

An Introduction to Computational Macroeconomics

Dynamic Programming: Chapter 6

John Stachurski

June – July 2022

Markov Decision Processes

Next we cover Markov decision processes (MDPs)

- A class of dynamic programs
- Broad enough to encompass many economic applications
- Includes optimal stopping problems as a special case
- Clean, powerful theory
- A range of important algorithms

Also a cornerstone for

- reinforcement learning, artificial intelligence, etc.

MDPs are dynamic programs characterized by two features

1. Rewards are additively separable:

$$\text{lifetime reward} = \mathbb{E} \sum_{t \geq 0} \beta^t R_t$$

2. The discount rate is constant

For now we restrict attention to finite state and action spaces

- Routinely used in quantitative applications
- Avoids technical issues we can put aside for later

Notation

Let X and A be any sets

- $\wp(A) :=$ the set of all subsets of A

A **correspondence** Γ from X to A is a map from X to $\wp(A)$

- called **nonempty** if $\Gamma(x) \neq \emptyset$ for all $x \in X$

Examples.

- $\Gamma(x) = [0, x]$ is a correspondence from \mathbb{R} to \mathbb{R}
- $\Gamma(x) = [-x, x]$ is a nonempty correspondence from \mathbb{R} to \mathbb{R}

We study a controller who, at each integer $t \geq 0$

1. observes the current state X_t
2. responds with an action A_t

Her aim is to maximize expected discounted rewards

$$\mathbb{E} \sum_{t \geq 0} \beta^t r(X_t, A_t), \quad X_0 = x_0 \text{ given}$$

We take as given

1. a finite set X called the **state space** and
2. a finite set A called the **action space**

The actions of the controller are limited by a **feasible correspondence** Γ

- A correspondence from X to A
- $\Gamma(x)$ is the set of actions available to the controller in state x

Given Γ , we set

$$G := \{(x, a) \in X \times A : a \in \Gamma(x)\}$$

- called the set of **feasible state-action pairs**

Reward $r(x, a)$ is received at feasible state-action pair (x, a) ,

A **stochastic kernel** from G to X is a map $P: G \times X \rightarrow \mathbb{R}_+$ satisfying

$$\sum_{x' \in X} P(x, a, x') = 1 \quad \text{for all } (x, a) \text{ in } G$$

Interpretation

- For each feasible state-action pair, $P(x, a, \cdot)$ is a distribution
- The next period state x' is selected from $P(x, a, \cdot)$

Now let's put it all together:

Given X and A , a **Markov decision process (MDP)** is a tuple (Γ, β, r, P) where

1. Γ is a nonempty correspondence from $X \rightarrow A$
2. β is a constant in $(0, 1)$
3. r is a function from G to \mathbb{R}
4. P is a stochastic kernel from G to X

In the foregoing,

- β is called the **discount factor**
- r is called the **reward function**

Algorithm 1: MDP dynamics: states, actions, and rewards

```
 $t \leftarrow 0$   
input  $X_0$   
while  $t < \infty$  do  
    observe  $X_t$   
    choose action  $A_t$  from  $\Gamma(X_t)$   
    receive reward  $r(X_t, A_t)$   
    draw  $X_{t+1}$  from  $P(X_t, A_t, \cdot)$   
     $t \leftarrow t + 1$   
end
```

Rules:

- Choose $(A_t)_{t \geq 0}$ to maximize $\mathbb{E} \sum_{t \geq 0} \beta^t r(X_t, A_t)$
- Actions don't depend on future outcomes

The **Bellman equation** is

$$v(x) = \max_{a \in \Gamma(x)} \left\{ r(x, a) + \beta \sum_{x' \in \mathbf{X}} v(x') P(x, a, x') \right\}$$

