

Today's central concepts

Two phases in solving an MDP.

i, Modeling. Define all the quantities (state-space, transitions, etc.)

ii, Use algorithm to obtain optimal policy

The second step is usually rather mechanical if the first is done well.

⇒ Don't rush step i, !!!

Finite-horizon, $T < \infty$:

(Stochastic) Dynamic Programming / Backward Induction:

• $u_T^*(s) = r_T(s)$

← the terminal reward-to-go

• $u_t^*(s) = \max_{a \in A_s} \left\{ r_t(s, a) + \sum_{s' \in S} p_t(s'|s, a) u_{t+1}^*(s') \right\}$

↖ optimal expected reward-to-go at time t

↖ optimal action

↖ immediate reward

↖ expected reward-to-go if we pick action a

• $a_t^*(s) \in \arg \max_{a \in A_s} \left\{ \dots \right\}$

if there are multiple optimal actions, then it's arbitrary which one we pick

Discounted ∞ -horizon:

- Optimality conditions / Bellman's equation:

$$u^*(s) = \max_{a \in A_s} \left\{ r(s, a) + \lambda \sum_{s' \in S} p(s'|s, a) u^*(s') \right\}$$

- Value iteration (VI) algorithm:

$$u_{k+1}(s) = \max_{a \in A_s} \left\{ r(s, a) + \lambda \sum_{s' \in S} p(s'|s, a) u_k(s') \right\}$$

will converge to $u^*(s)$ for any initial condition $u_0(s)$.

Note, this is essentially DP with a longer and longer horizon.

- Policy iteration (PI) algorithm:

0, Guess a policy π_0

i, (Policy evaluation)

compute the value of policy π_k by solving:

$$u^{\pi_k}(s) = r(s, \pi_k(s)) + \lambda \sum_{s' \in S} p(s'|s, \pi_k(s)) u^{\pi_k}(s')$$

ii, (Policy improvement)

Update the policy:

$$\pi_{k+1}(s) \in \arg \max_{a \in A_s} \left\{ r(s, a) + \lambda \sum_{s' \in S} p(s'|s, a) u^{\pi_k}(s') \right\}$$

Stop if $\pi_k = \pi_{k+1}$.

Remark:

For the systems we consider, deterministic policies are sufficient for optimality.

That is, a policy π is a function that takes a state $s \in \mathcal{S}$ and maps it into an action $a \in \mathcal{A}_s$.

Remark:

These algorithms are also valid for other objective functions. There are variations of these algorithms that are more efficient.

Remark:

It is possible to solve MDPs via linear programming, which allows constraints to be enforced.

Ex 3.10 | need to sell apartment within N days.

We receive an offer w_t every evening that we need to accept or reject the next day. Every offer is a multiple of 10 000 SEK, and they are i.i.d., positive and bounded. Fixed interest rate $\rho > 0$ once we sell. Compute an optimal policy. ^{daily}

Solution:

We derived the MDP-model in the previous session.

State-space:

Let's assume we always talk in terms of multiples of 10 000 SEK, and that the maximum bid is 10 000. b_{max} SEK. Then:

$$S = \{0, 1, \dots, b_{max}\} \cup \{sold\}$$

↖ current bid under consideration

Actions:

$A_b = \{R, \alpha\}$: the actions available when we are considering a bid b
↙ reject ↘ accept

$A_{sold} = \{C\}$: ————— " —————
↙ continue we have sold.

Time-horizon and objective:

Finite-horizon, N :

$$\mathbb{E} \left\{ \sum_{t=0}^{N-1} r_t(s_t, a_t) + r_N(s_N) \right\}$$

Rewards:

Terminal:

$$\bullet r_N(\text{sold}) = 0$$

$$\bullet r_N(b) = b$$

we are forced to sell at time N if we haven't sold

Non-terminal:

$$\bullet r_t(s=b, a=\mathcal{A}) = b(1+\rho)^{N-t}$$

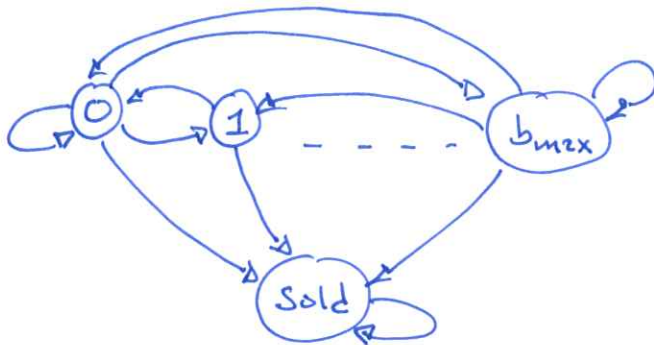
the current offer

the future interest that we will collect

$$\bullet r_t(s=b, a=\mathcal{R}) = 0$$

$$\bullet r_t(s=\text{Sold}, a=\mathcal{C}) = 0$$

Transitions:



Note:
we could take A_{sold} to be $\{\mathcal{A}, \mathcal{R}\}$ as we did last time. In that case, the choice of action in the policy will be arbitrary.

$$\bullet P_t(s'=b' | s=b, a=\mathcal{R}) = \Pr\{\omega_t = b'\}$$

if we reject, tomorrow's offer is drawn i.i.d.

$$\bullet P_t(s'=\text{Sold} | s=b, a=\mathcal{A}) = 1$$

we decide to sell

$$\bullet P_t(s'=\text{Sold} | s=\text{Sold}, a=\mathcal{C}) = 1$$

if we have sold

To solve the MDP, recall the backward induction:

- $u_N^*(s_N) = r_N(s_N)$
- $u_t^*(s_t) = \max_{a \in A_{s_t}} \left\{ r_t(s_t, a) + \sum_{s' \in \mathcal{S}'} P_t(s' | s_t, a) u_{t+1}^*(s') \right\}$
- $a_t^*(s_t) \in \arg \max_{a \in A_{s_t}} \left\{ \text{---} \parallel \text{---} \right\}$

Base case:

- $u_N^*(\text{sold}) = r_N(\text{sold}) = 0$
- $u_N^*(b) = r_N(b) = b$

Recursion:

For $s_t = \text{sold}$:

$$\begin{aligned} u_t^*(\text{sold}) &= \max_{a \in A_{\text{sold}}} \left\{ \underbrace{r_t(\text{sold}, a)}_{= \{d\}} + \underbrace{\sum_{s' \in \mathcal{S}'} P_t(s' | \text{sold}, a)}_{= 0} \underbrace{u_{t+1}^*(s')}_{= \begin{cases} 1 & \text{if } s' = \text{sold} \\ 0 & \text{o.w.} \end{cases}} \right\} \\ &= 0 + 1 \cdot u_{t+1}^*(\text{sold}) \\ &= u_{t+1}^*(\text{sold}). \end{aligned}$$

By induction, we conclude that

$$u_t^*(\text{sold}) = u_{t+1}^*(\text{sold}) = \dots = u_N^*(\text{sold}) = 0. \quad (*)$$

For $s_t = b$:

$$u_t^*(b) = \max_{a \in A_b} \left\{ r_t(b, a) + \sum_{s' \in S'} P_t(s' | b, a) u_{t+1}^*(s') \right\}$$

this is = 0 for $s' = \text{Sold}$ according to (*)

$$= \max_{a \in \{\text{A}, \text{R}\}} \left\{ r_t(b, a) + \sum_{i=0}^{b_{\max}} P_t(i | b, a) u_{t+1}^*(i) \right\} \quad (**)$$

Note: This is a sum over $S \setminus \{\text{Sold}\} = \{0, 1, \dots, b_{\max}\}$.

It's easiest to evaluate the two actions separately:

If $a = \text{A}$:

$$\underbrace{r_t(b, \text{A})}_{= b(1+p)^{n-t}} + \underbrace{\sum_{i=0}^{b_{\max}} P_t(i | b, \text{A}) u_{t+1}^*(i)}_{= 0 \text{ for all } i} = b(1+p)^{n-t}$$

If $a = \text{R}$:

$$\underbrace{r_t(b, \text{R})}_{= 0} + \underbrace{\sum_{i=0}^{b_{\max}} P_t(i | b, \text{R}) u_{t+1}^*(i)}_{= \Pr\{\omega_t = i\}}$$

$$\sum_{i=0}^{b_{\max}} \Pr\{\omega_t = i\} u_{t+1}^*(i) = \mathbb{E}_{\omega} \left\{ u_{t+1}^*(\omega) \right\}$$

Taken together in (**), we obtain:

$$u_t^*(b) = \max \left\{ \underbrace{b(1+p)^{N-t}}_A, \underbrace{\sum_{i=0}^{b_{max}} Pr\{\omega_t=i\} u_{t+1}^*(i)}_R \right\}$$

$$= \max \left\{ \underbrace{b}_A, \frac{\sum_{i=0}^{b_{max}} Pr\{\omega_t=i\} u_{t+1}^*(i)}{(1+p)^{N-t}} \right\} \times (1+p)^{N-t}$$

def. = α_t

$$\text{def.} \\ = \max \left\{ \underbrace{b}_A, \underbrace{\alpha_t}_R \right\} \times (1+p)^{N-t}$$

Hence, the optimal policy is:

If we receive an offer b at time t , then we should:

A (cept) if $b > \alpha_t$,
 R (ject) if $b \leq \alpha_t$,

if $b = \alpha_t$, the action is arbitrary.

