac{2}{\tau} \log\left(\frac{2\tau^2}{\delta}\right)}$

So, what is our target? We want to make probability $\bar{X}_i - \mu$ greater than equal to ϵ_i less than equal to δ by i square ok. Now, we have fixed. So, i is a fixed integer here ok, there is nothing random about this integer i . So, \bar{X}_i is simply now the sample mean of i the first iid samples from a Gaussian with a mean μ and so, here is where you can actually use Chernoff.

So, by Chernoff an upper bound, so, this is implied if the Chernoff upper bound which is $2e^{-i\epsilon_i^2/2}$ can be made less than δ by i square and this is just the same as saying as setting ϵ_i to be equal to $\sqrt{2 \log 2 i^2 / \delta}$ ok. So, right.

So, what we have shown by this union bounding argument along with Chernoff is that even for a random no matter what the value of the stopping time τ you are assured that the interval I_τ . So, you can set I_τ as \bar{X}_τ - the lower end point being the - and upper endpoint being the $+\sqrt{2 \log 2 \tau^2 / \delta}$ ok.

So, you have an interval at stopping and then you can basically make a decision based on whether 0 is in this interval or not ok. So, the usual theory carries forward. So, let us think about this union bound result we have shown from a slightly different viewpoint.

(Refer Slide Time: 29:45)

$$c_i = \sqrt{\frac{2}{i} \log\left(\frac{2i}{\delta}\right)}.$$

In other words, we've shown

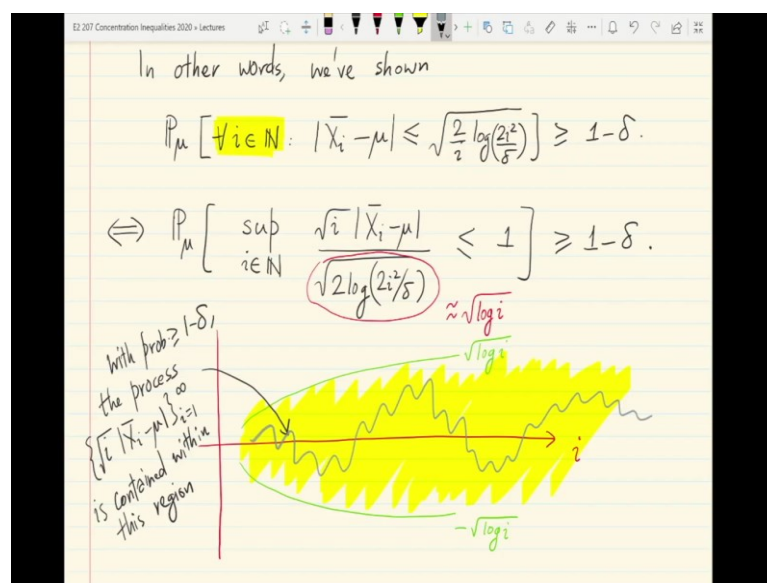
$$\mathbb{P}_\mu \left[\forall i \in \mathbb{N}: |\bar{X}_i - \mu| \leq \sqrt{\frac{2}{i} \log\left(\frac{2i}{\delta}\right)} \right] \geq 1 - \delta.$$
$$\Leftrightarrow \mathbb{P}_\mu \left[\sup_{i \in \mathbb{N}} \frac{\sqrt{i} |\bar{X}_i - \mu|}{\sqrt{2 \log(2i/\delta)}} \leq 1 \right] \geq 1 - \delta.$$

So, in other words what have we shown, what have we shown? We have shown that the probability that for any positive integer i , $\bar{X}_i - \mu$ less than so. In fact, for all i less than equal to n for all i natural numbers the probability of the deviation of \bar{X}_i from μ being less than $\sqrt{2 \log(2i/\delta)}$ is greater than equal to $1 - \delta$ ok.

So, uniformly over time, so, for all i natural numbers is just saying uniformly over time the sample mean \bar{X}_i ; $\bar{X}_i - \mu$ always stays within this fixed i dependent number on the right ok. Put slightly differently this is also equivalent to saying that the supremum taken over all i positive integers of; so, if you divide the left hand side by the right hand side and it is at most 1.

And it just says that the supremum of \sqrt{i} times $\bar{X}_i - \mu$ divided by $\sqrt{2 \log(2i/\delta)}$ is less than equal to 1. This entire event has a significantly large probability of $1 - \delta$ of at least $1 - \delta$ ok.

(Refer Slide Time: 31:44)



So, just to pictorially illustrate; so, this is; so, let us illustrate this statement with something that conveys the nature of fluctuations of the stochastic process \bar{X}_i ok. So, for instance here is one way of representing this result. So, on the X axis we have i . Let us so in fact, one remark here is that the denominator here. So, imagine δ to be fixed δ is a fixed number. So, order wise this is basically nothing but $\sqrt{\log i}$ ok. So, let us plot $\sqrt{\log i}$ on one hand.

