

Stochastic Modeling and the Theory of Queues
Prof. Krishna Jagannathan
Department of Electrical Engineering
Indian Institute of Technology – Madras

Lecture –35
G/G/1 Queue and Little's Theorem

(Refer Slide Time: 00:13)

Lec 29: G/G/1 Queue & Little's Theorem

Renewal arrivals $\{X_i, i \geq 1\}$ $E[X] = \frac{1}{\lambda} < \infty$

Arrival process - Renewal process $\{X_i, i \geq 1\}$ $E[X] = \frac{1}{\lambda} < \infty$

Service time of customer i $V_i \leftarrow iid$ $E[V] = \frac{1}{\mu} < \infty$

$\{X_i, i \geq 1\}$ indep $\{V_i, i \geq 0\}$

Today we will discuss the G / G / 1 Queue and a very important result in queuing known as Little's theorem which relates the average occupancy of a queuing system to the average waiting time in the queuing system. So, we will discuss the G / G / 1 Queue a little bit and we will prove Little's theorem using the renewal reward approach that we have developed for the G / G / 1 Queue.

Although, the Little's theorem itself is not just limited to G / G / 1 Queue so it is actually applicable even to multi server systems as we will see, but we will prove it only for G / G / 1 that this today's agenda. So, you guys already know what the G / G / 1 Queue is. So, this is a server and this is a queue. You have renewal arrivals according to some distribution x. So, x i is of inter arrivals.

And the service times of each job are also iid. So, that is actually what is usually known as GI / GI / 1 Queue. You have independence between consecutive service time independence between consecutive arrival times. So, some people write GI / GI / 1, but usually G / G / 1

means without any further qualifications usually means $G / G / 1$. So, you can look at iid so arrival times arrival process is a renewal process with inter arrival time x_i 's.

And expected x let us say $1 / \lambda$, λ is the arrival rate and this is assumed to be finite and the service process so each customer whenever this renewal process ticks there is a customer that arrives. So, if you look at the time axis let us say this is time 0. So, you have some renewal arrivals and whenever this arrival occurs there is a customer who comes to the system.

And each customer brings with him or her an inherent service time. So, we assume that the server works at a constant rate let us say 1 and you can think of this customers as bringing with them a certain file size or job size or whatever. So, some could be large and some could be small and what we are saying is that these random variables are also iid this file size or customer service time random variables are also iid.

So, these are the file sizes or let us say v_i so you can say that the servers serves at a constant rate 1 and each file of size V_i will take an amount of time V_i to get served at the server. So, we assume that the service time of customer i is also v_i which are also iid and we will say x_i so this sequence x_i is independent of this sequence V_i . Well this sequence V_i I will write i greater than or equal to 0 because for the $G / G / 1$ Queue convenient convention is to start ticking time when the first arrival comes to the system.

The system is empty for a very long time there is nobody at all and then suddenly the first arrival comes and you start counting time at that point. So, what you will do is you can take the first arrival at 0. So, in some sense it is renewal process after time 0 there is a first arrival at time 0 always because you are starting your clock at that time. So, this you can think of as customer with file size V_0 who comes to the system at time 0.

Starts getting served and then the subsequent arrivals are renewal arrivals and each of course brings in a job size or file size or customer service time V_i . So, this x_i are iid, V_i are iid again we can write expectation of V as $1 / \mu$ it is also assumed to be finite though and x_i are independent of V_i . So, this is the model of a $G / G / 1$ Queue. So, the arrival distribution can be anything, inter arrival distribution can be anything, service time distribution can be anything.

They just cannot be correlated among themselves or across the x_i and V_i that is the $G / G / 1$ Queue for you.

(Refer Slide Time: 06:20)

NPTEL

Renewal process
 $\sim X_i$

Arrival process - Renewal process $\{X_i, i \geq 1\}$ $E[X] = \frac{1}{\lambda} <= 1$

Service Time of customer i $V_i \leftarrow iid$ $E[V] = \frac{1}{\mu} <= 1$

$\{X_i, i \geq 1\}$ indep $\{V_i, i \geq 0\}$

Notation
 $W_i \leftarrow$ Total time spent in the system by customer i
 $W_i^q \leftarrow$ Time spent waiting in the queue by customer i

