

4. OPTIMAL CONTROL

4.1. The Problem, and Heuristics. Let $y(t)$ be a quantity (e.g. production rate at a factory, portfolio value, etc.), that satisfies an ordinary differential equation:

$$\partial_t y(t) = F(y(t), \alpha_1(t), \alpha_2(t), \dots).$$

Above, $\alpha_1, \alpha_2, \dots$ are parameters that may depend on time and which influence F .

Example:

$$(4.1) \quad \partial_t y(t) = ry(t) - \alpha(t), \quad y(0) = x.$$

This can of course be solved, and the solution

$$y(t) = e^{rt} \left(x - \int_0^t e^{-rs} \alpha(s) ds \right)$$

depends on $\alpha(s)$. Let now h be a utility function, which can depend on y and α ; for example, the utility function could be lower when α is large (cost of controlling something), and higher when y is high.

The *problem of optimal control* is this: find *control parameters* $\alpha_1(s), \alpha_2(s), \dots$ such that h is maximal when averaged over time. In other words, we have to find

$$(4.2) \quad u(x, t) = \max_{\alpha} \int_t^T h(y(s), \alpha(s)) ds,$$

where initially (at time t), the value of $y(t) = x$. Importantly, $y(s)$ itself depends both on x and on the control $\alpha(s)$ through the ODE, e.g. (4.1).

There are many variants of the optimal control problem. One can let h depend explicitly on time. The most important instance of this is *discounting*, when we have $e^{-rs} h(y(s), \alpha(s))$ under the integral. Or, we can have a final time utility $g(y(T))$ added to u .

The example (4.1) corresponds to an optimal consumption problem. There $y(t)$ is the wealth at time t , $\alpha(t)$ is the rate of consumption, y_0 is the initial wealth; wealth that has not been consumed earns interest at rate r . The control problem is to maximize

$$(4.3) \quad u(x, t) = \int_0^T e^{-\rho s} h(\alpha(s)) ds.$$

h is usually a concave function, such as $h(\alpha) = \alpha^{1/2}$ or, more generally, $h(\alpha) = \alpha^\gamma$ with $0 < \gamma < 1$. The concavity has a nice interpretation: while it may make you happier to spend more money per day, it will not make you twice as happy; or, the less you are used to get by with, the more you will appreciate a bit more money to spend. T is the final time by which you should have spent it all, wealth that is left at that time will not benefit you.

It is a bit confusing that (4.3) seems not to depend on y at all. But in fact it does. Firstly, x is the starting point of $y(s)$, $s \geq t$, and secondly, we have the condition that $y(s) \geq 0$ for all s , i.e. we adhere to the old-fashioned strategy that you cannot

spend more than you have. This in turn restricts the controls α , since a control that is too large at the beginning will have to be zero after $y(s)$ hits zero, which happens sooner if α is large. Also, of course $\alpha \geq 0$, i.e. the consumption should be non-negative.

In the following we will give a general strategy to solve optimal control problems. The problem of optimal consumption will then be done in an exercise.

4.2. Solving the optimal control problem: dynamic programming. A general optimal control problem consists of the following parts:

The ODE:

$$\partial_s \mathbf{y}(s) = F(\mathbf{y}(s), \boldsymbol{\alpha}(s)), \quad \mathbf{y}(t) = \mathbf{x},$$

with $\mathbf{y} : \mathbb{R} \rightarrow \mathbb{R}^n$, $\boldsymbol{\alpha} : \mathbb{R} \rightarrow \mathbb{R}^m$ and $F : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$.

The control constraint: $\boldsymbol{\alpha}(s) \in A \subset \mathbb{R}^m$ for all s . This constraint may not always be present, and A may even depend on time.

The state constraint: $y(t) \in Y \subset \mathbb{R}^n$ for all t . This is usually a difficult constraint. The way we deal with it in the optimal consumption problem is to ignore it and verify afterwards that the solution fulfils it. This is not always possible.