So, if i start from 1 then this is how $\sqrt{\log i}$ looks and on the negative side this is how $-\sqrt{\log i}$ looks ok and this is sort of the region in between. So, let us think about the paths of the entire stochastic process indexed by i which is \sqrt{i} times $\bar{X}_i - \mu$ ok.

So, all we are saying is that with probability at least $1 - \delta$ the stochastic process given by \sqrt{i} times $\bar{X}_i - \mu$ ok over all integers i is contained within this shaded yellow region ok. So, there is this tube if you wish. So, this is sort of a tube that is bounded by the functions roughly $\sqrt{\log i}$ and $-\sqrt{\log i}$.

And what we have established by this very basic union bounding argument is that if you basically take a scaled version of the difference between the sample mean and the true mean scaled by basically \sqrt{i} the number of samples and plot that trajectory ok then a typical trajectory is going to sort of wander always within this tube ok.

So, that is what we have shown ok. So, this is in some sense to be regarded as a statement about the trajectory of the entire stochastic process \sqrt{i} into $\bar{X}_i - \mu$ for i equal to 1 to infinity ok. So, it is an infinitely long infinite time stochastic process ok. So, that is what we have shown.

So, this raises the question of how much we can actually improve things ok. So, can we actually try to if you are given you know a $1 - \delta$ high probability target what is sort of the narrowest kind of green tubes that you can draw to bracket the entire trajectory of this stochastic process with probability at least $1 - \delta$ ok?

(Refer Slide Time: 35:09)

3) CONTROLLING (TIME) UNIFORM DEVIATIONS

Question: We saw that for $\{X_i\}_{i \geq 1} \stackrel{iid}{\sim} \mathcal{N}(0,1)$,

$$P \left[\sup_n \frac{|\sum_{i=1}^n X_i|}{\sqrt{2n \log \left(\frac{2n^2}{\delta} \right)}} \leq 1 \right] \geq 1 - \delta.$$

By how much can the denominator be improved (increased) for uniform ($n \in \mathbb{N}$) deviations to be under control?

And that is the that brings us finally, to the subject of controlling uniform deviations or time uniform deviations of a stochastic process right from time 1 to time infinity. So, here is a natural question we saw that if for instance in the μ equal to 0 case, so, when X_i are all drawn X_i is a sequence of iid infinite iid random variables drawn from the standard normal distribution then we saw that the probability that the largest possible value of this ratio which is absolute value of S_n divided by $\sqrt{2n \log 2n}$ square.

S_n is nothing but the sum of the first n random variables X_i ok. So, you have just represented \bar{X}_i as S_i divided by i ok. So, what we saw is that S_n ok you can think of S_n as a random walk. The magnitude of S_n if you scale it by $\sqrt{2n \log 2n}$ square by δ that quantity is never going to hit 1, hit cross the level 1 with probability at most $1 - \delta$ ok.

So, basically it says that if you draw in the S_n sense if you draw 2 tubes of - - $\sqrt{n \log n}$ then S_n is not going to wander outside this tube ever unless you are in an event of probability at most δ ok. So, the natural question here is can you do something about this denominator by how much can the denominator be improved or increased ok?

So, if you want tighter control of if you see tighter control of S_n let us say then it means that you would want to try to raise the level of the denominator (Refer Time: 37:03) as a function of n and δ . So, by how much can the denominator be improved or increased for uniform deviations to be under control ok. By uniform deviations we basically mean the supremum over all positive integers ok.

(Refer Slide Time: 37:26)

(increased) for uniform ($n \in \mathbb{N}$) deviations to be under control?

We will show: for $\{X_i\}_i \text{ iid } \mathcal{N}(0,1)$

$$\mathbb{P} \left[\sup_{n \in \mathbb{N}} \frac{S_n}{O\left(\sqrt{\frac{n \log \log n}{\delta}}\right)} \leq 1 \right] \geq 1 - \delta,$$

an exponential improvement over n^α .

Is this best possible?

So, this is what we will show to this effect. It turns out that there is sort of in some sense a large room for improvement compared to what we have got so far. So, we will show the following result in this lecture. So, we will show. In fact, that the usual situation let us say without loss of generality that all the X 's are drawn iid with the standard normal distribution.

Then the probability of the worst case ratio of between S_n and a quantity that is roughly order \sqrt{n} the \sqrt{n} is inevitable, but what we will be able to do is we will be able to replace the polynomial n sitting inside the log by a logarithmic n . So, we can actually get $\log \log n$ as a function of n exceeding a being less than 1 to be at most to be at least $1 - \delta$ ok.

