

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 reward-to-go at time t ← expected

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. 3.2 | Machine replacement problem.

Machine breaks down with prob. θ . Cost R to replace machine, and cost c of not being able to utilize it. Find optimal policy for T time-steps.

a, Model as MDP.

b, solve for $\theta = 0.1, R = 10, c = 5$ and $T = 3$.

Solution:

a, State-space: $S = \{ \text{Functional}, \text{Broken} \} = \{ F, B \}$.

Actions: $A = \{ \text{Continue}, \text{Replace} \} = \{ C, R \}$.

Time-horizon and objective:

Finite-horizon $T < \infty$:

$$E \left[\sum_{t=0}^{T-1} r_t(s_t, a_t) + r_T(s_T) \right]$$

Rewards:

Terminal:

$$r_T(\cdot) = 0$$

non-terminal:

$$r_t(s = \cdot, a = R) = -R$$

$$r_t(s = B, a = C) = -c$$

$$r_t(s = F, a = C) = 0$$

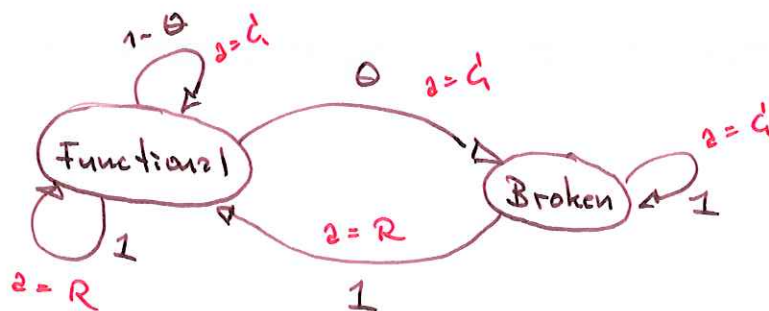
← chooses to replace

← not able to utilize the machine

← normal operation

Note: we allow for replacing a functional machine, and will let DP show that this would be nonsense. We could have incorporated prior knowledge by:
 $A(F) = \{ C \}$
 $A(B) = \{ C, R \}$

Transitions:



- $P_t(s' = B | s = B, a = C) = 1$ neglect to replace
- $P_t(s' = F | s = B, a = R) = 1$ replace
- $P_t(s' = F | s = F, a = C) = 1 - \theta$ does not break down
- $P_t(s' = B | s = F, a = C) = \theta$ breaks down
- $P_t(s' = F | s = F, a = R) = 1$ replace functional machine.

Note:
 (We assume the machine cannot break down immediately when it has been replaced.)

b, To solve the MDP, recall the backward induction:

- $u_T^*(s_T) = r_T(s_T)$
- $u_t^*(s_t) = \max_{a \in \mathcal{A}_{s_t}} \left\{ r_t(s_t, a) + \sum_{s' \in \mathcal{S}} P_t(s' | s_t, a) u_{t+1}^*(s') \right\}$
- $a_t^*(s_t) \in \arg \max_{a \in \mathcal{A}_{s_t}} \left\{ \text{---} \parallel \text{---} \right\}$

We are given that

$$\theta = 0.1, R = 10, C = 5 \text{ and } T = 3.$$

Backward - C2SC:

- $u_3^*(F) = r_3(F) = 0$
- $u_3^*(B) = r_3(B) = 0$

It is convenient to record the results in two tables:

u^* :

$t =$	0	1	2	3
$s = F$				0
$s = B$				0

a^* :

$t =$	0	1	2	3
$s = F$.
$s = B$.

At $t=T$, we are indifferent to what action is picked.

Recursion:

$$\begin{aligned}
 u_2^*(F) &= \max_{a \in \mathcal{A}} \left\{ r(F, a) + \sum_{s' \in \mathcal{S}} p_{t'}(s' | F, a) u_3^*(s') \right\} \\
 &= \max \left\{ \begin{aligned} &r(F, C) + p(F|F, C) u_3^*(F) + p(B|F, C) u_3^*(B), \\ &r(F, R) + p(F|F, R) u_3^*(F) + p(B|F, R) u_3^*(B) \end{aligned} \right\} \\
 &= \max \left\{ 0 + (1-\theta) \cdot 0 + \theta \cdot 0, -R + 1 \cdot 0 + 0 \cdot 0 \right\} \\
 &= \max \left\{ 0, -R \right\} = 0,
 \end{aligned}$$

indicates what action the term corresponds to

with $a_2^*(F) = C$.