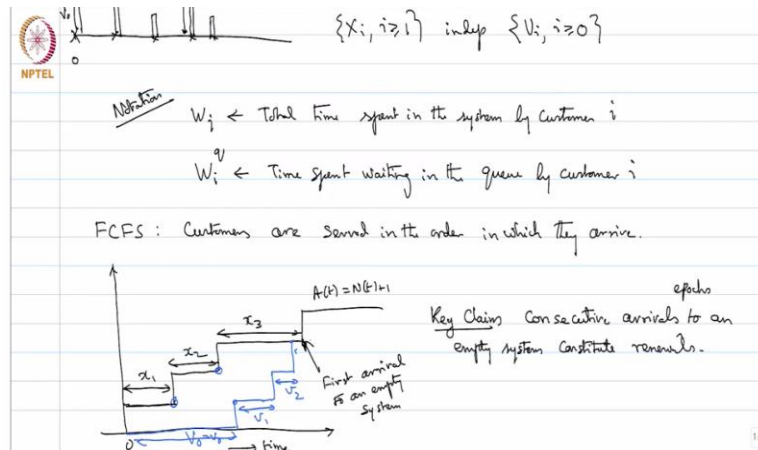
FCFS: Customers are served in the order in which they arrive.

Now, we are interested in looking at the waiting times of these files or these customers in queue. The usual notation is W_i is taken as the notation for total notation, total time spent in the system by customer i and the book uses W_i^q . So, these are random variables so these big capital letters. Just in case you are confused W_i^q is the time spent waiting in the queue by customer i .

So, what is the difference between W_i and W_i^q ? W_i is the total from the time customer i arrives at the system, wait, get into service, finish service and get out that is the total system time. W_i^q is simply the time spent waiting not including the time getting served. So, W_i is basically the W_i^q plus your own service time you spent getting served yourself. Now, we are going to do I have to specify the service discipline.

I am going to do only first come first served in class, but I will also point out why this FCFS is not I mean much of what we are going to say in today's lecture we will hold for other service disciplines as well, but FCFS is probably the most intuitive and most common. So, this is first come first served. So, customers are served in the order in which they arrive. So, person who came earlier will be served earlier.

(Refer Slide Time: 09:04)



This is not of course necessary you can do something else also. I am just going to analyze this situation FCFS situation. So, please remember that there is the time begins when the 0th customer comes at time 0 with some file size V_0 that is when we are starting time. So, I want to just plot for you the evolution of this FCFS queue $G / G / 1$ queue. So, the first arrival happens at time 0.

So, let me plot the arrival process has $1 + N t$ where $N t$ is the renewal process. So, it starts at time 0 and there is an arrival here let us say. So, this is just the arrival process I am plotting the arrival process. So, let us say this one more arrival and so on. So, this I can call it as arrival process this is time of course I am going to call it A of $t = N$ of $t + 1$ because N of t is the arrival process renewal process N of t .

Why am I adding + 1? There is arrival at 0 always after that you are counting renewals. Now, I also want to plot the departure process. So, this is the arrival process then this is FCFS and I want to plot the departure process. So, the customer who arrives at 0 starts getting served at 0 he comes to an empty system. Now how long will he be in the system the first customer to come at 0 his service time which is V_0 .

So, his departure will occur at time V_0 so this let us say that is my departure process I am drawing this in blue. So, this length if you look at it right this length will be equal to this length will be V_0 . Of course, these lengths are your x_i may be here I should not write apologies so these are realizations. So, maybe you should write so this is like x_1 is equal to little x_1 this is x_2 that is x_3 and so on these are the inter arrival times.

And this V_0 there is some little V_0 and when this guy leaves the second person gets into service. So, the second customer who gets into service at this point is the person who arrived over here is that right. So, that person get serve for some time and then he leaves and then the person who came here gets into service. The second arrival not counting the arrival at 0 gets into service and in this picture let us say he finish a service.

So, in my picture so this guy is v_1 , this guy is v_2 this interval is v_2 . So, v_1 corresponds to the service time of the customer who entered here and v_2 corresponds to the service time of the customer who entered here. So, at this point how many people are there in the system? System has emptied and then in this small duration system is empty and then there is another arrival and so on (()) (14:03) sort of thing repeats.