The control problem: Let $h : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ be the running utility function, and $g : \mathbb{R}^n \rightarrow \mathbb{R}$ be the final utility function. The control problem is to find the optimal value function

$$(4.4) \quad u(\mathbf{x}, t) = \max_{\boldsymbol{\alpha} \in A, \mathbf{y} \in Y} \left(\int_t^T h(\mathbf{y}_{\boldsymbol{\alpha}, \mathbf{x}}(s), \boldsymbol{\alpha}(s)) ds + g(\mathbf{y}_{\boldsymbol{\alpha}, \mathbf{x}}(T)) \right).$$

The subscript on the function \mathbf{y} is there to remind us that it actually depends on the starting point and the control, but it will be often omitted in the notation.

The idea when solving this is to work backwards from the final time T . We ignore the state constraint for now.

- 1) For $t = T$, it is clear that there is nothing to control, and $u(\mathbf{x}, T) = g(\mathbf{x})$.
- 2) Let δt be a very small time step. Let us consider $t = T - \delta t$. We try to optimize (4.4) by using a control $\boldsymbol{\alpha}$ that is constant on $[T - \delta t, T]$. In the limit when $\delta t \rightarrow 0$ there is hope that this is good enough. We approximate the solution of the ODE by Taylor expansion:

$$(4.5) \quad \mathbf{y}(t + s) \approx \mathbf{y}(t) + \partial_t \mathbf{y}(t)s = \mathbf{y}(t) + F(\mathbf{y}(t), \boldsymbol{\alpha}(t))s$$

for $s \in [0, \delta t]$. We also approximate $h(\mathbf{y}(t + s), \boldsymbol{\alpha}(t + s)) \approx h(\mathbf{y}(t), \boldsymbol{\alpha}(t))$ (zero order Taylor), and

$$(4.6) \quad \int_{t=T-\delta t}^T h(\mathbf{y}(s), \boldsymbol{\alpha}(s)) ds \approx (T - t)h(\mathbf{y}(t), \boldsymbol{\alpha}(t)) = \delta t h(\mathbf{y}(t), \boldsymbol{\alpha}(t)).$$

Thus

$$\begin{aligned}
 (4.7) \quad u(\mathbf{x}, T - \delta t) &= \max_{\boldsymbol{\alpha}(s)} \left(\int_{T-\delta t}^T h(\mathbf{y}(s), \boldsymbol{\alpha}(s)) \, ds + g(\mathbf{y}(T)) \right) \\
 &\approx \max_{\boldsymbol{\alpha}} (h(\mathbf{x}, \boldsymbol{\alpha})\delta t + g(\mathbf{x} + F(\mathbf{x}, \boldsymbol{\alpha}))\delta t) \\
 &= \max_{\boldsymbol{\alpha}} \left(h(\mathbf{x}, \boldsymbol{\alpha})\delta t + u(\mathbf{x} + F(\mathbf{x}, \boldsymbol{\alpha}), T)\delta t \right).
 \end{aligned}$$

So, given that we know the optimal value function at time T (which we do, it is g), we also know it at time $T - \delta t$. Indeed, we get it by maximizing the known function on the right hand side of (4.7) over $\boldsymbol{\alpha}$, for all \mathbf{x} , e.g. by differentiating and looking for zeros. This even gives a numerical scheme for finding the optimal value function and the optimal control, but that scheme is very impractical in higher dimensions. However, the above reasoning will soon lead to a PDE that one can treat easier, at least numerically.

3) Step 2 is now repeated: we know the value function at time $T - \delta t$, so we get it (by step 2) for time $T - 2\delta t$, and so on until we reach the initial time.

The insight that knowing the optimal utility at a time t_1 helps us determine it at an earlier time is important. Indeed, the equation

$$u(\mathbf{x}, t) = \max_{\boldsymbol{\alpha}(s)} \int_t^{t_1} h(\mathbf{y}(s), \boldsymbol{\alpha}(s)) \, ds + u(\mathbf{y}(t_1), t_1),$$

for $t_1 > t$, can be seen to be true in a similar way as above. This equation is called *dynamic programming principle* because it is at the basis of the algorithm given above.