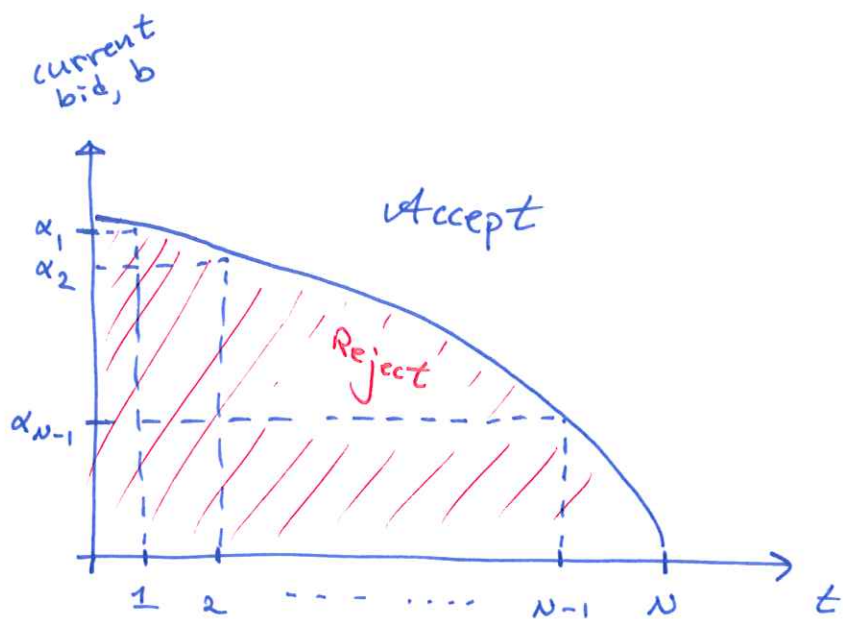
where

$$\alpha_t = \frac{\sum_{i=0}^{b_{max}} Pr\{\omega_t=i\} u_{t+1}^*(i)}{(1+p)^{N-t}} = \mathbb{E} \left\{ u_{t+1}^*(\omega) \right\} \cdot (1+p)^{t-N}$$

↗
 expected reward-to-go

↑
 discount by interest

This is a time-dependent threshold policy:



The more time we have ahead of us, the more greedy we can be with the bid we choose to accept.

Ex 3.2 | At every time a job, with salary from the set $\{w_1, \dots, w_n\}$, is offered and has to be accepted or rejected. The offers are iid. An unemployment compensation c is given at every time (if unemployed). Assume future discount factor λ .

- a, Show that there is a threshold \bar{w} over which offers should be accepted. Characterize \bar{w} .
- b, Assume the worker is fired from job i w.p. p_i . Show that a, holds if $p_i = p$ for all i .

Solution:

State-space:

$$S = \{s_1, \dots, s_n\} \cup \{s'_1, \dots, s'_n\}$$

s_i : not employed, considering an offer from job i with salary w_i

s'_i : employed at job i with salary w_i

Actions:

$$A_{s_i} = \{A, R\}$$

↗ accept
↖ reject

$$A_{s'_i} = \{C\}$$

↖ continue working

Time-horizon and objective:

∞ -horizon, discounted:

$$E \left\{ \sum_{t=0}^{\infty} \lambda^t r_t(s_t, a_t) \right\}$$

Remark:

How do we interpret λ ?

One interpretation is that we are maximizing

$$\mathbb{E} \left\{ \sum_{t=0}^T r_t(s_t, z_t) \right\}$$

for a random T : $\Pr\{T=k\} = \lambda^{k-1}(1-\lambda)$.