So, this is the result that we will be able to show in contrast to the previous result. So, instead of $\log n$ in the previous result, so, there is a polynomial n sitting inside the log which is roughly equal into $\log n$. We will be able to actually get $\log \log n$ deviation as a multiplier for the \sqrt{n} ok. And this is actually if you think about it this $\log n$ compared to the n square on top is an exponential improvement over n to the alpha or polygon.

So, n square like this for instance ok that you typically get by doing a union bound over the entire time horizon ok. So, one may wonder is this the only is this the limit of how you can how far you can go.

(Refer Slide Time: 39:58)

Over n^{α} .

Is this best possible?

YES, in a sense, due to the Law of the Iterated Logarithm (LIL):

For $\{X_i\}_i$ iid $N(0,1)$,

$$\mathbb{P}\left[\limsup_{n \rightarrow \infty} \frac{S_n}{\sqrt{2n \log \log n}} = 1\right]$$

$$= \mathbb{P}\left[\liminf_{n \rightarrow \infty} \frac{S_n}{\sqrt{2n \log \log n}} = -1\right] = 1.$$

Well in fact, in some sense the answer is yes this is the best possible result that you can show to that you can hope to achieve in the sense of controlling uniform deviations of S_n . So, in a sense because there is in some sense a lower limit or a fundamental limit to how far you can go to control uniform deviations due to what is called the law of the iterated logarithm LIL.

And here is the statement of the LIL for the same setting where you are dealing with the sums of iid standard normals. So, the LIL basically says that for X_i drawn iid from normal 0, 1 the probability that the limsup as n tends to infinity of S_n . So, if you rescale if you scale S_n by $\sqrt{n \log \log n}$ basically roughly root $n \log \log n$ ok.

The probability that this limsup equal to 1 or the on the symmetric side if you take the liminf if you take the probability of the event that liminf as n becomes large of S_n scaled by $\sqrt{2n \log \log n}$ equal to -1 . These are actually almost surely true ok. So, these are actually events that happen with probability 1 ok.

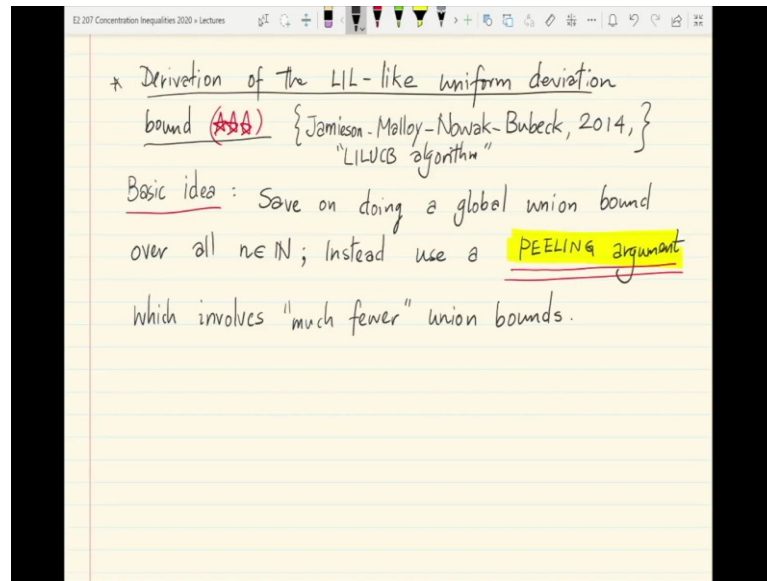
So, what this is again with respect to the figure is that if you; so, it does not matter how you draw things. If you draw a $\sqrt{2n \log \log n}$ and on the bottom you plot the envelope of $-\sqrt{2n \log \log n}$ ok, what it says is that what this law of iterated logarithm says is that infinitely often the process S_n is going to touch both these boundaries ok. We are going to come arbitrarily close to both these boundaries ok.

So, basically if you bump things down by a factor of $1 - \epsilon$ where ϵ is small then it means that infinitely often with probability 1 the process S_n is going to wander out of this tube infinitely many times. It cannot be constrained within this narrower tube ok. So, this sets an absolute lower limit because there is an almost sure event that the process S_n is going to breach these the boundary of such a tube with such bits.

And so, one cannot hope to really improve the scaling of the denominator better than $\sqrt{n \log \log n}$ ok. So, what we are going to show here this result can be thought of as a you know uniform deviations or finite time analogue of the law of the iterated logarithm. The law iterated logarithm is really a statement about its an asymptotic statement about the nature of the sequence S_n .

So, let us go ahead and try to derive this kind of best possible bound which is basically expressing uniform deviations of the S_n process ok.