4.3. Hamilton-Jacobi-Bellman (HJB) equation. Let us replace T by $s + \delta t$ in (4.7). We then have

$$u(\mathbf{x}, s) = \max_{\boldsymbol{\alpha}} \left(h(\mathbf{x}, \boldsymbol{\alpha})\delta t + u(\mathbf{x} + F(\mathbf{x}, \boldsymbol{\alpha})\delta t, s + \delta t) \right).$$

We now do a Taylor expansion of the function $\delta t \mapsto u(\mathbf{x} + F(\mathbf{x}, \boldsymbol{\alpha})\delta t, s + \delta t)$, and find

$$u(\mathbf{x}, s) \approx \max_{\boldsymbol{\alpha}} \left(h(\mathbf{x}, \boldsymbol{\alpha})\delta t + u(\mathbf{x}, s) + \nabla u(\mathbf{x}, s) \cdot F(\mathbf{x}, \boldsymbol{\alpha})\delta t + \partial_t u(\mathbf{x}, s)\delta t \right).$$

We can now subtract $u(\mathbf{x}, s)$ from both sides. The \approx sign really means that there are terms of order $(\delta t)^2$ that we ignored. Dividing by δt and sending $\delta t \rightarrow 0$ gives the *Hamilton-Jacobi-Bellman equation* (HJB equation):

$$(4.8) \quad \partial_s u(\mathbf{x}, s) + \max_{\boldsymbol{\alpha} \in A} (\nabla u(\mathbf{x}, s) \cdot F(\mathbf{x}, \boldsymbol{\alpha}) + h(\mathbf{x}, \boldsymbol{\alpha})) = 0.$$

There is a special notation that is often used for the HJB equation: we define the *Hamiltonian* $H(\mathbf{p}, \mathbf{x})$ as

$$(4.9) \quad H(\mathbf{p}, \mathbf{x}) = \max_{\boldsymbol{\alpha} \in A} (\mathbf{p} \cdot F(\mathbf{x}, \boldsymbol{\alpha}) + h(\mathbf{x}, \boldsymbol{\alpha})),$$

with $\mathbf{x}, \mathbf{p} \in \mathbb{R}^n$. Then (4.8) reads

$$(4.10) \quad \partial_s u + H(\nabla u, \mathbf{x}) = 0.$$

It is clear from (4.10) that for solving a control problem, we have to maximize over $\boldsymbol{\alpha}$ in (4.9) first, and then solve (4.10). In particular, to find the optimal control at space point \mathbf{x} and 'momentum' \mathbf{p} , we need not know the solution u of (4.10). However, of course since we then replace \mathbf{p} by ∇u , the solution will be fed back into H .

We have now found the equation for the optimal value function: It is the solution of (4.8) with final data $u(\mathbf{x}, T) = g(\mathbf{x})$. To know the value function is good enough for e.g. option pricing, when we allow the buyer to change some control parameter: we now know how much an investor can make at most, if they play the game in an optimal way. This should then be the fair price of the option.

But what about actually finding the optimal strategy $\boldsymbol{\alpha}(s)$ for $t \leq s \leq T$? For this, we have to plug the solution $u(\mathbf{x}, s)$ back into the second term of (4.8), and find the argmax for each time. More precisely: from (4.9) we obtain both $H(\mathbf{p}, \mathbf{x})$ for all \mathbf{p} and all \mathbf{x} , and $\boldsymbol{\alpha}_*(\mathbf{p}, \mathbf{x})$ as the argmax of the right hand side. Now we solve (4.10) with the H that we just obtained, and the correct final condition. Once this is done, we know $u(\mathbf{x}, s)$ for all \mathbf{x} and all s , and thus also $\nabla u(\mathbf{x}, s)$. The ODE for the optimally controlled \mathbf{y} is then

$$\partial_s \mathbf{y}(s) = F(\boldsymbol{\alpha}_*(\nabla u(\mathbf{y}(s), s), \mathbf{y}(s)), s),$$

where now the right hand side only contains known functions of s and $\mathbf{y}(s)$. Once we have obtained the solution $\mathbf{y}_*(s)$ to this ODE (with initial condition $\mathbf{y}(t) = \mathbf{x}$), we can finally determine the purely time dependent optimal control for the case when $\mathbf{y}(t) = \mathbf{x}$. It is given by $\boldsymbol{\alpha}_*(\nabla u(\mathbf{y}_*(s), s), \mathbf{y}_*(s))$.