Reduces an infinite horizon problem to a two period problem

In the two period problem, the controller trades off

1. current rewards and
2. expected discounted value from future states

Current actions influence both terms

Example. Cake eating (no labor income), where

$$W_{t+1} = RW_t - C_t \quad (t = 0, 1, \dots)$$

- Investing d dollars today returns Rd next period
- $C_t, W_t \geq 0$ are current consumption and wealth
- Assume wealth takes values in a finite set $W \subset \mathbb{R}_+$

The agent seeks to maximize

$$\mathbb{E} \sum_{t \geq 0} \beta^t u(C_t) \quad \text{given } W_0 = w$$

Bellman equation:

$$v(w) = \max_{0 \leq w' \leq Rw} \{u(Rw - w') + \beta v(w')\}$$

This model can be framed as an MDP with W as the state space

- The action is $S_t = RW_t - C_t = \text{next-period wealth}$
- The action space is also W
- The feasible correspondence is

$$\Gamma(w) = \{s \in W : s \leqslant Rw\}$$

- Current reward is

$$r(w, s) = u(Rw - s) \quad (s \in \Gamma(w))$$

- The stochastic kernel is $P(w, s, w') = \mathbb{1}\{w' = s\}$

Example. All optimal stopping problems can be framed as MDPs

- See the text for details

This is important from a theoretical perspective

- illustrates the generality of MDPs

However, expressing optimal stopping problems as an MDP requires an extra state variable

- status $S_t = 1$ if stopped else zero

Hence treating optimal stopping problems separately is neater

Policies

Actions will be governed by policies

- maps from states to actions
- today's action is a function of today's state!

The set of **feasible policies** is

$$\Sigma := \text{all } \sigma \in A^X \text{ s.t. } \sigma(x) \in \Gamma(x) \text{ for all } x \in X$$

Meaning of selecting σ from Σ :

- respond to state X_t with action $A_t := \sigma(X_t)$ at all t

Dynamics

What happens when we always follow $\sigma \in \Sigma$?

Now

$$X_{t+1} \sim P(X_t, \sigma(X_t), \cdot) \quad \text{at every } t$$

Thus, X_t updates according to the stochastic matrix

$$P_\sigma(x, x') := P(x, \sigma(x), x') \quad (x, x' \in \mathbf{X})$$

The state process becomes P_σ -Markov

- Fixing a policy “closes the loop” in the state dynamics
- Solving an MDP means choosing a Markov chain!

Rewards

Under the policy σ , rewards at x are $r(x, \sigma(x))$

Let

$$r_\sigma(x) := r(x, \sigma(x)) \quad (x \in \mathbf{X})$$

Now set

$$\mathbb{E}_x := \mathbb{E}[\cdot \mid X_0 = x]$$

Then the expected time t reward is

$$\mathbb{E}_x r(X_t, A_t) = \mathbb{E}_x r_\sigma(X_t) = (P_\sigma^t r_\sigma)(x)$$

Let $(X_t)_{t \geq 0}$ be P_σ -Markov with $X_0 = x$

The lifetime value of σ starting from x is

$$v_\sigma(x) := \mathbb{E}_x \sum_{t \geq 0} \beta^t r_\sigma(X_t)$$

Since $\beta < 1$, we have $r(\beta P_\sigma) < 1$ and hence

$$v_\sigma = \sum_{t \geq 0} \beta^t P_\sigma^t r_\sigma = (I - \beta P_\sigma)^{-1} r_\sigma$$

The **value function** is defined as

$$v^*(x) = \max_{\sigma \in \Sigma} v_\sigma(x) \quad (x \in X)$$

Recall that the Bellman equation is

$$v(x) = \max_{a \in \Gamma(x)} \left\{ r(x, a) + \beta \sum_{x' \in X} v(x') P(x, a, x') \right\}$$

The **Bellman operator** for the MDP is the self-map T on \mathbb{R}^X defined by

$$(Tv)(x) = \max_{a \in \Gamma(x)} \left\{ r(x, a) + \beta \sum_{x' \in X} v(x') P(x, a, x') \right\}$$