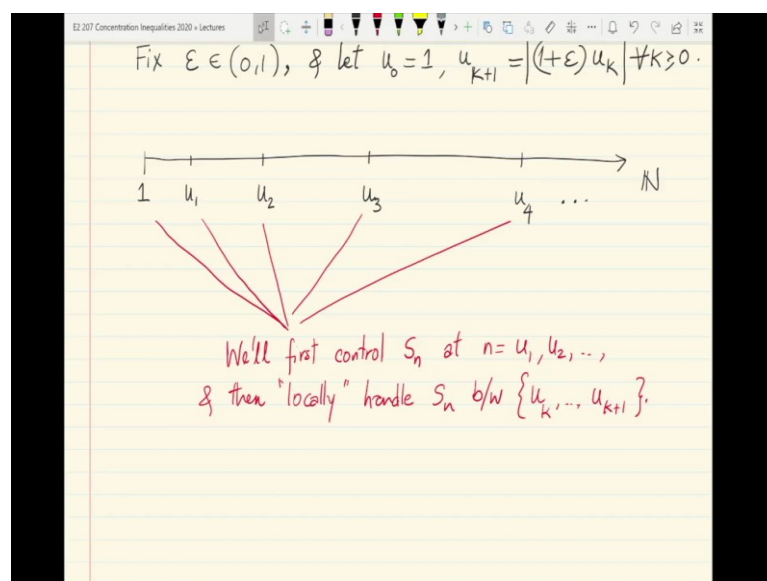
(Refer Slide Time: 43:54)



So, that will be the last part of this lecture. So, what is the basic idea here? The basic idea is. So, recall how we got the worst the weaker bound with the $\sqrt{n \log 2 n}$ square by δ . We basically did a union bound over all values of time t . So, in contrast the idea that I would like to present here to get sort of almost optimal LIL type uniform deviation bounds is to save on doing so many union bound.

So, save on doing a global union bound over all values of small n over natural numbers, instead we will use what is called a peeling argument or a peeling trick. This is sort of the name that has become sort of standard for this kind of procedure by now which involves broadly speaking much fewer union bonds or much fewer instants of time t .

(Refer Slide Time: 45:29)



So, here is so maybe let us set up some notation before i describe the idea in more detail. So, let us fix some ε constant between 0 and 1. You can think of ε equal to half for instance for a remainder of the proof and that so, let us define u_0 as time instant 1 and let us define the k -1st instant next time instant as $1 - \varepsilon$ times u_k , for all k greater than 0 greater than equal to 0 ok.

So, u_0 is 1, u_1 is $1 - \varepsilon$, u_2 is $1 - \varepsilon$ of u_1 and so on and so forth. In fact, technically we need all of these to be integers. So, its best to take the ceiling ok. So, this is sort of the precise definition of these in case. These are also a subset of the natural numbers. And to convey the entire idea of peeling using a picture here is the set of all times at which you want to control uniformly the deviations of S_n . So, it starts with time 1, ok.

So, that is u_0 , u_1 is here, u_2 is basically a geometric is sort of a scale factor $1 - \varepsilon$ time u_1 . And so, this grid u_k basically grows geometrically as a function of k ok it explores a rapidly. So, it is a geometrically increasing grid u_4 ok. So, in order to so, we have a process S_1, S_2, S_3 and so on which is taking values at every possible natural number. The u_k 's are a subset of these natural numbers.

Imagine that you can control imagine that you want to first you want to be able to control the fluctuations of the stochastic process S_n by first controlling S_n at these chosen points u_k 's ok. So, first let us say that we want to be able to control a sense that only these u_k 's these u_k 's are much in some sense far fewer numbers and the entire set of

natural numbers. So, we first control S_n at n equal to u_1, u_2 and so on and then locally handle S_n between each let us say u_k to u_{k+1} , ok.

So, in some sense it is a hierarchical approach. This peeling approach that is why it has the name peeling because you sort of its a hierarchical approach where you first control you basically block you bunch time up into exponentially increasing segments you control the process at the end points of each segment and then do sort of a separate local procedure to control the deviations of S_n within these segments ok.

So, we will see exactly how this idea pans out. So, at least in the process the total number of union bounds at the outer level over the u_k 's is in some sense going to be much smaller than taking union bound over all natural numbers and that is sort of the reason behind why you get this exponential improvement in the end ok.

So, this proof that we will present using the peeling argument it occurs in several parts several places in the literature on statistics in sequential statistics and online learning, but we will use the argument that is given in this paper by Jamieson et al, which is basically an online learning paper featuring what is called the LILUCB algorithm for multi armed bandits.

(Refer Slide Time: 49:43)

Let $\psi(x) := \sqrt{2x \log(\log x)}$.

