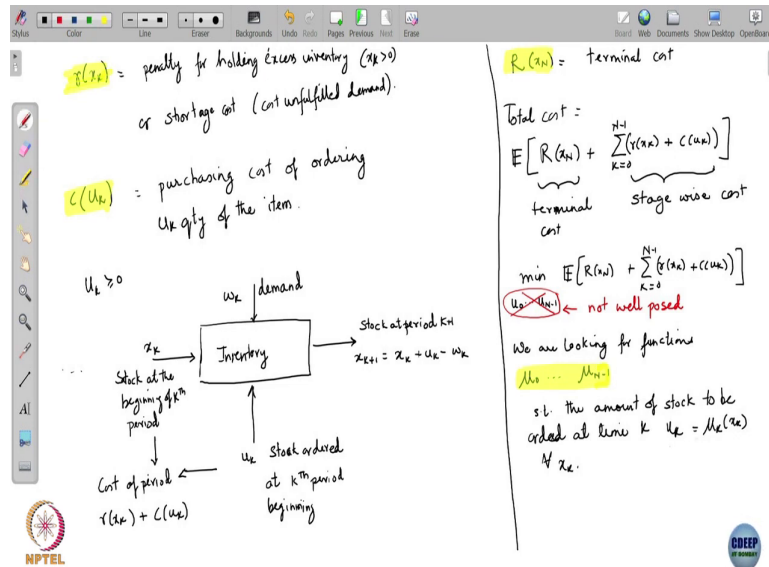


Optimization from Fundamentals
Prof. Ankur Kulkarni
Department of Systems and Control Engineering
Indian Institute of Technology, Bombay

Lecture - 23B
Policy and value function

(Refer Slide Time: 00:20)



So, this way of writing the problem as I said is not, this is not well posed here, this is not well posed. So, because this is not well posed the right way of posing the problem is that we are not looking for actions u_0 to u_{n-1} , but rather functions. We are looking for functions and now let me denote these functions in the following way μ_0 to μ_{n-1} right.

So, μ_0 to μ_{n-1} such that, the amount of stock to be ordered at to be ordered at time k , which is u_k is equal to then μ_k of x_k , for all values of x_k . So, whatever be the value of

So, what we are looking for are these functions ok . These functions μ_0 to μ_{n-1} and as these functions would then once these functions are fully specified, what they let you do is pick your action, which is the amount of stock to be ordered as the function of the information that would be best.

Correct way of posing the problem is

$$\min_{\mu_0 \dots \mu_{N-1}} E \left[R(x_N) + \sum_{k=0}^{N-1} c(u_k) + c(w_k) \right]$$

$\mu_0 \dots \mu_{N-1}$ = policy (π)
↑ ↑
"strategy"
 $u_0 \dots u_{N-1}$ = "actions".

1) State' of the system.
 (x_k)
"configuration" of the system

2) State evolution or dynamics
 $x_{k+1} = f_{x_k}(x_k, u_k, w_k)$
 ↑ ↖ noise @ k
 action @ k

w_k = independent random variables.

3) Constraint on actions ($u_k \geq 0$)
 $u_k \in U_k(x_k)$

4) Additive cost from
 $E \left[g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, u_k, w_k) \right]$
 $(g_N(x_N) = R(x_N))$, $g_k(x_k, u_k, w_k) = r(x_k) + c(u_k)$

So, therefore, the problem the correct way of posing the problem, the correct way of posing the problem is that you want to minimize the expected cost; includes the terminal cost, plus all these running cost, starting from 0 to n minus 1 r of x_k plus c of u_k and you want to minimize this by choosing μ_0 to μ_{n-1} , these functions.

So, mathematically this is what is being done, you are getting, you are trying to optimize this particular cost function, now by choosing a sequence of functions ok. Now, the reason we needed to do this was, because demand itself is something that was going to be realized in the future alright. And, it was not known its value was not known to us.

So, we had to essentially plan for every possible value and which and when you are planning for every possible value, what we are effectively doing is finding this these the sequence of such functions right. Now, there would be another case, suppose in which someone told you what the exact demand is going to be ok. Hypothetically say, suppose someone gave you the exact value of the demand.

In which case this the expectation here would be moved, the expectation in this that I have written here, would have would be moved. And, then what you would know the exact value of the information or the exact value of the amount of stock that would be that would be present at the beginning of each time period.