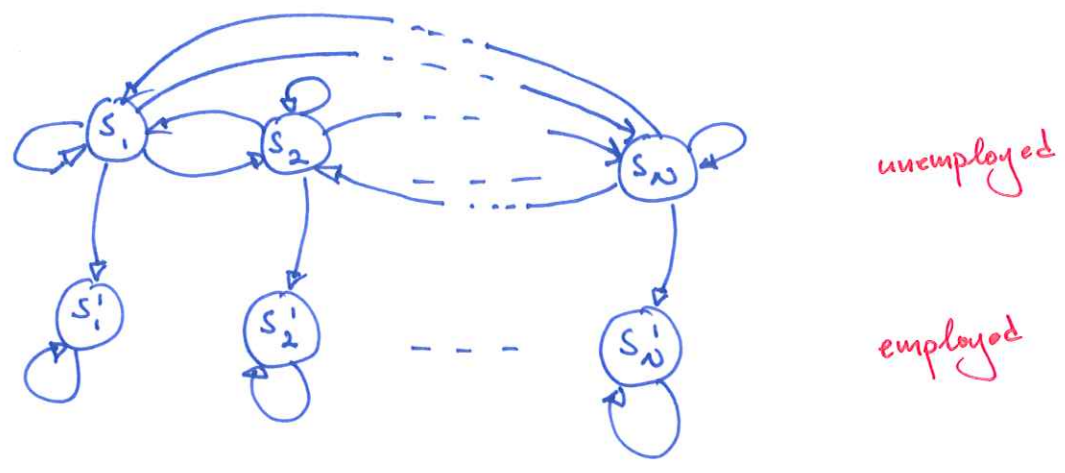
$1-\lambda$ can be interpreted as the worker's risk of dying each time unit. The worker's expected life-time is $\mathbb{E}\{T\} = \frac{1}{1-\lambda}$ time units. The worker tries to maximize the total money earned, subject to knowing that she might die.

Rewards:

- $r_t(s_i, \alpha) = w_i$ \leftarrow she accepts the offer from job i
- $r_t(s_i, R) = c$ \leftarrow she gets the unemployment compensation
- $r_t(s_i', C_i) = w_i$ \leftarrow she works at job i

Transitions:

let q_i denote the probability that the worker receives an offer from job i .



The non-zero transitions are:

- $P_t(s'_i | s'_i, C) = 1$
- $P_t(s'_i | s_i, A) = 1$
- $P_t(s_j | s_i, R) = q_j$

Bellman equation:

$$u^*(s) = \max_{a \in A_s} \left\{ r_t(s, a) + \lambda \sum_{s' \in S} p_t(s' | s, a) u^*(s') \right\}$$

Remark:
 Everything is stationary, so we drop the time indices on r_t and p_t .

Let's evaluate this for the different states!

If $s = s_i'$:

$$u^*(s_i') = \max_{a \in A_{s_i'}} \left\{ r(s_i', a) + \lambda \sum_{s' \in \mathcal{S}} P(s' | s_i', a) u^*(s') \right\}$$

$= \{c_i\}$
 $= \begin{cases} 1 & \text{if } s' = s_i' \\ 0 & \text{o.w.} \end{cases}$

$$= \underbrace{r(s_i', c_i)}_{= \omega_i} + \lambda \cdot 1 \cdot u^*(s_i')$$

$$= \omega_i + \lambda u^*(s_i')$$

$$\Rightarrow u^*(s_i') - \lambda u^*(s_i') = \omega_i$$

$$\Rightarrow u^*(s_i') = \frac{\omega_i}{1 - \lambda} \quad (*)$$

If $s = s_i$:

$$u^*(s_i) = \max_{a \in A_{s_i}} \left\{ r(s_i, a) + \lambda \sum_{s' \in \mathcal{S}} P(s' | s_i, a) u^*(s') \right\}$$

Let's consider the two actions separately:

If $a = \mathcal{A}$:

$$\underbrace{r(s_i, \mathcal{A})}_{= \omega_i} + \lambda \sum_{s' \in \mathcal{S}} P(s' | s_i, \mathcal{A}) \underbrace{u^*(s')}_{= \begin{cases} 1 & \text{if } s' = s_i' \\ 0 & \text{o.w.} \end{cases}} = \omega_i + \lambda u^*(s_i')$$

If $\lambda = R$:

$$\underbrace{r(s_i, R)}_{= c} + \lambda \underbrace{\sum_{s' \in S'} p(s' | s_i, R) u^*(s')}_{= \begin{cases} q_j & \text{if } s' = s_j \\ 0 & \text{o.w.} \end{cases}} = c + \lambda \sum_{j=0}^N q_j u^*(s_j)$$

Taken together, we have that

$$u^*(s_i) = \max \left\{ \underbrace{\omega_i}_{\checkmark} + \lambda u^*(s_i'), c + \lambda \sum_{j=0}^N \underbrace{q_j}_{R} u^*(s_j) \right\}$$

$$= \left\{ \begin{array}{l} \text{From (*) we know that} \\ u^*(s_i') = \frac{\omega_i}{1-\lambda} \end{array} \right\}$$

$$= \max \left\{ \underbrace{\omega_i}_{\checkmark} + \lambda \frac{\omega_i}{1-\lambda}, c + \lambda \sum_{j=0}^N \underbrace{q_j}_{R} u^*(s_j) \right\}$$

$$= \max \left\{ \underbrace{\frac{\omega_i}{1-\lambda}}_{\checkmark}, c + \lambda \sum_{j=0}^N \underbrace{q_j}_{R} u^*(s_j) \right\}$$

$$= \max \left\{ \underbrace{\omega_i}_{\checkmark}, \underbrace{(1-\lambda) \left[c + \lambda \sum_{j=0}^N q_j u^*(s_j) \right]}_{\stackrel{\text{def.}}{=} \bar{\omega}} \right\} \cdot \frac{1}{1-\lambda}$$

$$= \max \left\{ \underbrace{\omega_i}_{\checkmark}, \underbrace{\bar{\omega}}_R \right\} \cdot \frac{1}{1-\lambda}$$