Step ①: Control of $S_{u_k} \forall k \geq 1$:

$$\mathbb{P}[\exists k \geq 1: S_{u_k} \geq \sqrt{1+\varepsilon} \cdot \psi(u_{k+1})]$$

$$\leq \sum_{k=1}^{\infty} \exp\left(-(1+\varepsilon) \log\left(\frac{\log u_{k+1}}{\delta}\right)\right) \quad \left(\text{Union bound + Chernoff for Gaussian}\right)$$

So, let us do the peeling argument in detail. So, let us define this function for convenience $\psi(x)$ as roughly $\sqrt{x \log \log x}$ by δ . This is the kind of uniform deviations

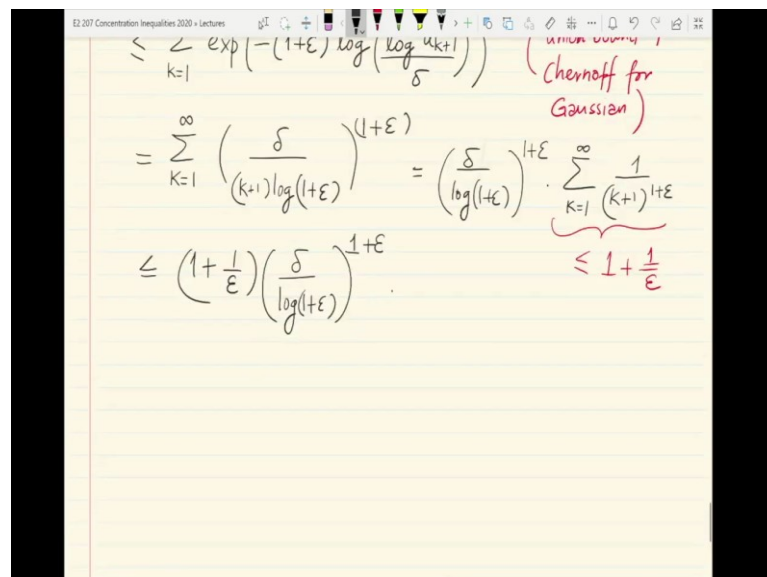
growth that we want in a result. So, there are three steps to the argument. So, step 1; step 1 is the control the first level of the hierarchy. So, this is the control of $S u_k$, for all k and let us see how we do this.

So, we just want to bound let us say the probability that there ever exists a value of k for which $S u_k$ exceeds let us say $\sqrt{1 - \epsilon}$. So, you should just think of the $\sqrt{1 - \epsilon}$ as just some number larger than 1 ok that is all the meaning of, that is all there is to the meaning of this ϵ . ϵ you can just think of ϵ conveniently as a half into ψ of u of $k - 1$, ok.

ψ of u of $k - 1$ is in some sense a threshold that we want to control every $S u_k$ by. So, we can first split this by the union bound over these k 's ok which are much smaller than the total number of I mean in some sense they are smaller than the total number of larger.

They avoid many natural numbers and many union bounds potentially wasteful. e raised to $-1 - \epsilon \log \log u_{k-1}$ divided by δ . This is just by using the union bound and also simultaneously applying the Chernoff bound for Gaussians. So, first using the union bound - the Chernoff bound for Gaussians ok, so, this is just a simple union bound followed by Chernoff.

(Refer Slide Time: 51:55)



$$\begin{aligned} &\leq \sum_{k=1}^{\infty} \exp\left(- (1+\epsilon) \log\left(\frac{\log u_{k+1}}{\delta}\right)\right) \quad (\text{Chernoff for Gaussian}) \\ &= \sum_{k=1}^{\infty} \left(\frac{\delta}{(k+1)^{\log(1+\epsilon)}}\right)^{1+\epsilon} = \left(\frac{\delta}{\log(1+\epsilon)}\right)^{1+\epsilon} \cdot \sum_{k=1}^{\infty} \frac{1}{(k+1)^{1+\epsilon}} \\ &\leq \left(1 + \frac{1}{\epsilon}\right) \left(\frac{\delta}{\log(1+\epsilon)}\right)^{1+\epsilon} \quad \leq 1 + \frac{1}{\epsilon} \end{aligned}$$

So, this is the same as the sum. So, by the way I am avoiding the I am avoiding sort of more precise computation because u_{k-1} is the ceiling of $1 - \epsilon$ times u_k . So, we presume that we can use fractions for you can use real numbers for integers, but with

small modifications the entire argument goes through in exactly the same way as we are doing now.

So, this is exactly summation δ by $k - 1 \log 1 - \epsilon$ the whole thing raise to $1 - \epsilon$. And the only part that depends on k if you remove you are left outside with δ by \log . So, let me write it here. This is δ by $\log 1 - \epsilon$ raise to $1 - \epsilon$ that comes out followed by a sum of 1 by $k - 1$ to the $1 - \epsilon$. And by some bound integral bounds you can easily show that this is at most $1 - 1$ by ϵ .