Obviously

- $Tv = v$ iff v satisfies the Bellman equation
- T is order-preserving on \mathbb{R}^X

Fix $v \in \mathbb{R}^X$

A policy $\sigma \in \Sigma$ is called **v -greedy** if

$$\forall x \in X, \sigma(x) \in \operatorname{argmax}_{a \in \Gamma(x)} \left\{ r(x, a) + \beta \sum_{x' \in X} v(x') P(x, a, x') \right\}$$

A policy $\sigma \in \Sigma$ is called **optimal** if

$$v_\sigma = v^*$$

Thus,

σ is optimal \iff lifetime value is maximal at each state

Proposition. For the MDP described above

1. v^* is the unique fixed point of T in \mathbb{R}^X
2. T is a contraction of modulus β on \mathbb{R}^X under the norm $\|\cdot\|_\infty$
3. A feasible policy is optimal if and only if it is v^* -greedy
4. At least one optimal policy exists

Proof:

- similar to that for optimal stopping
- full details deferred until we study RDPs

Ex. Show that optimal iff v^* -greedy implies that at least one optimal policy exists

Proof: Given $v \in \mathbb{R}^X$, a v -greedy policy is a $\sigma \in \Sigma$ such that

$$\sigma(x) \in \operatorname{argmax}_{a \in \Gamma(x)} \left\{ r(x, a) + \beta \sum_{x' \in X} v(x') P(x, a, x') \right\} \quad \text{for all } x \in X$$

For each $v \in \mathbb{R}^X$, a v -greedy policy exists

- just select a maximizer a_x^* from the nonempty set $\Gamma(x)$ at each x in X

Hence a v^* -greedy policy exists

Hence an optimal policy exists

Ex. Show that optimal iff v^* -greedy implies that at least one optimal policy exists

Proof: Given $v \in \mathbb{R}^X$, a v -greedy policy is a $\sigma \in \Sigma$ such that

$$\sigma(x) \in \operatorname{argmax}_{a \in \Gamma(x)} \left\{ r(x, a) + \beta \sum_{x' \in X} v(x') P(x, a, x') \right\} \quad \text{for all } x \in X$$

For each $v \in \mathbb{R}^X$, a v -greedy policy exists

- just select a maximizer a_x^* from the nonempty set $\Gamma(x)$ at each x in X

Hence a v^* -greedy policy exists

Hence an optimal policy exists

Computing the Value of a Policy

How should we compute the value v_σ of a given policy σ ?

We saw above that

$$v_\sigma = (I - \beta P_\sigma)^{-1} r_\sigma$$

- Computationally helpful when X is small
- But problematic for large dynamic programs

If, say, X has 10^6 elements, then $I - \beta P_\sigma$ is $10^6 \times 10^6$

Matrices of this size are difficult invert—or even store in memory

Another way to compute v_σ : use the **policy operator**

$$(T_\sigma v)(x) = r(x, \sigma(x)) + \beta \sum_{x' \in X} v(x') P(x, \sigma(x), x')$$

- Defined at all $v \in \mathbb{R}^X$
- Analogous to T_σ for the optimal stopping problem

In vector notation, we can write

$$T_\sigma v = r_\sigma + \beta P_\sigma v$$

- T_σ is order-preserving on \mathbb{R}^X — why?

Ex. Show that T_σ is a contraction of modulus β on \mathbb{R}^X

For any v, w in \mathbb{R}^X we have

$$\begin{aligned} |T_\sigma v - T_\sigma w| &= \beta |P_\sigma v - P_\sigma w| \\ &= \beta |P_\sigma(v - w)| \\ &\leq \beta P_\sigma |v - w| \\ &\leq \beta P_\sigma \|v - w\|_\infty \mathbb{1} \\ &= \beta \|v - w\|_\infty \mathbb{1} \end{aligned}$$