Since $\bar{\omega}$ is a constant (does not depend on i), this is threshold policy:

if we are offered a salary ω_i , then we should:

Accept if $\omega_i > \bar{\omega}$,

Reject if $\omega_i \leq \bar{\omega}$,

← arbitrary if $\omega_i = \bar{\omega}$.

where

$$\bar{\omega} = (1-\alpha) \left[c + \lambda \sum_{j=0}^{\infty} \alpha^j u^*(s_j) \right].$$

Remark:

How does this help us — $u^*(s)$ is still unknown...? There are only $n+1$ possible threshold policies.

This limits the policy space we need to search over vastly. If n is reasonably large, we can simply compute u for each threshold policy: u^* is the best one.

Part b, : The worker gets fired w.p. p_i .
 We need to remodel the problem:

State-space:

$$S = \underbrace{\{s_1, \dots, s_N\}}_{\text{unemployed}} \cup \underbrace{\{s'_1, \dots, s'_N\}}_{\text{employed}}$$

Actions:

$$A_{s_i} = \{A, R\} \quad A_{s'_i} = \{C_i\}$$

Rewards:

- $r(s_i, A) = w_i$
- $r(s_i, R) = c$
- $r(s'_i, C_i) = w_i$

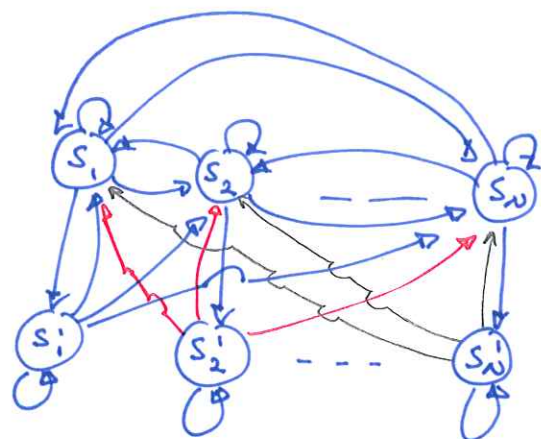
Time-horizon and objective:

∞ -horizon, discounted:

$$E \left\{ \sum_{t=0}^{\infty} \lambda^t r(s_t, a_t) \right\}$$

Transitions:

- $p(s'_i | s'_i, C_i) = 1 - p_i$
 not fired
- $p(s_j | s'_i, C_i) = p_i \cdot q_j$
 fired, and offered job j
- $p(s_j | s_i, R) = q_j$
 reject, and be offered job j
- $p(s'_i | s_i, A) = 1 - p_i$
 accept job i , and not fired this time period (e.g., month)
- $p(s_j | s_i, A) = p_i \cdot q_j$
 accept job i , but fired and then offered job j .



Bellman equation:

$$u^*(s) = \max_{a \in A_s} \left\{ r(s, a) + \sum_{s' \in S} p(s'|s, a) u^*(s') \right\}$$

Again, we consider the states separately:

If $s = s_i'$:

$$u^*(s_i') = \max_{a \in A_{s_i'}} \left\{ \underbrace{r(s_i', a)}_{= \{c_i\}} + \underbrace{\lambda \sum_{s' \in S} p(s'|s_i', a) u^*(s')}_{= \omega_i} \right\}$$

$$= \omega_i + \lambda \left[(1-p_i) u^*(s_i') + \sum_{j=0}^N p_i q_j u^*(s_j) \right]$$

If we solve for $u^*(s_i')$:

$$u^*(s_i') - \lambda(1-p_i) u^*(s_i') = \omega_i + \lambda p_i \sum_{j=0}^N q_j u^*(s_j) \implies$$

$$u^*(s_i') = \frac{\omega_i + \lambda p_i \sum_{j=0}^N q_j u^*(s_j)}{1 - \lambda(1-p_i)} \quad (*)$$

If $s = s_i$:

$$u^*(s_i) = \max_{a \in A_{s_i}} \left\{ \underbrace{r(s_i, a)}_{= \{c_i, R_i\}} + \lambda \sum_{s' \in S} p(s'|s_i, a) u^*(s') \right\}$$