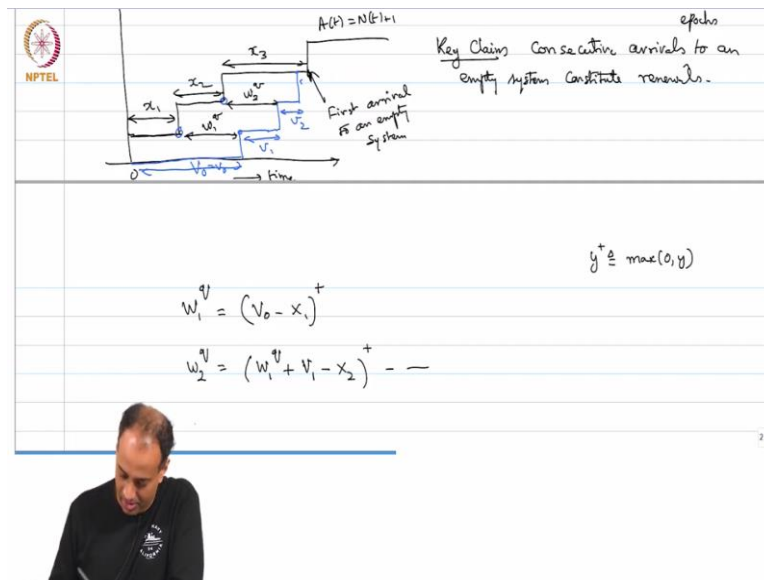
Now a key fact that we will argue prove at least prove using pictures we can make this very rigorous easily is that consecutive arrivals epochs to an empty system constitute renewals. This is a key fact. So, what do I mean? So, if you look at this picture at time 0 you had an arrival to a empty system and what is the next arrivals that occurs to an empty system this guy.

This is the first arrival to an empty system of course not counting the 0th arrival to the empty system. The time at 0 is the beginning of time a bunch of arrivals come and they leave and then there is a another arrival to an empty system. Now of course the similar sort of thing plays out. This guys getting served more people will come and so on and let us say this system empties again and one more time there will be an arrival to the empty system.

What we are saying is that these durations constitutes renewals so that is what we are going to show. This is not all the difficult to see because in some sense if your system emptied out statistically what is going to happen after if you have an arrival to an empty system statistically what is going to play out again is the same as what played out starting a $t = 0$ you again have some arrivals which have inter arrival time the renewals and service times are iid.

I am not saying they will be equal of course they will be different, but statistically they will be the same and statistically they will also be independent because all these x_i and v_i are independent you see what I mean so that is what we have to argue.

(Refer Slide Time: 17:08)



Now to understand why this is true I want to be able to track how long these customers wait before they get into service. Of course, the arrival at $t = 0$ waits for 0 time. The second arrival which occurs after $t = 0$ how long does that person wait? The x_1 arrival so to speak waits till the time 0 arrival complete service. So, this is the amount of time that the customer number 1 waits in queue before getting into service.

So, this should be w_1^q . So, in this picture I can write $w_1^q =$ I can write this in terms of random variables. I have drawn specific realizations I can write $w_1^q =$ how much $v_0 - x_1$, but see in my picture v_0 is bigger than x_1 , but if the amount of time if x_1 the first arrival took place after v_0 completed then what happens the waiting time will again be 0. So, if $v_0 - x_1$ is negative then I should count this as 0 I should not count negative time.

So, I can put plus max of whatever this is, 0. So, basically I am using the notation $y^+ = \max(0, y)$. Now what happens to the second customer? The x_2 arrival so to speak waits for that much time. So, what is w_2^q can you tell me some of those equations? I would like it in terms of the previous waiting time so it will be $w_1^q + v_1 - x_2$. So, what am I saying so it is w_2^q is simply it is the time that it arrived.

So, this total length is $w_1^q + v_1 - x_2$ will be w_2^q . So, let me just write it out. I have written it is $w_1^q + v_1 - x_2$. The plus will come for the same reason if this is negative you will have to count it as 0. So, you can generally say you can keep doing this.

(Refer Slide Time: 20:24)

$y^+ \triangleq \max(0, y)$

$$w_1^q = (V_0 - X_1)^+$$

$$w_2^q = (w_1^q + V_1 - X_2)^+ \quad \leftarrow \text{Lindley's Recursion}$$

$$w_i^q = (w_{i-1}^q + V_{i-1} - X_i)^+ \quad i \geq 2$$

Index of the first arrival (after $t=0$) which enters an empty system is the smallest 'n' that satisfies $w_n^q = 0 \Leftrightarrow$
 smallest 'n' for which $\sum_{i=1}^n X_i \geq \sum_{i=1}^n V_{i-1}$.

You can write for the i th customer you can write $w_i^q = w_{i-1}^q + v_{i-1} - x_i$ whole plus. If this turned out to be negative then you will put 0. This is true for i greater than or equal to 2. So, if you take so this bunch of equations actually I can actually just delete the second equation it is just a consequence of it is contained in the third. So, this bunch of recursive equations are known as Lindley's Recursion.

It is a recursive equation you can write to characterize the waiting time of the i th customer in terms of the waiting time of the $i - 1$ customer and the x_i and v_i . See this x_i and v_i are given to you. I mean those distributions are known to you and you will be able to write these waiting time in terms of this x_i and v_i and the previous waiting time. This is a first come first serve $G / G / 1$ Queue.

Now, we want to see the 0 that arrival at time 0 had 0 waiting time then wait at all. Now, what is the index of the next arrival for which faces 0 waiting time. So, in our picture for example the x_1 arrival the first arrival not counting the arrival at 0 has some positive waiting time w_1^q and then w_2^q was also positive and the third arrival which happened to an empty system had 0 waiting time.

If you think about why this happened? What is it about that the third arrival in terms of the x_i and v_i that helped it to face 0 waiting time. See, it is the first arrival for which the sum of the x_i is greater than or equal to the sum of the previous v_i . See for the first arrival v_0 was bigger than x_1 . For the second arrival also w_2^q was positive because $v_0 + v_1$ was bigger than $x_1 + x_2$.

But for the third arrival that x_3 is quite large in the way I have drawn it. It turns out that the $x_1 + x_2 + x_3$ which is the epoch x_3 of the third arrival epoch is larger than the sum of the service time of all the customers who have come before. In fact you can see this from the Lindley's Recursion $w_i \geq 0$ for the first time what is the smallest i for which $w_i \geq 0$ will also be the first index for which the sum of x_i is still that particular index is exceeds the sum of the v_{i-1} I guess.

So, I am going to write it like this so index of the first arrival so after $t = 0$ which enters empty system is the smallest n that satisfies $w_n \geq 0$. This is the thing about the first arrival sees an empty system this is equivalent to the smallest n for which sum over $i = 1$ to n x_i is greater than or equal to sum over $i = 1$ to n v_{i-1} . In the picture I have drawn this n is 3. Now, I will argue so what does this depend on?

So you think about this. So, you have a bunch of random variables further over the queuing system, you have a bunch of iid random variable x_i . We have given a bunch of iid random variables v_{i-1} and I am asking you keeping looking at these x_i and v_i and tell me the first time that sum over x_i exceeds sum over v_{i-1} . Now what does this sound like? This sounds like something we have already studied.

(Refer Slide Time: 26:25)

NPTEL

Index of the first arrival (after $t=0$) which enters an empty system is the smallest n that satisfies $w_n \geq 0 \Leftrightarrow$
 smallest n for which $\sum_{i=1}^n x_i \geq \sum_{i=1}^n v_{i-1}$.

\Leftrightarrow Smallest n for which $\sum_{i=1}^n (x_i - v_{i-1}) \geq 0$

$T = \min \{n \mid \sum_{i=1}^n (x_i - v_{i-1}) \geq 0\}$

T is a stopping rule (possibly defective) for $(x_i, v_{i-1}, i \geq 1)$.

2/2



So, let me put it this way so I can write this as the smallest n equivalently the smallest n for which sum over $i = 1$ to n $x_i - v_{i-1}$ is greater than or equal to 0. See you call $x_i - v_{i-1}$ as some other variable y_i they are also iid. So, we are looking at the first time that these y_i

could be positive or negative I am looking at the first n such that $\sum_{i=1}^n y_i$ is now negative. What does this sound like?

See this is like a stopping rule except that there are two sequences of random variables so you know what sequence of I mean stopping rule J for a sequence of random variable x_i you look at each n you look at just the first n realization and decide to stop or not. Now, I am going to look at not just x_i I have a bunch of x_i , I have a bunch of v_i and I am going to look at both and I am going to stop whenever $\sum_{i=1}^n x_i$ exceeds greater than or equal to $\sum_{i=1}^n v_i - 1$.

So, I am write a stopping rule J as the minimum of n such that. Now, clearly the event you can argue you are only looking at x_1 through whether you stop at n depends only on x_1 through x_n and v_0 through v_{n-1} . So, is J a stopping rule it could be defective. It is not clear that you will eventually see a arrival to an empty system what if queue keeps building up or something.

It could be defective possibly defective for not just x_i , but $x_i v_{i-1}$. So, you can define this is like a slightly generalized version of stopping rule where you do not have just one sequence you have two sequences. It does not matter you are only looking at so whether you stop a time in or not is dependent only on the information available at that time you are not looking ahead.

You are looking at x_1 through x_n and v_0 through v_{n-1} . Remember, J is not a time J is an index that is why I am writing stopping rule instead of saying stopping time. It is the index of the arrival it is not actual time. Now, what we will see is that J is a genuine stopping rule not a defective stopping rule if the average inter arrival duration is greater than the average service duration if expected x is greater than expected v .

Then J is a legitimate stopping rule which means that if you will stop in finite you will always have an arrival with probability 1 which enters in empty system. Furthermore, if expected x is greater than expected v you can show that expected J is finite slightly non trivial, but you can show it and what we can also show is that once the stopping rule realizes as $J = \text{some } n$ then these subsequent arrivals let us say condition of $J = n$.

The subsequent inter arrival time will be x_n, x_{n+1} etcetera in our case $n = 3$ then the subsequent inter arrival times are x_n, x_{n+1} etcetera and this subsequent service times are v_{n-1}, v_n etcetera and again Lindley's equation can be written for the subsequent waiting time and all that. Now, you can also argue that condition on $J = n$ what are the inter arrival times x_n, x_{n+1} etcetera which are independent of the previous inter arrival times in the previous busy period so to speak when the server were busy.

And likewise for the v_n, v_{n-1}, v_{n+1} etcetera are independent of the previous service times and all these for every n your conditioning on $J = n$ these random variables x_n, x_{n+1} etcetera and v_n etcetera are independent of J and independent of the previous x_i and v_i in the previous busy period. So, I have argued independence and identical distribution condition on $J = n$, but this is true for every n .

So, you are basically argued that these are renewals. What do I mean by these are renewals? The times at which that are consecutive customer entries to a empty system these constitute renewals.

(Refer Slide Time: 32:37)

Theorem: For a $G/G/1$ queue with $E[X] > E[V]$ (ie, $\lambda < \mu$), the subsequence of arrivals that see an empty system form a renewal process. The expected number $E[D]$ of arrivals between arrivals to an empty system is finite & the expected time between the commencement of successive busy periods is equal to $E[D] E[X]$.

So, we will state this as a theorem for $G/G/1$ Queue with expected x greater than expected v i. e. $\lambda < \mu$. Thus, subsequence of that see an empty system in the renewal process. The expected number of arrivals between arrivals to an empty system is finite and the expected time between successive busy periods between the commencement of successive busy periods is equal to expected J times expected x .

This is the expected inter renewal duration. So, I have used the term here busy periods. So, you go back to this picture so this particular picture shows one busy period. So, you have an arrival to a empty system then there is a whole queue build up wait, wait, wait everybody leaves and then the system empties then there is another arrival. So, I am saying these successive intervals constitute renewal period and this is one busy period basically.

So, what we are saying is that so we are saying two things let us go back to the theorem now. We are saying that the subsequence that see an empty system form a renewal process. So, what do I mean by subsequence? You have the arrival sequence x_1, x_2 this is the renewal process to begin with, but not all of them see a empty system. There is a subsequence that sees an empty system.

So, in our picture the first arrival of course is an empty system I mean first arrival at 0th arrival sees an empty system then the third arrival sees an empty system then maybe the 7th and 13th or whatever see an empty system. We are saying that particular sequence that particular subsequence of this x_i itself constitute a renewal process and we are also saying that J which is the stopping rule which is the minimum x which is the first index for which you see a empty system.

This J is a legitimate stopping rule with finite expectation that we have not proved. It is not difficult to prove it I think there is an exercise in your book it uses a truncation argument we have to use inter arrival distribution which are bounded prove it for that and then use a truncation argument for proving expected J is finite and then the expected real time, the continuous time between one arrives which is an empty system and the next arrival which is an empty system.

The expected time will simply it is like expected $S J$ which is just by Wald. Expected J times expected x which is also be finite. If expectation of x is greater than expectation of b that is the arrival rate should be less than the service rate. λ is the arrival rate which is 1 over expected x and μ is 1 over expected v so λ is less than μ then we are guaranteed that you will have infinitely many of these entries to empty systems.

And the expected duration between these renewals these expected renewal time will also be finite which is just expected J time expected x . Now we will use this bigger renewal process

this what I call the subsequence what sees empty systems to prove Little's theorem that will be a discussion next class.