Now use $|a| \leq |b|$ implies $\|a\|_\infty \leq \|b\|_\infty$

Ex. Show that T_σ is a contraction of modulus β on \mathbb{R}^X

For any v, w in \mathbb{R}^X we have

$$\begin{aligned} |T_\sigma v - T_\sigma w| &= \beta |P_\sigma v - P_\sigma w| \\ &= \beta |P_\sigma(v - w)| \\ &\leq \beta P_\sigma |v - w| \\ &\leq \beta P_\sigma \|v - w\|_\infty \mathbb{1} \\ &= \beta \|v - w\|_\infty \mathbb{1} \end{aligned}$$

Now use $|a| \leq |b|$ implies $\|a\|_\infty \leq \|b\|_\infty$

Ex. Show that v_σ is the unique fixed point of T_σ in \mathbb{R}^X

Proof: If $v = T_\sigma v$, then

$$v = r_\sigma + \beta P_\sigma v$$

Since $\beta < 1$, we then have

$$v = (I - \beta P_\sigma)^{-1} r_\sigma$$

Hence $v = v_\sigma$

Ex. Show that v_σ is the unique fixed point of T_σ in \mathbb{R}^X

Proof: If $v = T_\sigma v$, then

$$v = r_\sigma + \beta P_\sigma v$$

Since $\beta < 1$, we then have

$$v = (I - \beta P_\sigma)^{-1} r_\sigma$$

Hence $v = v_\sigma$

Algorithms

Previously we used value function iteration (VFI) to solve optimal stopping problems

Here we

1. present a generalization suitable for arbitrary MDPs
2. introduce two other important methods

The two other methods are called

1. Howard policy iteration and
2. Optimistic policy iteration

Algorithm 2: VFI for MDPs

input $v_0 \in \mathbb{R}^X$, an initial guess of v^*

input τ , a tolerance level for error

$\varepsilon \leftarrow \tau + 1$

$k \leftarrow 0$

while $\varepsilon > \tau$ **do**

for $x \in X$ **do**

$v_{k+1}(x) \leftarrow (Tv_k)(x)$

end

$\varepsilon \leftarrow \|v_k - v_{k+1}\|_\infty$

$k \leftarrow k + 1$

end

Compute a v_k -greedy policy σ

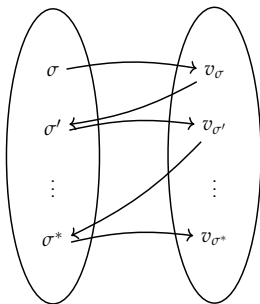
return σ

VFI is

- robust
- easy to implement and
- works relatively well in high dimensions after certain modifications

However, we can often find faster methods with a bit of effort

Howard Policy Iteration



Iterates between computing the value of a given policy and computing the greedy policy associated with that value

Algorithm 3: Howard policy iteration for MDPs

input $\sigma_0 \in \Sigma$, an initial guess of σ^*

$k \leftarrow 0$

$\varepsilon \leftarrow 1$

while $\varepsilon > 0$ **do**

$v_k \leftarrow$ the σ_k -value function $(I - \beta P_{\sigma_k})^{-1} r_{\sigma_k}$

$\sigma_{k+1} \leftarrow$ a v_k greedy policy

$\varepsilon \leftarrow \|\sigma_k - \sigma_{k+1}\|_\infty$

$k \leftarrow k + 1$

end

return σ_k

Advantages:

1. in a finite state setting, the algorithm always converges to the exact optimal policy in a finite number of steps
2. the rate of convergence is faster than VFI

We prove these facts in a more general setting when we discuss RDPs

Optimistic Policy Iteration

OPI borrows from both value function iteration and Howard policy iteration

The same as Howard policy iteration (HPI) except that

- HPI takes σ and obtains v_σ
- OPI takes σ and iterates m times with T_σ

Recall that $T_\sigma^m \rightarrow v_\sigma$ as $m \rightarrow \infty$

Hence OPI replaces v_σ with an approximation

Algorithm 4: Optimistic policy iteration for MDPs

input $v_0 \in \mathbb{R}^X$, an initial guess of v^*

input τ , a tolerance level for error

input $m \in \mathbb{N}$, a step size

$k \leftarrow 0$

$\varepsilon \leftarrow \tau + 1$

while $\varepsilon > \tau$ **do**

$\sigma_k \leftarrow$ a v_k -greedy policy

$v_{k+1} \leftarrow T_{\sigma_k}^m v_k$

$\varepsilon \leftarrow \|v_k - v_{k+1}\|_\infty$

$k \leftarrow k + 1$

end

return σ_k

Regarding m ,

- If $m = \infty$, OPI is identical to HPI
- If $m = 1$, OPI is identical to VFI

Usually, an intermediate value of m is better than both

- We investigate this in the applications below

The sequence $(\sigma_k)_{k \geq 1}$ always converges to an optimal policy

Application: Optimal Inventories

Previously we analyzed S-s inventory dynamics

- our aim was to understand Markov chains

But are such dynamics realistic?

We now investigate whether S-s behavior arises naturally in optimizing model

- firm chooses its inventory path to maximize firm value

We assume for now that the firm only sells one product

Given a demand process $(D_t)_{t \geq 0}$, inventory $(X_t)_{t \geq 0}$ obeys

$$X_{t+1} = m(X_t - D_{t+1}) + A_t$$

where

- $m(y) := y \vee 0$
- A_t is units of stock ordered this period
- The firm can store at most K items at one time

The state space is $X := \{0, \dots, K\}$

We assume $(D_t) \stackrel{\text{iid}}{\sim} \varphi \in \mathcal{D}(\mathbb{Z}_+)$

Profits are given by

$$\pi_t := X_t \wedge D_{t+1} - cA_t - \kappa \mathbb{1}\{A_t > 0\}$$

- Orders in excess of inventory are lost (rather than backfilled)
- c is unit product cost
- κ is a fixed cost of ordering inventory

With $\beta := 1/(1+r)$ and $r > 0$, the value of the firm is

$$V_0 = \mathbb{E} \sum_{t \geq 0} \beta^t \pi_t$$

Managers of the firm try to maximize shareholder value

Let

$$r(x, a) := \sum_{d \geq 0} (x \wedge d) \varphi(d) - ca - \kappa \mathbb{1}\{a > 0\}$$

The Bellman equation for this dynamic program is

$$v(x) = \max_{a \in \Gamma(x)} \left\{ r(x, a) + \beta \sum_{d \geq 0} v(m(x - d) + a) \varphi(d) \right\}$$

Since K is maximum inventory capacity, the set of feasible order sizes at x is

$$\Gamma(x) := \{0, \dots, K - x\}$$

Ex. Write down the Bellman operator for this model and prove it is a contraction on \mathbb{R}^X when paired with the supremum norm

An MDP with state space X and action space $A := X$

- Γ , r and β are as given above
- The stochastic kernel is

$$P(x, a, x') := \mathbb{P}\{m(x - D) + a = x'\} \quad \text{when } D \sim \varphi$$

Since the inventory model is an MDP, all optimality results apply

- the unique fixed point of the Bellman operator is v^*
- a policy σ^* is optimal if and only if it is v^* -greedy
- etc.

using Distributions, OffsetArrays

$m(x) = \max(x, 0)$ *# Convenience function*

```
function create_inventory_model(;  $\beta=0.98$ ,      # discount factor  
                                 $K=40$ ,          # maximum inventory  
                                 $c=0.2$ ,  $\kappa=2$ , # cost parameters  
                                 $p=0.6$ )        # demand parameter  
     $\phi(d) = (1 - p)^d * p$  # demand pdf  
    return (;  $\beta$ ,  $K$ ,  $c$ ,  $\kappa$ ,  $p$ ,  $\phi$ )  
end
```

"The function $B(x, a, v) = r(x, a) + \beta \sum_{x'} v(x') P(x, a, x')$."

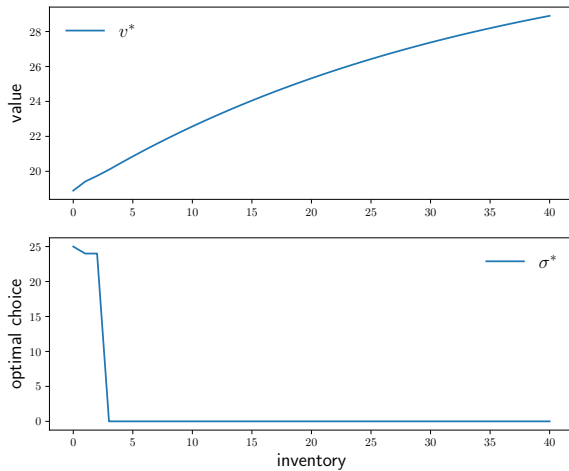
```
function B(x, a, v, model;  $d_{\max}=100$ )  
    (;  $\beta$ ,  $K$ ,  $c$ ,  $\kappa$ ,  $p$ ,  $\phi$ ) = model  
    reward = sum(min(x, d)* $\phi(d)$  for  $d$  in  $0:d_{\max}$ ) -  $c * a - \kappa * (a > 0)$   
    continuation_value =  $\beta * \text{sum}(v[m(x - d) + a] * \phi(d)$  for  $d$  in  $0:d_{\max}$ )  
    return reward + continuation_value  
end
```

"The Bellman operator."

```
function T(v, model)
    (;  $\beta$ , K, c,  $\kappa$ , p,  $\phi$ ) = model
    new_v = similar(v)
    for x in 0:K
         $\Gamma$ x = 0:(K - x)
        new_v[x], _ = findmax(B(x, a, v, model) for a in  $\Gamma$ x)
    end
    return new_v
end
```

"Get a v-greedy policy. Returns a zero-based array."

```
function get_greedy(v, model)
    (;  $\beta$ , K, c,  $\kappa$ , p,  $\phi$ ) = model
     $\sigma$ _star = OffsetArray(zeros{Int32, K+1}, 0:K)
    for x in 0:K
         $\Gamma$ x = 0:(K - x)
        _, a_idx = findmax(B(x, a, v, model) for a in  $\Gamma$ x)
         $\sigma$ _star[x] =  $\Gamma$ x[a_idx]
    end
    return  $\sigma$ _star
end
```



Ex. Try to replicate these plots

- Use the code given above (or at least the same parameters)
- Use value function iteration

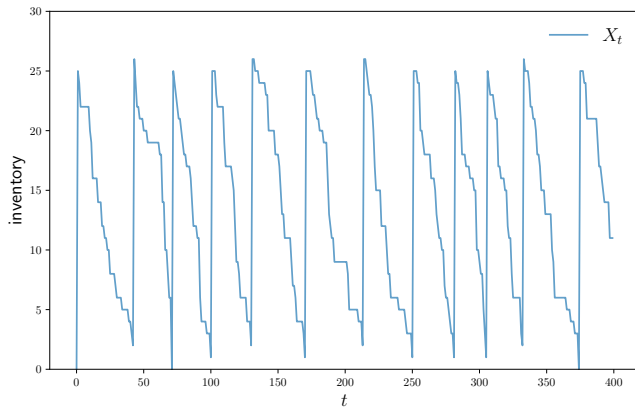


Figure: Optimal inventory dynamics

Optimal Savings with Labor Income

Wealth evolves according to

$$W_{t+1} = RW_t + Y_t - C_t \quad (t = 0, 1, \dots)$$

- (W_t) takes values in finite set $W \subset \mathbb{R}_+$
- (Y_t) is Q -Markov chain on finite set Y
- $C_t, W_t \geq 0$

The household maximizes

$$\mathbb{E} \sum_{t \geq 0} \beta^t u(C_t)$$

The model is an MDP:

We set the state to $x := (w, y) \in X := W \times Y$

The feasible correspondence is

$$\Gamma(w, y) = \{s \in W : s \leq R w + y\}$$

The current reward is $r(w, s) = u(R w + y - s)$

The stochastic kernel is

$$P((w, y), s, (w', y')) = \mathbb{1}\{w' = s\} Q(y, y')$$

Hence all of the MDP optimality results apply

Bellman operator:

$$(Tv)(w, y) =$$

$$\max_{w' \in \Gamma(w, y)} \left\{ u(Rw + y - w') + \beta \sum_{y' \in Y} v(w', y') Q(y, y') \right\}$$

The policy operator for given $\sigma \in \Sigma$ is

$$(T_\sigma v)(w, y) =$$

$$u(Rw + y - \sigma(w, y)) + \beta \sum_{y' \in Y} v(\sigma(w, y), y') Q(y, y')$$

How to solve $v_\sigma = (I - \beta P_\sigma)^{-1} r_\sigma$?

We set

$$P_\sigma((w, y), (w', y')) := \mathbb{1}\{\sigma(w, y) = w'\} Q(y, y')$$

and

$$r_\sigma(w, y) := u(Rw + y - \sigma(w, y))$$

How to use matrix algebra routines?

We map indices i, j for state (w_i, y_j) into a single index m , as in

$$x_m = (w_i, y_j)$$

```
using QuantEcon, LinearAlgebra, IterTools
```

```
function create_savings_model(; R=1.01,  $\beta$ =0.99,  $\gamma$ =2.5,  
                             w_min=0.01, w_max=5.0, w_size=200,  
                              $\rho$ =0.9,  $\nu$ =0.1, y_size=5)  
    w_grid = LinRange(w_min, w_max, w_size)  
    mc = tauchen(y_size,  $\rho$ ,  $\nu$ )  
    y_grid, Q = exp.(mc.state_values), mc.p  
    return (;  $\beta$ , R,  $\gamma$ , w_grid, y_grid, Q)  
end
```

```
"B(w, y, w') = u(R*w + y - w') +  $\beta \sum_{y'} v(w', y') Q(y, y')$ ."
```

```
function B(i, j, k,  $\nu$ , model)  
    (;  $\beta$ , R,  $\gamma$ , w_grid, y_grid, Q) = model  
    w, y, w' = w_grid[i], y_grid[j], w_grid[k]  
    u(c) = c^(1- $\gamma$ ) / (1- $\gamma$ )  
    c = R*w + y - w'  
    @views value = c > 0 ? u(c) +  $\beta$  * dot( $\nu$ [k, :], Q[j, :]) : -Inf  
    return value  
end
```

"The Bellman operator."

```
function T(v, model)
    w_idx, y_idx = (eachindex(g) for g in (model.w_grid, model.y_grid))
    v_new = similar(v)
    for (i, j) in product(w_idx, y_idx)
        v_new[i, j] = maximum(B(i, j, k, v, model) for k in w_idx)
    end
    return v_new
end
```