If $a = \alpha$:

$$\underbrace{r(s_i, \alpha)}_{= \omega_i} + \lambda \sum_{s' \in S} p(s'|s_i, \alpha) u^*(s') =$$

$$\omega_i + \lambda \left[(1-p_i) u^*(s_i') + \sum_{j=0}^N p_i q_j u^*(s_j) \right]$$

If $a = R$:

$$r(s_i, R) + \lambda \sum_{s' \in S'} p(s'|s_i, R) u^*(s') = c + \lambda \sum_{j=0}^N q_j u^*(s_j)$$

$\underbrace{r(s_i, R)}_{=c}$

Taken together, we have:

$$u^*(s_i) = \max \left\{ \begin{array}{l} \omega_i + \lambda [(1-p_i)u^*(s_i) + \sum_{j=0}^N p_j q_j u^*(s_j)] \\ c + \lambda \sum_{j=0}^N q_j u^*(s_j) \end{array} \right\} \quad (**)$$

Let's analyze this. Assume ^(wlog) that the salaries are sorted:

$$\omega_1 < \omega_2 < \dots < \omega_N.$$

It is also given that we should assume $p_i = p$.

Then (*) : constant (in i)

$$u^*(s_i) = \frac{\omega_i + \lambda p \sum_{j=0}^N q_j u^*(s_j)}{1 - \lambda(1-p)} \propto \omega_i$$

is an increasing function in i .

(That is,
 $u^*(s_1) < u^*(s_2) < \dots < u^*(s_N).$)

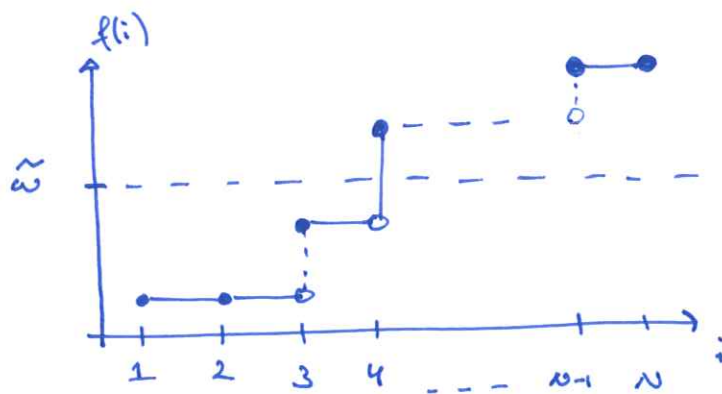
And from $(**)$:

$$u^*(s_i) = \max \left\{ \begin{array}{l} \alpha \\ c + \lambda \left[\underbrace{(1-p)u^*(s_i)}_{\text{increasing in } i} + \underbrace{\sum_{j=0}^N p q_j u^*(s_j)}_{\text{constant (in } i)}} \right] \end{array} \right\}$$

$\stackrel{\text{def}}{=} \tilde{\omega}, \text{ constant (in } i)$

$$\stackrel{\text{def.}}{=} \max \left\{ \begin{array}{l} f(i) \\ \tilde{\omega} \end{array} \right\},$$

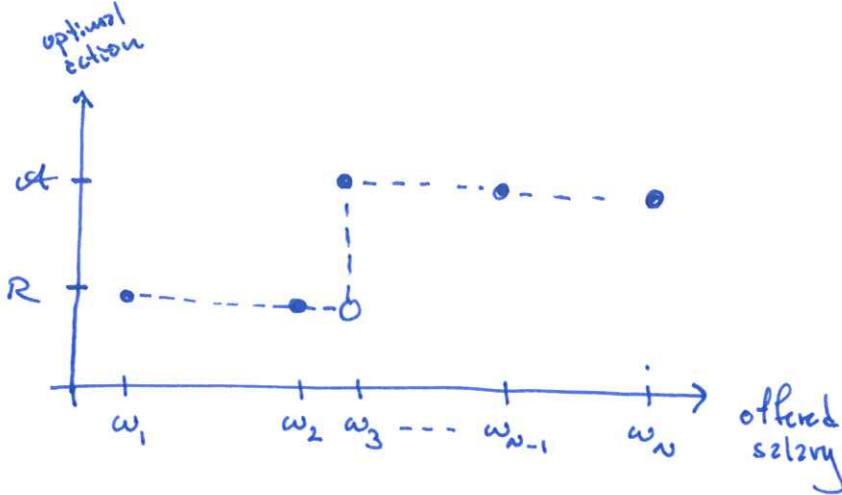
where $f(i)$ is an increasing function in i .



The optimal policy is to α (cept) if $f(i) > \tilde{\omega}$.
 Since $f(i)$ is increasing, this means that if
 $f(i) > \tilde{\omega} \implies f(i+1) > \tilde{\omega}$.

In other words, if there is a salary w_i such that $f(i) > \tilde{w}$, then it should be accepted. But then so should the salary w_{i+1} since $f(i+1) > \tilde{w}$, etc.

This implies that the solution is a threshold policy:



Ex 3.3 | An order is received w.p. p at each time step.

One can choose to process all orders (setup cost $K > 0$), or wait (costs $c > 0$ per order on hold). The maximum number of unfilled orders is n . There is a discount factor λ . Characterize an optimal processing policy.

Solution:

State-space: $S = \{0, 1, \dots, n\}$

of unfilled orders when starting the period

Actions:

$A_n = \{P\}$
process all orders

$A_i = \{W, P\}$ $i = 1, \dots, n-1$
wait *process all orders*

Rewards:

In this case, we deal with costs:

- $r(i, P) = -K$
- $r(i, W) = -ci$ $i = 1, \dots, n-1$
- $r(n, P) = -K$

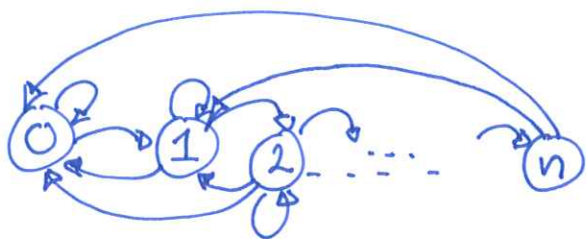
Time-horizon and objective:

∞ -horizon, discounted:

$$\mathbb{E} \left\{ \sum_{t=0}^{\infty} \lambda^t r(s_t, a_t) \right\}$$

Remark:
we can interpret $1-\lambda$ as the probability of "going out of business" (e.g., bankruptcy).

Transitions:



- $P(i|i, W) = 1-p$
wait, and no new order
- $P(i+1|i, W) = p$
wait, and a new order $i = 1, \dots, n-1$
- $P(0|i, P) = 1-p$
process everything, and no new order
- $P(1|i, P) = p$
process everything, and receive a new order
- $P(0|n, P) = 1-p$
- $P(1|n, P) = p$

Bellman equation:

$$u^*(s) = \max_{a \in A_s} \left\{ r(s, a) + \lambda \sum_{s' \in S} p(s'|s, a) u^*(s') \right\}$$

If $s=n$:

$$u^*(n) = \max_{a \in A_n} \left\{ \underbrace{r(n, a)}_{= -K} + \lambda \sum_{s' \in S} p(s'|n, a) u^*(s') \right\}$$

$$= -K + \lambda [(1-p)u^*(0) + pu^*(1)] \quad (*)$$

If $s = i \wedge n$:

$$u^*(i) = \max_{a \in A_i} \left\{ r(i, a) + \lambda \sum_{s' \in S} p(s' | i, a) u^*(s') \right\}$$

$\underbrace{\hspace{10em}}_{= \{P, W\}}$

Evaluate the actions separately:

If $a = P$:

$$\underbrace{r(i, P)}_{=-K} + \lambda \sum_{s' \in S} p(s' | i, P) u^*(s') =$$

$$-K + \lambda [(1-p)u^*(0) + pu^*(1)] \stackrel{\text{def.}}{=} \delta$$

If $a = W$:

$$\underbrace{r(i, W)}_{=-ci} + \lambda \sum_{s' \in S} p(s' | i, W) u^*(s') =$$

$$-ci + \lambda [(1-p)u^*(i) + pu^*(i+1)]$$

So that:

$$u^*(i) = \max_{\substack{P \\ W}} \left\{ \delta, -ci + \lambda [(1-p)u^*(i) + pu^*(i+1)] \right\}. \quad (**)$$

It is reasonable to assume that the optimal policy is of threshold type. Can we prove that it is?

A threshold policy is such that

$$a^*(i) = \bar{p} \Rightarrow a^*(i+1) = \bar{p}.$$

This is equivalent to

$$\begin{aligned} \delta > -c_i + \lambda [(1-p)u^*(i) + pu^*(i+1)] &\Rightarrow \\ \delta > -c_{(i+1)} + \lambda [(1-p)u^*(i+1) + pu^*(i+2)] \end{aligned}$$

Note that this holds if

$$u^*(i) \geq u^*(i+1), \quad (\square)$$

since $c > 0$.

Because then:

$$\begin{aligned} \delta > -c_i + \lambda [(1-p)u^*(i) + pu^*(i+1)] &\quad \leftarrow \text{since } c > 0 \\ > -c_{(i+1)} + \lambda [(1-p)u^*(i) + pu^*(i+1)] \\ &\quad \underbrace{\hspace{1.5cm}}_{\geq u^*(i+1)} \quad \underbrace{\hspace{1.5cm}}_{\geq u^*(i+2)} \\ \geq -c_{(i+1)} + \lambda [(1-p)u^*(i+1) + pu^*(i+2)] \end{aligned}$$

Can we prove that (\square) holds?

We will use value iteration to prove (\square) .

Base case:

Let $u_0(s) = 0$ for all $s \in S$.

Then $u_0(i) \geq u_0(i+1)$ holds trivially.

Induction:

Assume that $u_k(i) \geq u_k(i+1)$.

(We will show that then $u_{k+1}(i) \geq u_{k+1}(i+1)$ holds.)

If $s = i+1 \leq n$:

It is clear that the following holds:

$$\underbrace{-c(i+1)}_{\substack{\leq -c \\ \text{since} \\ c > 0}} + \lambda \left[\underbrace{(1-p)u_k(i+1)}_{\substack{\leq u_k(i) \\ \text{by assumption}}} + \underbrace{pu_k(i+2)}_{\substack{\leq u_k(i+1) \\ \text{by assumption}}} \right] \leq -c_i + \lambda \left[(1-p)u_k(i) + pu_k(i+1) \right] \quad (\Delta)$$

Let $F_k(x)$ be defined as

$$F_k(x) = \max \left\{ -K + \lambda \left[(1-p)u_k(0) + pu_k(1) \right], x \right\}.$$

This function is increasing in x .

One iteration of value iteration is (compare the computations that led to $(**)$):

$$\begin{aligned} u_{k+1}(i+1) &= \max \left\{ -K + \lambda \left[(1-p)u_k(0) + pu_k(1) \right], \right. \\ &\quad \left. -c(i+1) + \lambda \left[(1-p)u_k(i+1) + pu_k(i+2) \right] \right\} \\ &= F_k(-c(i+1) + \lambda \left[(1-p)u_k(i+1) + pu_k(i+2) \right]) \end{aligned}$$

$$\leq \bar{F}_k(-ci + \lambda[(1-p)u_k(i) + pu_k(i+1)])$$

by
(A)
and that
 $F_k(x) \geq x$

$$= \max \left\{ -K + \lambda[(1-p)u_k(0) + pu_k(1)], \right. \\ \left. -ci + \lambda[(1-p)u_k(i) + pu_k(i+1)] \right\}$$

$$= u_{k+1}(i).$$

by definition of one
iteration of value iteration,
compare (**).

If $s = n$:

Again, one iteration of value iteration (compare
computation of (*)) yields:

$$u_{k+1}(n) = -K + \lambda[(1-p)u_k(0) + pu_k(1)]$$

$$\leq \max \left\{ -K + \lambda[(1-p)u_k(0) + pu_k(1)], \right. \\ \left. -c(n-1) + \lambda[(1-p)u_k(n-1) + pu_k(n)] \right\}$$

obviously
smaller than
the max of itself
and something
else

$$= u_{k+1}(n-1).$$

compare
(**)

In summary, we have now shown that

$$u_k(i) \geq u_k(i+1) \Rightarrow u_{k+1}(i) \geq u_{k+1}(i+1)$$

for all $i = 0, \dots, n-1$.

Taking the limit:

We have shown that

- $u_0(i) \geq u_0(i+1)$
- $u_k(i) \geq u_k(i+1) \Rightarrow u_{k+1}(i) \geq u_{k+1}(i+1)$.

Since it is well known that

$$u^*(s) = \lim_{k \rightarrow \infty} u_k(s),$$

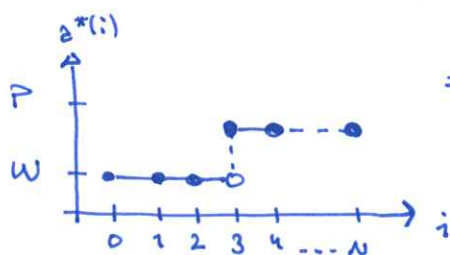
when u_k is computed using value iteration for discounted problems, we conclude that

$$u^*(i) \geq u^*(i+1),$$

which is \square .

See, e.g.,
p. 84 in Vol. 2
of Bertsekas'
"Dynamic Program.
and Optimal Control".

In summary, this proves that $a^*(i) = P \Rightarrow a^*(i+1) = P$



\Rightarrow To find the optimal policy, we just need to check which of the $N+1$ possible threshold policies is the best.

(This entails performing the policy evaluation step of policy iteration $N+1$ times.)