$$\begin{aligned}
\bullet u_2^*(B) &= \max_{a \in \mathcal{A}} \left\{ r(B, a) + \sum_{s' \in \mathcal{S}} p_t(s'|B, a) u_3^*(s') \right\} \\
&= \max \left\{ r(B, C) + p(F|B, C) u_3^*(F) + p(B|B, C) u_3^*(B), \right. \\
&\quad \left. r(B, R) + p(F|B, R) u_3^*(F) + p(B|B, R) u_3^*(B) \right\} \\
&= \max \left\{ -C + 0 \cdot 0 + 1 \cdot 0, -R + 1 \cdot 0 + 0 \cdot 0 \right\} \\
&= \max \left\{ -C, -R \right\} = \max \left\{ -5, -10 \right\} = -5,
\end{aligned}$$

with $a_2^*(B) = C$.

Continue the recursion one more step:

$$\begin{aligned}
\bullet u_1^*(F) &= \max \left\{ r(F, C) + p(F|F, C) u_2^*(F) + p(B|F, C) u_2^*(B), \right. \\
&\quad \left. r(F, R) + p(F|F, R) u_2^*(F) + p(B|F, R) u_2^*(B) \right\} \\
&= \max \left\{ 0 + (1-\theta) \cdot 0 + \theta \cdot (-5), -R + 1 \cdot 0 + 0 \cdot (-5) \right\} \\
&= \max \left\{ -5\theta, -R \right\} = \max \left\{ -0.5, -10 \right\} = -0.5,
\end{aligned}$$

with $a_1^*(F) = C$.

$$\begin{aligned}
\bullet u_1^*(B) &= \max \left\{ r(B, C) + p(F|B, C) u_2^*(F) + p(B|B, C) u_2^*(B), \right. \\
&\quad \left. r(B, R) + p(F|B, R) u_2^*(F) + p(B|B, R) u_2^*(B) \right\} \\
&= \max \left\{ -C + 0 \cdot 0 + 1 \cdot (-5), -R + 1 \cdot 0 + 0 \cdot (-5) \right\} \\
&= \max \left\{ -C - 5, -R \right\} = \max \left\{ -10, -10 \right\} = -10,
\end{aligned}$$

with $a_1^*(B) \in \{C, R\}$.

Note that both actions yields the same value, so we are free to chose one.

Finally, the last stage in the recursion:

$$\begin{aligned}
 \bullet u_0^*(F) &= \max \left\{ \begin{array}{l} r(F, C) + p(F|F, C) u_1^*(F) + p(B|F, C) u_1^*(B), \\ r(F, R) + p(F|F, R) u_1^*(F) + p(B|F, R) u_1^*(B) \end{array} \right\} \\
 &= \max \left\{ \begin{array}{l} 0 + (1-\theta)(-0.5) + \theta \cdot (-10), \\ -R + 1 \cdot (-0.5) + 0 \cdot (-10) \end{array} \right\} \\
 &= \max \left\{ \begin{array}{l} -0.9 \times 0.5 - 1, \\ -10 - 0.5 \end{array} \right\} \\
 &= \max \left\{ \begin{array}{l} -1.45, \\ -10.5 \end{array} \right\} = -1.45,
 \end{aligned}$$

with $a_0^*(F) = C$.

$$\begin{aligned}
 \bullet u_0^*(B) &= \max \left\{ \begin{array}{l} r(B, C) + p(F|B, C) u_1^*(F) + p(B|B, C) u_1^*(B), \\ r(B, R) + p(F|B, R) u_1^*(F) + p(B|B, R) u_1^*(B) \end{array} \right\} \\
 &= \max \left\{ \begin{array}{l} -C + 0 \cdot (-0.5) + 1 \cdot (-10), \\ -R + 1 \cdot (-0.5) + 0 \cdot (-10) \end{array} \right\} \\
 &= \max \left\{ \begin{array}{l} -C - 10, \\ -R - 0.5 \end{array} \right\} = \max \left\{ \begin{array}{l} -15, \\ -10.5 \end{array} \right\} = -10.5,
 \end{aligned}$$

with $a_0^*(B) = R$.

In summary, we have that

$\epsilon =$	0	1	2	3
$u^*:$ $s = F$	-1.45	-0.5	0	0
$s = B$	-10.5	-10	-5	0

$\epsilon =$	0	1	2	3
$d^*:$ $s = F$	G	G	G	•
$s = B$	R	G/R	G	•

either action
is optimal

Note: The computations above are essentially matrix-vector multiplications and additions. As such, they can be computed much more succinctly than above. Try to figure out how on your own! This is also good to keep in mind when implementing on a computer.