So, finally, we get the bound of δ by; so, let us say let us write $1 - 1$ by ϵ into this quantity δ by $\log 1 - \epsilon$ into $1 - \epsilon$ ok. So, you should really think of this complicated looking right hand side as just upper bounded by some constant times δ ok. So, because δ is a number smaller than 1, δ to the $1 - \epsilon$ can; obviously, be upper bounded by δ itself and the rest is just a constant $1 - 1$ by ϵ in the denominator which is a function of ϵ . So, it is just a constant times δ finely ok.

(Refer Slide Time: 54:02)

Step 2: Control of S_n , $n \in (u_k, u_{k+1})$

We'll use the following result from probability:

HOEFFDING'S MAXIMAL INEQUALITY: $\forall m \geq 1, x > 0$,

$$\mathbb{P} \left[\exists n \in [m]: S_n \geq x \right] \leq \exp\left(-\frac{x^2}{2m}\right).$$

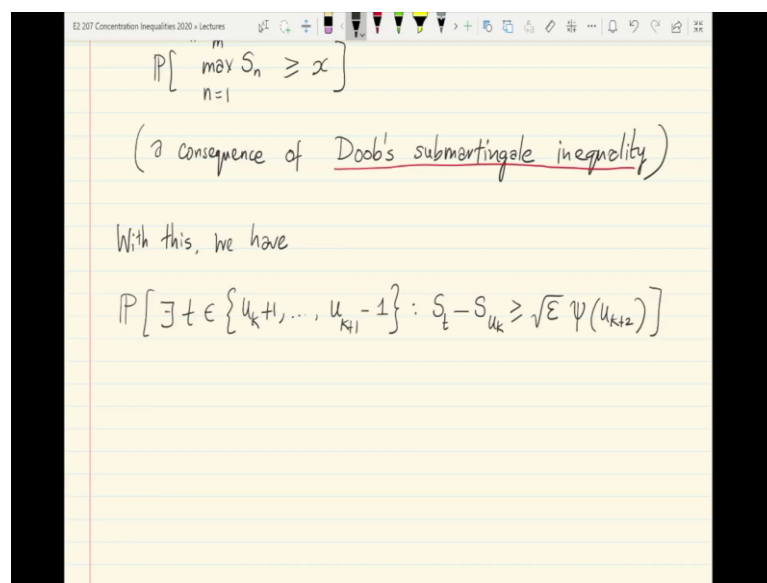
$$\mathbb{P} \left[\max_{n=1}^m S_n \geq x \right]$$

So, we now come to the step, 2nd step which is descending one step lower into the hierarchy and trying to do local control. So, this is about controlling S_n for n in the interval u_k to u_{k-1} right. So, how do we do this? So, we will take recourse to the following result from probability called Hoeffding's maximal inequality.

So, we will use the following result from probability called Hoeffding's maximal inequality which in our case applied to standard normal iid random variables sees the following. So, for every integer m and x greater than 0, the probability that there exists in n in 1 through m for which exceeds x ok.

This is really the same as the probability that the max over n equal to 1 to m of S_n is larger than equal to x is actually controlled by basically the same kind of Chernoff exponent Chernoff rate as what you would expect for S_m the last element S_n , ok.

(Refer Slide Time: 55:59)



Handwritten notes on a yellow background:

$$\mathbb{P}\left[\max_{n=1}^m S_n \geq x\right]$$

(a consequence of Doob's submartingale inequality)

With this, we have

$$\mathbb{P}\left[\exists t \in \{u_k+1, \dots, u_{k+1}-1\} : S_t - S_{u_k} \geq \sqrt{\epsilon} \psi(u_{k+2})\right]$$

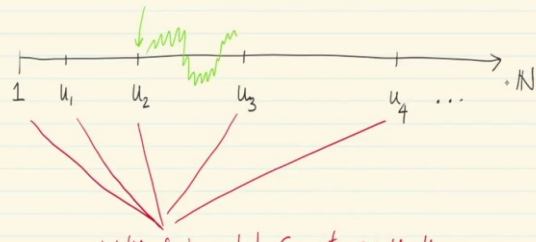
And the reason this holds you can go and look this up, this is a consequence of an even simply simpler stated, but a very powerful result called Doob's martingale inequality, Doob's sub martingale inequality in fact, ok. So, we are not going to prove this inequality in this class, but we will use it as a statement of fact. So, what it says is basically it allows you to control over a given interval of time m uniformly the supremum of a random walk is n , so, ok.

So, using this we can write the following. So, we have that the probability. So, using this we have that the probability that there exists a time t in the interval between 2 epochs $u_k - 1$ all the way up to so, this is a set of integers up to $u_{k+1} - 1$ ok. Since we have already controlled for u_k and u_{k+1} separately at the top level of the hierarchy, S_t relative to S_{u_k} ok. So, think of S_{u_k} . So, going back to the figure here think of suppose k is equal to 2.

(Refer Slide Time: 57:41)

which involves "much fewer" union bounds.

Fix $\varepsilon \in (0,1)$, & let $u_0 = 1$, $u_{k+1} = \lceil (1+\varepsilon)u_k \rceil \forall k \geq 0$.



We'll first control S_n at $n = u_1, u_2, \dots$,
& then "locally" handle S_n b/w $\{u_k, \dots, u_{k+1}\}$.

Let $\psi(x) := \sqrt{2x \log(\log x)}$.

So, you already controlled S_{u_2} and we are trying to sort of understand how much fluctuation excess fluctuation there is in this interval between u_2 to u_3 . So, $S_t - S_{u_k}$ exceeds let us say something again $\sqrt{\varepsilon}$ into $\psi(u_{k+2})$. So, u_{k+2} is chosen conveniently as you will see soon.

(Refer Slide Time: 58:06)

$$\begin{aligned}
 & \mathbb{P}[\exists t \in \{u_{k+1}, \dots, u_{k+1} - 1\} : S_t - S_{u_k} \geq \sqrt{\varepsilon} \psi(u_{k+2})] \\
 &= \mathbb{P}[\exists t \in [u_{k+1} - 1 - u_k] : S_t \geq \sqrt{\varepsilon} \psi(u_{k+2})] \\
 & \quad \text{(H.M.I.)} \quad \quad \quad \text{(by recentering time)} \\
 &\leq \exp\left(-\varepsilon \cdot \cancel{u_{k+2}} \log\left(\frac{\log u_{k+2}}{\delta}\right) \cdot \frac{1}{\cancel{u_{k+1} - u_k}}\right) \\
 &\leq \exp(-(1+\varepsilon))
 \end{aligned}$$

This is equal to the probability that there exist t . So, just by shift by recentering time you can just recenter everything from time u_k onwards, it does not change the nature of the random work. So, there exists a t in the interval. So, basically in 1 through $u_{k+1} - 1 - u_k$

ok such that S_t exceeds $\sqrt{\epsilon} u_{k-2}$ as before. This is the same because all we have done is we have basically started time for the random work from the bar the time u_k .

So, this is by recentering time and now we can apply Hoeffding's maximal inequality to get an upper bound of ϵ . So, the square; so, whatever is here get squared. So, ϵ the square of u_{k-2} is precisely $2 u_{k-2} \log \log u_{k-2}$ by δ . And finally, we have to divide by this by 2 times the number of steps in the random work which for us is just $u_{k-1} - u_k$ ok, so, right.

So, we have this upper bound. Now, we can take this ϵ here and this u_{k-2} and this difference essentially between u_{k-1} and u_k and show that this ratio is lower bounded by $1 - \epsilon$ ok. So, this is an easy exercise. So, you can replace with $1 - \epsilon$.

So, this is basically the reason why the geometry grading works essentially says that you know there is a factor of $1 - \epsilon$ increase in the sort of scale and at the same time It also saves you many many union bounds and the top level there is an exponential reduction.

(Refer Slide Time: 60:41)

$$\leq \exp\left(-\epsilon \cdot \frac{u_{k+2}}{\log \frac{u_{k+2}}{\delta}} \cdot \frac{1}{u_{k+1} - u_k}\right)$$

$$\leq \exp\left(-(1+\epsilon) \log\left(\frac{\log u_{k+2}}{\delta}\right)\right) \quad \left\{ \because \frac{\epsilon u_{k+2}}{u_{k+1} - u_k - 1} \geq (1+\epsilon) \right\}$$

$$= \left(\frac{\delta}{(k+2) \log(1+\epsilon)}\right)^{1+\epsilon}$$

So, $\log \log u_{k-2}$ by δ ; the reason for this is in short that ϵu_{k-2} divided by $u_{k-1} - u_k - 1$ can be shown to be greater than equal to $1 - \epsilon$. There is always the ratio $1 - \epsilon$ here ok. And this in turn is just simply expressed as δ over using the definition of u_{k-2} . u_{k-2} again is basically $1 - \epsilon$ raise to $k-2$ up to a small approximation.

So, ignoring that approximation error we can just write this as δ over. So, $\log u_{k-2}$ is just $k-2 \log 1 - \epsilon$ ok and the entire thing is raised to the exponent $1 - \epsilon$ ok. So, this is again essentially a constant times δ , this is another constant time δ . And so, we have essentially managed to bound both local and global fluctuations of S_n by basically constant time δ and the last step is to just put both of these together.

(Refer Slide Time: 62:06)

$((K+2)\log(1+\epsilon))$

Step ③: Putting everything together

By steps ① & ②, with prob. at least

$$1 - \left\{ \left(1 + \frac{1}{\epsilon}\right) + \frac{1}{\epsilon} \right\} \left(\frac{\delta}{\log(1+\epsilon)} \right)^{1+\epsilon}$$

So, to get the LIL type result. So, putting everything together, so, this is the last step. So, by steps 1 and 2 we get that so, steps 1 and 2 essentially give you probabilities of bad events upper bounded by some numbers. So, here is the good event, this is the complement with probability at least $1 -$ the total bad events. So, step 1 gave you a bad event probability of this much and step 2 gave you a bad event probability of this much ok.

So, we can argue that if you sum up the bad event probabilities and subtract them from 1, so, step 1 gives you basically $1 - 1$ by ϵ into some expression and then a further 1 by ϵ into this thing δ by $\log 1 - \epsilon$ base 2 $1 - 1 - \epsilon$. So, recall that in step 2 you have to sum over all the all case to get the net bad event of step 2 that any of the local fluctuations is bounded.

(Refer Slide Time: 63:43)

$$\forall k \geq 0, \forall t \in \{u_{k+1}, \dots, u_{k+1}\},$$

$$S_t = S_t - S_{u_k} + S_{u_k}$$

$$\leq \sqrt{\varepsilon} \psi(u_{k+2}) + \sqrt{1+\varepsilon} \psi(u_{k+1})$$

So, with this much probability for all epochs k and for all time steps local time steps within the epoch k , so, the right endpoint is now included thanks to the step 1 result. We must have that S_t . So, can we bound S_t ? Find the S_t is actually $S_t - S_{u_k} - S_{u_k}$ and each of these terms is bounded the first 1 by step 2 and the second term the difference by step 2 and the last term by step 1.

So, we know that by step 2 this is at most $\sqrt{\varepsilon} \psi(u_k - 2)$ and the second deviation is bounded in height with high probability by $\sqrt{1 - \varepsilon} \psi(u_k - 1)$.

(Refer Slide Time: 64:40)

$$\leq \sqrt{\varepsilon} \psi(u_{k+2}) + \sqrt{1+\varepsilon} \psi(u_{k+1})$$

$$\leq \sqrt{\varepsilon} \psi((1+\varepsilon)^2 t) + \underbrace{\sqrt{1+\varepsilon}}_{\leq 1+\sqrt{\varepsilon}} \psi((1+\varepsilon)t)$$

$$\leq (1+2\sqrt{\varepsilon}) \psi((1+\varepsilon)^2 t).$$

And we just need to do some very basic algebra and bounding. So, we know that t is within u_k to u_{k-1} and so, u_{k-2} must be at most $1 - \epsilon$ square into t and ψ is a mono term function $-\sqrt{1 - \epsilon}$. Again you can bound u_{k-1} by at most $1 - \epsilon \times t$, ok.

You can use a further bound on $\sqrt{1 - \epsilon}$ of $1 - \sqrt{\epsilon}$ and finally, you get something like $1 - 2\sqrt{\epsilon}$ ψ of $1 - \epsilon$ the whole square which is another constant times t ok. So, the orders are all essentially in the right place. So, for instance, so, this essentially completes the argument.

(Refer Slide Time: 65:34)

$$\leq (1+2\sqrt{\epsilon}) \psi((1+\epsilon)^2 t).$$

So, if we take, say, $\epsilon = 1/2$, then

$$\mathbb{P}\left[\exists n \in \mathbb{N}: S_n > (1+\sqrt{2}) \sqrt{2 \cdot \frac{9}{4} n \log\left(\frac{\log \frac{9}{4} n}{\delta}\right)}\right]$$

$$\leq (1+4) \cdot \delta \cdot \frac{\delta^\epsilon}{(\log \frac{3}{2})^{3/2}} = (\text{const.}) \delta.$$

If you specialize and write this for instance with ϵ equal to half if we take say ϵ equal to half and write down this result, then the probability that there is ever a time uniformly over the entire time horizon, where S_n exceeds $1 - \sqrt{1 - 2\sqrt{\epsilon}}$ that is $1 - \sqrt{2}$ into the $\sqrt{\text{of}}$ let us say twice 9 by 4 that is $1 - \epsilon$ the whole square $n \log \log$ of 9 by 4 n divided by δ .

So, that gives you the $\log \log n$ type term is no more than $1 - 4$ into some constant into δ into δ to the ϵ by $\log 3$ by 2 raised to 3 by 2 which is ultimately some constant some explicit numerical constant times δ ok and which is essentially in the form of the it is a finite time version of the law of iterated logarithm ok.

And the key idea here was to use this kind of hierarchical peeling idea or peeling trick as its commonly called to essentially try to save on the number of union bounds getting as close to optimal as possible.

Thank you.