"The policy operator."

```
function T_σ(v, σ, model)
    w_idx, y_idx = (eachindex(g) for g in (model.w_grid, model.y_grid))
    v_new = similar(v)
    for (i, j) in product(w_idx, y_idx)
        v_new[i, j] = B(i, j, σ[i, j], v, model)
    end
    return v_new
end
```

```
include("finite_opt_saving_0.jl")

"Compute a v-greedy policy."
function get_greedy(v, model)
    w_idx, y_idx = (eachindex(g) for g in (model.w_grid, model.y_grid))
     $\sigma$  = Matrix{Int32}(undef, length(w_idx), length(y_idx))
    for (i, j) in product(w_idx, y_idx)
        _,  $\sigma$ [i, j] = findmax(B(i, j, k, v, model) for k in w_idx)
    end
    return  $\sigma$ 
end
```

```

"Get the value  $v_\sigma$  of policy  $\sigma$ ."
function get_value( $\sigma$ , model)
    # Unpack and set up
    (;  $\beta$ , R,  $\gamma$ , w_grid, y_grid, Q) = model
    wn, yn = length(w_grid), length(y_grid)
    n = wn * yn
    u(c) =  $c^{(1-\gamma)} / (1-\gamma)$ 
    # Function to extract (i, j) from  $m = i + (j-1)*wn$ 
    single_to_multi(m) = (m-1)%wn + 1, div(m-1, wn) + 1
    # Allocate and create single index versions of  $P_\sigma$  and  $r_\sigma$ 
    P_ $\sigma$  = zeros(n, n)
    r_ $\sigma$  = zeros(n)
    for m in 1:n
        i, j = single_to_multi(m)
        r_ $\sigma$ [m] = u(R * w_grid[i] + y_grid[j] - w_grid[ $\sigma$ [i, j]])
        for m' in 1:n
            i', j' = single_to_multi(m')
            if i' ==  $\sigma$ [i, j]
                P_ $\sigma$ [m, m'] = Q[j, j']
            end
        end
    end
    # Solve for the value of  $\sigma$ 
    v_ $\sigma$  = (I -  $\beta$  * P_ $\sigma$ ) \ r_ $\sigma$ 
    # Return as multi-index array
    return reshape(v_ $\sigma$ , wn, yn)
end

```

"Value function iteration routine."

```
function value_iteration(model, tol=1e-5)
    vz = zeros(length(model.w_grid), length(model.y_grid))
    v_star = successive_approx(v -> T(v, model), vz, tolerance=tol)
    return get_greedy(v_star, model)
end
```

"Howard policy iteration routine."

```
function policy_iteration(model)
    wn, yn = length(model.w_grid), length(model.y_grid)
     $\sigma$  = ones(Int32, wn, yn)
    i, error = 0, 1.0
    while error > 0
        v_ $\sigma$  = get_value( $\sigma$ , model)
         $\sigma_{\text{new}}$  = get_greedy(v_ $\sigma$ , model)
        error = maximum(abs.( $\sigma_{\text{new}}$  -  $\sigma$ ))
         $\sigma$  =  $\sigma_{\text{new}}$ 
        i = i + 1
        println("Concluded loop $i with error $error.")
    end
    return  $\sigma$ 
end
```

"Optimistic policy iteration routine."

```
function optimistic_policy_iteration(model; tolerance=1e-5, m=100)
    v = zeros(length(model.w_grid), length(model.y_grid))
    error = tolerance + 1
    while error > tolerance
        last_v = v
         $\sigma$  = get_greedy(v, model)
        for i in 1:m
            v = T_ $\sigma$ (v,  $\sigma$ , model)
        end
        error = maximum(abs.(v - last_v))
    end
    return get_greedy(v, model)
end
```

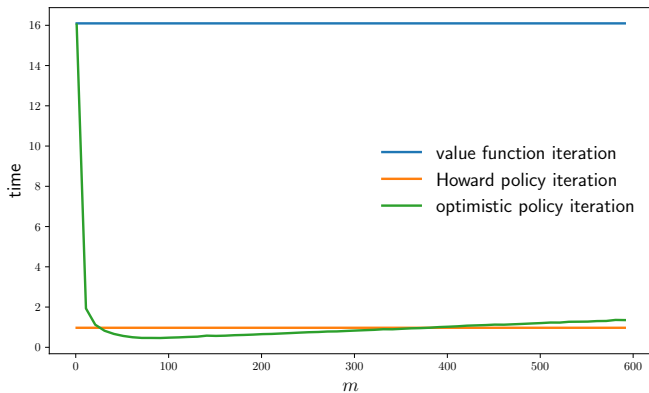


Figure: Timings for alternative algorithms