Ex 3.3 | 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 $p > 0$ 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 b_{max} SEK. Then:

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

↖ current bid under consideration

Actions:

$A_b = \{R, A\}$: 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:

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

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

Non-terminal:

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

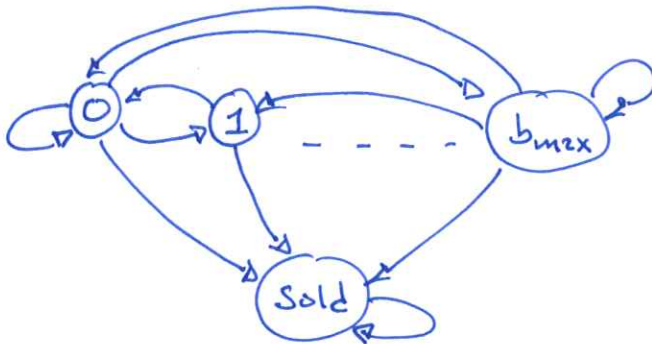
the current offer

the future interest that we will collect

- $r_t(s=b, a=R) = 0$

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

Transitions:



Note: we could take A_{sold} to be $\{0, R\}$ as we did last time. In that case, the choice of action in the policy will be arbitrary.

- $P_t(s'=b' | s=b, a=R) = P(\omega_t = b')$

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

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

we decide to sell

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

if we have sold

To solve the MDP, recall the backward induction:

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

Base case:

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

Recursion:

For $s_t = \text{sold}$:

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

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 (*)

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

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 $a = \alpha$:

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

If $a = R$:

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

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

Taken together in (**), we obtain:

$$u_t^*(b) = \max \left\{ \underset{\mathcal{A}}{b(1+p)^{N-t}}, \underset{\mathcal{R}}{\sum_{i=0}^{b_{max}} P_r \{ \omega_t = i \} u_{t+1}^*(i)} \right\}$$

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

def. = α_t

def.

$$= \max \left\{ \underset{\mathcal{A}}{b}, \underset{\mathcal{R}}{\alpha_t} \right\} \times (1+p)^{N-t}$$

Hence, the optimal policy is:

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

- \mathcal{A} (cept) if $b > \alpha_t$,
- \mathcal{R} (eject) if $b \leq \alpha_t$,

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

←

where

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

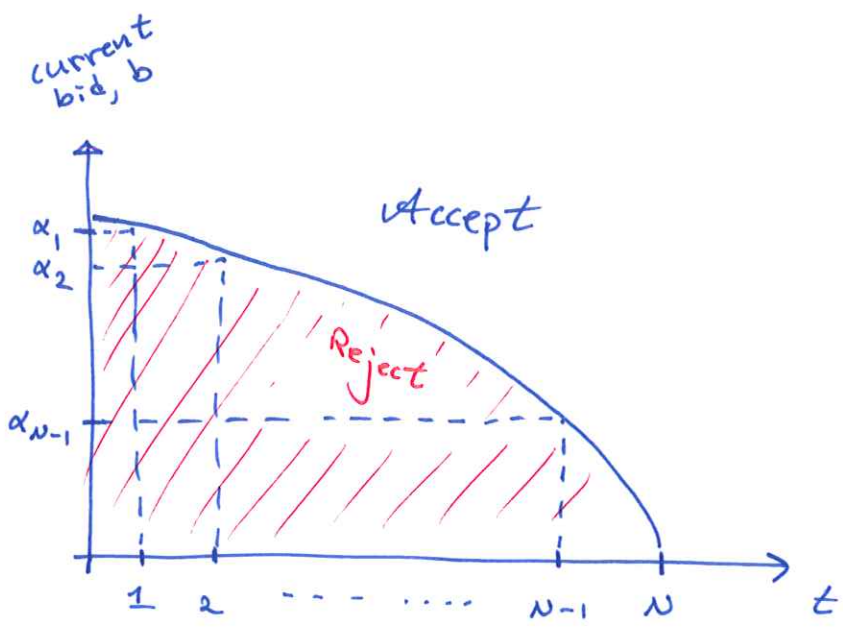
↗

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 3.6 | 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 : $\Pr\{T=k\} = \lambda^{k-1}(1-\lambda)$.