The thing we now have to check is that $\mathbf{y}(t)$ actually fulfils the state constraints. If it does, everything is fine. But if it does not, we have to go back to the dynamic programming approach and incorporate this at each step. There is no nice theory for this, and we will not do it.

We have derived the HJB equation using a lot of heuristic and non-rigorous steps. But once we have it, we can actually prove that it gives the optimal value function. This is the content of the following

Theorem: *Consider the general optimal control problem introduced at the beginning of this subsection. Assume that $w(\mathbf{x}, t)$ solves the HJB equation (4.8), with final condition $w(\mathbf{x}, T) = g(\mathbf{x})$. Assume also that \mathbf{y} derived from w in the way discussed above fulfils the state constraint. Then $w(\mathbf{x}, t) = u(\mathbf{x}, t)$, where u is defined as the solution of the control problem.*

Proof. Assume that w is a solution of the HJB equation with final condition $w(\mathbf{x}, T) = g(\mathbf{x})$. Let $\boldsymbol{\alpha}_0(s)$ be an arbitrary control, and let $\mathbf{y}_0(s)$ be the solution of the controlled ODE with start $\mathbf{y}_0(t) = \mathbf{x}$ and control $\boldsymbol{\alpha}_0(s)$. Let us investigate

the function $s \mapsto w(\mathbf{y}_0(s), s)$. We find

$$\begin{aligned} \frac{d}{ds}w(\mathbf{y}(s), s) &= \partial_s w(\mathbf{y}(s), s) + \nabla w(\mathbf{y}(s), s) \cdot \partial_s \mathbf{y}(s) \\ &= \underbrace{\partial_s w(\mathbf{y}(s), s) + \nabla w(\mathbf{y}(s), s) \cdot F(\mathbf{y}(s), \boldsymbol{\alpha}(s)) + h(\mathbf{y}(s), \boldsymbol{\alpha}(s))}_{(*)} - h(\mathbf{y}(s), \boldsymbol{\alpha}(s)). \end{aligned}$$

Now for an arbitrary control, the expression $(*)$ is smaller or equal to zero. This is because, by the HJB equation that w satisfies, when we take the maximum over all possible values of $\boldsymbol{\alpha}(s)$ we get $(*) = 0$, so any other $\boldsymbol{\alpha}$ will give less. We can then integrate the last equation from t to T and get

$$w(\mathbf{y}(T), T) - w(\mathbf{y}(t), t) \leq - \int_t^T h(\mathbf{y}(s), \boldsymbol{\alpha}(s)) ds,$$

and using $w(\mathbf{y}, T) = g(\mathbf{y})$ as well as $\mathbf{y}(t) = \mathbf{x}$, we find

$$g(\mathbf{y}(T)) + \int_t^T h(\mathbf{y}(s), \boldsymbol{\alpha}(s)) ds \leq w(\mathbf{x}, t).$$

So any arbitrary control can make the value function at most as large as $w(\mathbf{x}, t)$, and maximizing over the possible controls we find $u(\mathbf{x}, t) \leq w(\mathbf{x}, t)$. What remains to show is that $w(\mathbf{x}, t)$ is indeed a value function; so far we only know it to be the solution of a HJB equation. To see that w is indeed a value function, notice that when we choose the feedback control that we get from w (see the discussion before the theorem), then $(*)$ above is identically zero (independent of whatever $\mathbf{y}(s)$ happens to be). Thus, again by integrating, we see that w is indeed a value function. \square