In that case, in that case you can actually plan to for a, for taking a particular action, because you know the exact state right. So, there are no you do not need to plan for every possible state, you can simply try to find the value of the action for the particular state that is going to be realized alright. For the, or the particular stock level that is going to be realized. That is that, kind of case where the demand is known would be the case where the noise is deterministic.

So, it is really not noise any more the demand is deterministic it is known in its value is known in advance and all you are doing is planning for specific events, that are going to happen in the future, specific known events that are going to happen in the future. That is a much easier version of this particular problem, but the more general version of the problem is involves taking these a taking these decisions without knowing, what the information is going to be alright.

So, this μ_0 to μ_{n-1} that, I have written here this is what is called a policy ok. And, it is often denoted by π ok. So, the distinction here that I have let me mention 1 more a point here, the distinction between what is called, what I have been saying as taking a particular action, which is the amount of stock to be ordered versus planning for every possible, level of stock. The distinction between these two is the distinction between what we call in stochastic control and in games, as the distinction between actions and strategy.

What we have effectively done is because we do not know the information that is going to be realized in the future, we are not commit committing or deciding on what particular action is to be taken. Rather what we are doing is we are coming up with a strategy. A strategy which says that, you know if I had this information this is what I would do, if I have that information that is what I would do, alright.

So, this here these functions μ_0 to μ_{n-1} , they to or your policy effectively constitutes decides doing what is called strategy. It involves coming up with these plans, whereas, u_0 to u_{k-1} , these which are the actual decision they are going to take these constitute actions alright.

So, our problem therefore, is to come up with is to come up with these strategies or these policies right. Now, the every time a problem involves noise it whether you whether; however, simple or complex the noise may be the problem shifts from the space of actions to the space of strategies. Because, we do not know the value of the noise in advance, it means that the, you cannot plan for any specific actions.

So, problem in of choosing actions is simply a problem of choosing these vectors right. So, sorry this is not $k-1$ this should be $N-1$, the problem of choosing these action is simply a problem of choosing these N vectors u_0 to u_{N-1} alright. And, I can stack them up and essentially think of this as one composite decision problem involving one long vector right, so, u_0 to u_{N-1} .

However, this problem here the problem above is not a problem of choosing vectors, it is a problem of choosing functions. So, the space of the problem itself has changed. The, the problem of choosing actions which is the problem that you would have if say the w s were deterministic, that problem is the problem of simply the problem of vector optimization whereas, this problem here is the problem of optimization over functions right.

It is a problem of finding the right sequence of functions, not just as sequence of vectors ok. So, this distinction is a, is actually what makes this dynamic programming or dynamic optimization, significantly harder than static optimization alright. Now, the r as I go a little further I will explain to you how we can actually reduce, somehow this problem which is involves trying to decide these N functions to a setting to somehow I setting where we have to decide only actions ok.

So, actions, but actions and from there from those actions somehow deduce what the function should be? So, all of that will happen in a subsequently as we go further down in this course ok. So, now, to summarize what, we said in this example. Let me, write for you a general dynamic optimization problem or a general dynamic programming problem and it is main constituents ok.

So, the main constituents of a dynamic programming problem involve, first a state the first component is state of the system ok. So, the state of the system simply is your x_k in this case, which was the inventory level. So, this ok so, it is it will be denoted by x_k it was in the previous problem it was the inventory level. It captures, it is whatever is needed, it is whatever is needed to capture the configuration of the system.

It is some description that we have of the configuration of the system that we are that we are dealing with. In this case since the, system we were in the example we saw the system we were dealing with was the inventory in a shop. So, it was enough for us to keep track of just simply the amount of inventory present at the beginning of the time period right.

So, in so, as a result of this the, we took that as the state of the system. Now, the state of the system is a part of problem modeling, what exactly want to define as the state is a it is a bit of a there are usually more than one ways of defining the state, but try to always keep the state as simple as possible.

(Refer Slide Time: 12:37)

Dynamic Programming

Inventory control:
 Shop keeper needs to decide the amount of qty of a certain item to be ordered over some time instants $0, \dots, N$ (N time instants)
 goal is to meet the demand
 cost optimize costs
 x_k = stock available at the beginning of the k^{th} period.

Diagram showing time instants $0, 1, 2, 3, \dots, N$.
 u_k = stock ordered at the beginning of the k^{th} period, Delivered immediately
 w_k = demand for the item during the k^{th} period.

Equation: $x_{k+1} = x_k + u_k - w_k$
 Annotations: x_k is 'stock @ begin of k ', u_k is 'Stock added over the k^{th} period', and w_k is 'consumed over the k^{th} period'.

Assumption: $w_k \dots w_{k+1} \rightarrow$ independent random variable "Noise"

For instance in the previous problem we could have, because we were talking of inventory at I am trying to decide the stock over these n time periods. We could have taken the state as the entire history of you know the stock levels up till time k .

Now, that entire history could be taken, but it has no bearing on trying to decide it does not help in you know in trying to decide what we are, what we are actually looking for which is the amount of stock that is to be ordered.

Because, that the amount of for in order to decide, what the stock that we have to order, it is enough to simply know what the level of the stock is at the latest time period, not and the entire history is actually not relevant for that.

So, as a result we have picked a parsimonious definition of the state, we have taken this state as simply the current the level of stock at the current time period o_k . So, now, the state itself, state evolution or dynamics. Now the state itself evolves based on the action that you would take and the noise that, comes up in the system.

So, the way we express this is that we write x_{k+1} as some f_k of x_k comma, u_k comma, w_k right. So, here this is the action at time period k , this here is the noise, this is the noise at time period k alright. And, so, as a function of the action that you plan to take at time period k and noise that evolves at time period k , and the previous, and the state at time period k , you get the next state o_k , there is state at the next time period right.

Now, because this w_k is random this sequence will be random alright. So, these sequence of states will be random it is not something whose value you would know at the beginning of the time, at the beginning of the problem o_k . And, so, as a result the actions that you would end up taking would also be random, because they would be as would be chosen as a function of the information that you have.

We so, this here is these are this w_k is noise; usually we take w_k as independent random variables. In many cases also identically distributed or 0 mean and so on, but that depends on the problem. Usually, we take these as independent random variables alright. And, they are independent of each other and they their distribution cannot be decided you know as a function of as a function of using through your actions o_k .

Now, there would usually be a control constraint on actions. So, that constraint on actions is in this in our case was, in our case was say u_k greater than equal to 0 o_k . So, for example, that was our constraint, but more generally it could be any constraint and moreover that

constraint could depend also on the state that you are in. Say for example, it could be something as something like this u_k belongs to capital U_k of x_k ok.

This would be the most general way of writing the constraint. So, as a function of the state x_k that you are in the kind of the actions that you can potentially take have to be in this set capital U_k of x_k ok, alright. And, the additive cost form. The cost we would the cost form is denoted in this sort of way, here is your terminal cost g_N of x_N remember we are taking actions at starting at time 0 till time $N-1$.

The problem ends at time N . So, at time N whatever is the state based on that you incur a terminal cause g_N of x_N , we are not taking any further actions at time N . So, although I am saying I have said that we take actions at N time instance, the N time instance themselves are denoted 0 to $N-1$ alright.

So, this is so, the entire cost is therefore, g_N of x_N , which is your terminal cost plus g_k of x_k , u_k sometimes we can we also include w_k here, but it does not matter d_k you it can this inclusion here is optional it is you can simply write g_k comma u_k , that is also without loss of generality alright.

So, as you can see we are in our in the earlier in the problem above g_N of x_N was this was simply R of x_N and g_k of x_k , u_k , w_k this was r of x_k plus c of u_k ok alright. Now, ok. So, having formulated this the basic problem, then are the, are goal becomes the following, that we would like to decide, we would like to decide a policy.

(Refer Slide Time: 19:40)

$\pi = (\mu_0 \dots \mu_{N-1})$
 s.t. $u_k = \mu_k(x_k) \quad \forall x_k, k$
 admissible policy: $\mu_k(x_k) \in U_k(x_k) \quad \forall x_k, k$
 Initial state $x_0, \pi(\mu_0 \dots \mu_{N-1})$
 $x_{k+1} = f_k(x_k, \mu_k(x_k), \omega_k)$
 $E \left[g_0(x_0) + \sum_{k=1}^{N-1} g_k(x_k, \mu_k(x_k), \omega_k) \right] = \text{depends on } \pi \text{ \& \& on } x_0$
 $= J_\pi(x_0)$
 $= \text{cost incurred by policy } \pi, \text{ starting from state } x_0$

Find $\pi^* \in \Pi$, ($\Pi = \text{set of admissible policies}$)
 s.t.
 $J_{\pi^*}(x_0) = \min_{\pi \in \Pi} J_\pi(x_0)$
 $J^*(x_0) = J_{\pi^*}(x_0)$
 $= \text{"optimal value function"}$
 $= \text{optimal cost as a function of initial state}$

So, we would like to decide a π , which is denoted by μ_0 to μ_{N-1} , what we would a policy μ_0 to μ_{N-1} such that, it would map the state to control actions or 2 actions that you want to take ok. So, u_k would be equal to μ_k of x_k , for all values of x_k for all values of k right. And, now moreover this now what kind of policies are admissible?

Now, remember since, we mentioned, we have control; we have constraints here on control actions, that control actions needed have to satisfy this particular constraint. The kind of policies that are admissible are an admissible policy is only one, where an admissible policy is one, which is in which u_k or in which as μ_k of x_k always belongs to U of u capital U_k of x_k . So, this is to for all comma and all k .

So, an admissible policy is one which satisfies the constraints the control constraints that we have. So, now, given the given so, how does the problem begin? The problem begins with an

initial state. So, you are given an initial state x_0 . The it is say the inventory level at the start, at the start of the time horizons, say an initial state.

You given that initial state, you would then and given the initial given the, a choice of a policy, given an initial state and a policy μ_0 to μ_{N-1} , given the initial state and a policy, what how does the how does this system behave? Well, it takes the initial state uses the policy at times 0 to decide the action to be taken at time 0 right?

So, it takes the policy at time 0 to decide the action to be taken at time 0, which is that would be u_0 , u_0 , then feeds into the state dynamics and what we get from there is that you get this state at time 1. So, you get more generally x_{k+1} at time k a time x_{k+1} emerges as a function f_k of x_k and u_k , where u_k itself is just simply μ_k of x_k alright and the noise w_k right.

So, because you are choosing now the action as a function of the state x_k , because u_k itself is being chosen as a function as a function of the state, you can now substitute this here. And, what you find is that the next state comes up is defined through the dynamics and the policy that you have chosen, which is μ_k right.

So, after the substitution essentially, what I what one can do is you can do this substitution everywhere, in fact, you can do this substitution even in the cost function, if your cost function was g_N of; g_N of x_N plus. So, g_N of x_N plus g_k of u_k w_k . So, I have replaced the u_k there by simply the by simply μ_k of x_k right. So, as you can see the cost function now depends on two things.

So, this depends on π which is your policy and depends also on the initial state x_0 . It depends on where you are starting, because after all the state evolution that you that your system would see, would now get determined by the dynamics and the policy that has been stated.

So, once I fixed my initial state which is the state where I am starting from, and the kind of policy that I have chosen the noise distribution is what will determine the evolution of your of the state sequence x_k . And, so, what we have therefore, is the expected cost of the

associated and evaluated over this particular noise sequence, and that depends as you can see on π as well as x_0 .

So, we denote this we denote this therefore, by J_{π} of x_0 . So, this is therefore, this, what is this quantity this is the cost incurred by policy π starting from state, starting from state x_0 , this that is what this would be. So, what is therefore, the problem? The problem then is to find policy π^* in say a space Π . Now, what this Π here is the set of admissible policies, admissible set of admissible policies, sorry J_{π^*} .

Such so, find a π^* , which is in the set of admissible policies such that the cost of the cost that you incurred by π^* is the least among the cost of all policies, in this set of admissible policies ok. Now, the cost the notice that, this is you it appears like this particular problem is actually a problem that we would now need to solve for each value of x_0 , because we have not gotten rid of the initial state as yet.

So, we would have a cost for each initial state. And, it appears a prime of see that the π^* which is the optimal policy would vary could vary with x_0 . Because, I have just fixed 1 x_0 and found a π^* ; however, the nature of the problem is such that it turns out at you know this it is typically possible to find a π^* , that does not depend on x_0 itself.

So, you can choose any token x_0 to start with and the policy that you would find you end up finding a policy that would work for any initial state ok. So, that is how that is how it often happens is it is because the policy itself is something that is a plan for every possible state. So, it does not it is not it is not merely for this any the specific state that you are looking for it you usually get the policy for every possible initial every possible initial state.

So, this particular quantity which is the J_{π^*} this is also denoted by J^* of x_0 , this sometimes denoted by J^* of x_0 that is simply another way of writing the cost of the optimal policy. This has a name it is what is called the optimal value function, optimal value

function. So, the optimal value function is simply a function that tells you what the optimal value of the dynamic optimization is going to be as a function of the initial state ok.

This is called the optimal value function or optimal cost function ok, this optimal cost as a function of initial state ok. So, this is so, what we will now what we will do in the next class is find ways of computing this particular the optimal value the or the optimal value function ok.

So, I will pause here and we will resume again in the next class.