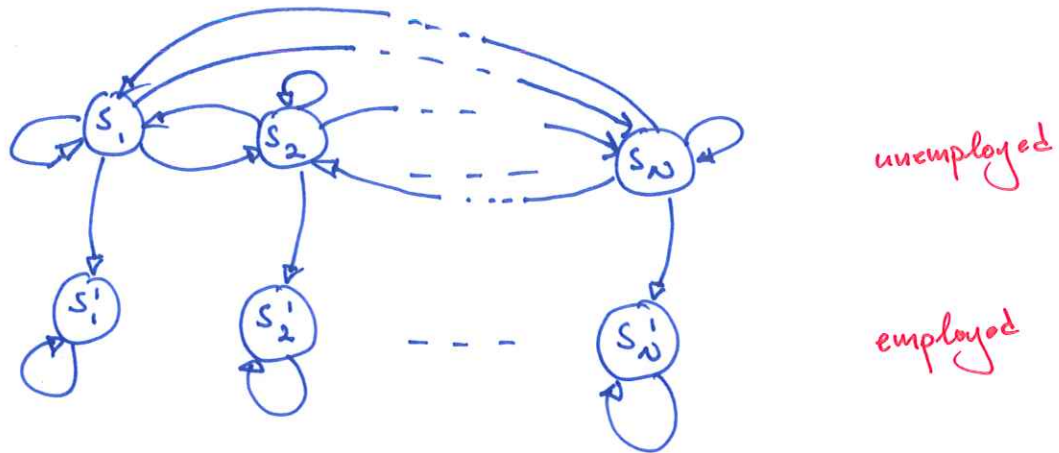
$1-\lambda$ can be interpreted as the worker's risk of dying each time unit. The worker's expected life-time is $\mathbb{E}\{T\} = \frac{1}{1-\lambda}$ time units. The worker tries to maximize the total money earned, subject to knowing that she might die.

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

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

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

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

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

$$= \max \left\{ \underbrace{w_i}_{\mathcal{A}}, \underbrace{(1-\lambda) \left[c + \lambda \sum_{j=0}^N q_j u^*(s_j) \right]}_{\substack{\text{def.} \\ = \bar{w}}} \right\} \cdot \frac{1}{1-\lambda}$$

$$= \max \left\{ \underbrace{w_i}_{\mathcal{A}}, \underbrace{\bar{w}}_{\mathcal{R}} \right\} \cdot \frac{1}{1-\lambda}$$

15
21

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

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

Remark:
How does this help us $\rightarrow u^*(s)$ 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:

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

Rewards:

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

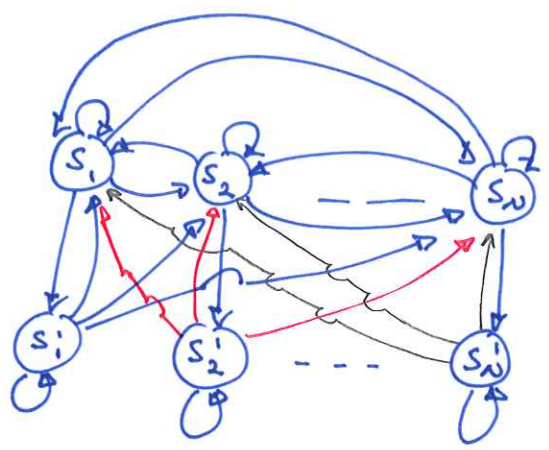
Time-horizon and objective:

∞ -horizon, discounted:

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

Transitions:

- $p(s'_i | s'_i, C) = 1 - p_i$
not fired
- $p(s_j | s'_i, C) = p_i \cdot q_j$
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 $u^*(s_i')$:

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

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

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 $a = R$:

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

Taken together, we have:

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

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 (*) :

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

constant (in i)

is an increasing function in i .

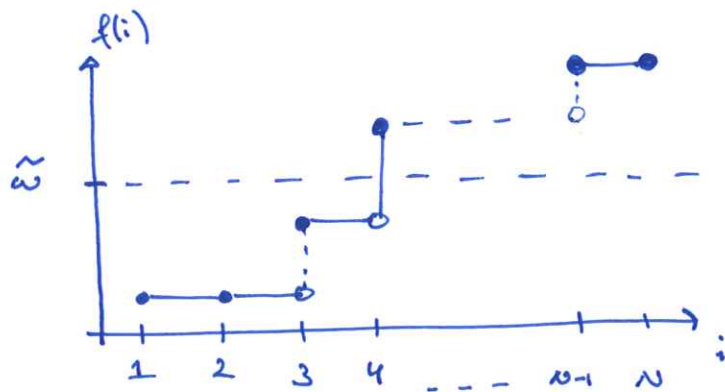
(That is,
 $u^*(s_1) < u^*(s_2) < \dots < u^*(s_N).$)

And from $(**)$:

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

$$\stackrel{\text{def.}}{=} \max \left\{ \underset{\text{at}}{f(i)}, \underset{R}{\tilde{\omega}} \right\},$$

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



The optimal policy is to accept 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:

