

中山大学 DCS5706 《随机过程及应用》期末作业

 $M/M/1$ 排队系统的控制研究

何瑞杰 25110801

目录

| | |
|------------------------------|----|
| 1. 问题描述 | 1 |
| 2. 模型建立 | 1 |
| 3. 最优策略计算 | 2 |
| 3.1. 代价函数 | 2 |
| 3.2. 值迭代 | 2 |
| 3.3. 策略迭代 | 3 |
| 3.4. 不同损失组合下的迭代结果 | 4 |
| 3.5. $M/M/1$ 排队系统的仿真模拟 | 4 |
| 4. 模型参数和折扣参数对最优策略的影响 | 5 |
| 5. 代码附录 | 6 |
| 5.1. $M/M/1$ 仿真模拟 | 6 |
| 5.2. 策略迭代和价值迭代 | 10 |

1. 问题描述

$M/M/1$ 排队系统广泛存在于生产生活中，它指的是一个先到先服务的单服务台的服务系统。顾客按照参数为 λ 的 Poisson 过程到达；服务台的服务时间服从参数为 μ （即服务速率）的指数分布，且和顾客的到达过程独立。

现考虑带有服务速率控制的 $M/M/1$ 排队系统，其服务速率 $\mu(i) \in (\lambda, \bar{\mu}]$ 取决于系统中的顾客数目 i ，该数目包括等待的顾客和正在服务的顾客。系统有两重成本：第一重为单位时间的服务成本 $q(\mu)$ ，其满足 $q(0) = 0$ ；第二重为顾客等待成本 $c(i)$ 。对该 $M/M/1$ 系统的控制目标为对系统内不同顾客数量时采用不同服务速率，以期最小化单位总成本。

2. 模型建立

带有控制的 $M/M/1$ 排队系统可使用连续时间 Markov 决策过程建模，其各参数如下：

| CTMDP 资料 | $M/M/1$ 系统中的元素 |
|----------------------|-------------------------|
| 状态 $x(t)$ | 系统中该时刻的顾客数目 i |
| 动作 $u(t)$ | 系统该时刻的服务速率 μ |
| 代价函数 $g(x(t), u(t))$ | 单位时间总成本 $q(\mu) + c(i)$ |
| 策略 μ_k | 系统的服务速度策略 $\mu(i)$ |

若转移速度对所有状态和动作均匀，则有

$$J_{\pi}(x_0) = \mathbb{E} \left[\sum_{k=0}^{\infty} \left(\frac{\nu}{\beta + \nu} \right)^k \frac{g(x_k, \mu_k(x_k))}{\beta + \nu} \right] = \mathbb{E} \left[\sum_{k=0}^{\infty} \alpha^k \cdot \tilde{g}(x_k, \mu_k(x_k)) \right],$$

对应的 Bellman 方程为

$$J(i) = \frac{1}{\beta + \nu} \min_{u \in U(i)} \left[g(i, u) + \nu \sum_j p_{i,j}(u) J(j) \right]$$

考虑转移速度对所有状态和动作不均匀, 但存在上界 ν , 若对状态 i 和动作 u , 有转移速度 $\nu_i(u)$, 考虑下面拥有新的转移概率的均匀转移速度的 CTMDP:

$$\tilde{p}_{i,j} = \begin{cases} \frac{\nu_i(u)}{\nu} p_{i,j}(u) & \text{if } i \neq j \\ \frac{\nu_i(u)}{\nu} p_{i,i}(u) + 1 - \frac{\nu_i(u)}{\nu} & \text{if } i = j \end{cases}$$

因此新的 CTMDP 的 Bellman 方程为

$$J(i) = \frac{1}{\beta + \nu} \min_{u \in U(i)} \left[g(i, u) + (\nu - \nu_i(u)) J(i) + \nu_i(u) \sum_j p_{i,j}(u) J(j) \right]$$

在 $M/M/1$ 队列中, 转移速率 $\nu_i(\mu)$ 在系统中无顾客 ($i = 0$) 时为 λ , 在有顾客时为 $\lambda + \mu$, 则依照上述结果的转移速率上界为 $\nu = \lambda + \bar{\mu}$ 。由于该系统的状态只可能向相邻状态转移, 且当系统中没有顾客时, 规定 $\mu(0) = 0$, 因此可以得到其 Bellman 方程为

$$J(i) = \begin{cases} \frac{1}{\beta + \nu} [c(0) + (\nu - \lambda) J(0) + \lambda J(1)] & i = 0 \\ \frac{1}{\beta + \nu} \min_{\mu} [c(i) + q(\mu) + (\nu - \lambda - \mu) J(i) + \lambda J(i+1) + \mu J(i-1)] & 1 \leq i < M \\ \frac{1}{\beta + \nu} \min_{\mu} [c(N) + q(\mu) + (\nu - \mu) J(N) + \mu J(N-1)] & i = N \end{cases}$$

注意系统中转移概率 $p_{i,i+1}(u)$ 对应着新顾客进入系统, 其值为 $\frac{\lambda}{\lambda + \mu}$, 而 $p_{i,i-1}(u)$ 对应着顾客服务完成离开系统, 其值为 $\frac{\mu}{\lambda + \mu}$ 。值得注意的是, 实际模拟中系统不可能有无限容量, 在此令系统容量为 M , 则当 $i = M$ 时, 对应于到达过程的等效速率变为 0。

3. 最优策略计算

本节介绍代价函数的取法和求解 Bellman 方程用到的算法。

3.1. 代价函数

本项目研究排队代价和服务代价分别为线性、二次函数、指数函数情况时的最优控制策略, 共有九种组合。具体地, 线性、二次代价和指数代价分别取

$$f_{\text{linear}}(x) = x, \quad f_{\text{quad}}(x) = \frac{1}{2}x^2, \quad f_{\text{exp}}(x) = e^{0.1x}.$$

3.2. 值迭代

第一种求解方法是值迭代, 其原理为直接应用 Bellman 方程的定义, 并用其迭代让边界处的值逐渐传导到其他各个状态, 直至收敛:

$$J_{k+1}(i) = \min_{u \in U(i)} \left[g(i, u) + \sum_{j=1}^n p_{i,j}(u) J_k(j) \right]$$

在 $M/M/1$ 系统中, 值迭代算法可以写为

Algorithm 1 Value Iteration for Controlled M/M/1 Queue

```

1: procedure Value-Iteration( $c(i)$ ,  $q(\mu)$ ,  $\lambda$ ,  $\bar{\mu}$ ,  $\beta$ ,  $\varepsilon$ )
2:    $\nu \leftarrow \lambda + \bar{\mu}$ 
3:    $J_0(i) \leftarrow 0, \forall i \in \{0, \dots, N\}$ 
4:    $k \leftarrow 0$ 
5:   while true do
6:      $\triangleright i = 0$ 
7:      $J_{k+1}(0) \leftarrow \frac{1}{\beta + \nu} [c(0) + (\nu - \lambda)J_{k(0)} + \lambda J_{k(1)}]$ 
8:      $\triangleright i = N$ 
9:      $J_{k+1}(N) \leftarrow \frac{1}{\beta + \nu} \min_{\mu} [c(N) + q(\mu) + (\nu - \mu)J(N) + \mu J(N - 1)]$ 
10:     $\triangleright$  Others
11:    for  $i \leftarrow 1, \dots, N - 1$  do
12:       $J_{k+1}(i) \leftarrow \frac{1}{\beta + \nu} \min_{\mu \in (\lambda, \bar{\mu}]} [c(i) + q(\mu) + \mu J_{k(i-1)} + (\nu - \lambda - \mu)J_{k(i)} + \lambda J_{k(i+1)}]$ 
13:    end
14:    if  $\max_i |J_{k+1}(i) - J_{k(i)}| < \varepsilon$  then
15:      break
16:    end
17:     $k \leftarrow k + 1$ 
18:  end
19:   $\mu^*(i) \leftarrow \arg \min_{\mu} [q(\mu) - \mu(J_{k+1}(i) - J_{k+1}(i - 1))], \forall i \geq 1$ 
20:  return ( $J_{k+1}, \mu^*$ )
21: end

```

3.3. 策略迭代

还可以通过策略迭代算法解 Bellman 方程。其核心为从一个初始策略 $\mu^{(0)}$ 出发，每个迭代循环中，通过策略评估得到当前策略对应的价值函数 $J_{\mu^{(i-1)}}$ ，然后根据这个价值函数贪心地取得新的策略 $\mu^{(i)}$ ，直至策略收敛。在策略评估过程中，将 \min_{μ} 直接替换为 $\mu_i^{(k)}$ 即在当前策略下的 μ 。这样 Bellman 方程就变为一个线性方程组 $AJ = b$ ，其中 A 是三对角矩阵，可以使用数值计算包高效求解。

Algorithm 2 Policy Iteration for Controlled M/M/1 Queue

```

1: procedure Policy-Iteration( $c(i)$ ,  $q(\mu)$ ,  $\lambda$ ,  $\bar{\mu}$ ,  $\beta$ )
2:    $\mu_0(i) \leftarrow \bar{\mu}, \forall i \in \{1, \dots, N\}$ 
3:    $k \leftarrow 0$ 
4:   while true do
5:     Solve linear system  $AJ_{\mu_k} = b$ .
6:     for  $i \leftarrow 1, \dots, N$  do
7:        $\Delta J \leftarrow J_{\mu_k}(i) - J_{\mu_k}(i - 1)$ 
8:        $\mu_{k+1}(i) \leftarrow \arg \min_{\mu \in (\lambda, \bar{\mu}]} [q(\mu) - \mu \cdot \Delta J]$ 
9:     end
10:    if  $\mu_{k+1}(i) = \mu_k(i), \forall i$  then
11:      break
12:    end
13:     $k \leftarrow k + 1$ 
14:  end
15:  return ( $J_{\mu_k}, \mu_k$ )
16: end

```

3.4. 不同损失组合下的迭代结果

实际测试中，到达速率 $\lambda = 10$ ，系统最大容量 $N = 100$ ，折扣参数 $\beta = 0.01$ 。值迭代和策略迭代收敛到同样的策略，但值迭代相比策略迭代慢得多，因此这里使用策略迭代。通过最优策略可以计算得到 Q 矩阵，进而计算出系统的稳态分布，最后得到系统的平均代价，实测这九个损失组合的最优策略的平均代价与实际模拟得出的代价基本一致（见附录）。

将不同损失组合下迭代得到的最优策略函数绘制如下，可见当服务损失时线性，而排队损失时线性、二次或指数时，最优策略在系统顾客数量较低时迅速增长至最大服务速率。当排队损失为线性时，不论服务损失时二次或是指数，最优策略随系统中顾客数量增长对应的系统服务速度增长较为缓慢。其他情况下，随系统中顾客数量的增长，最优策略下的服务速率先是立刻以大斜率增加，然后缓慢或线性增加，至系统中顾客人数在总容纳量一半左右时到达最高服务速率。

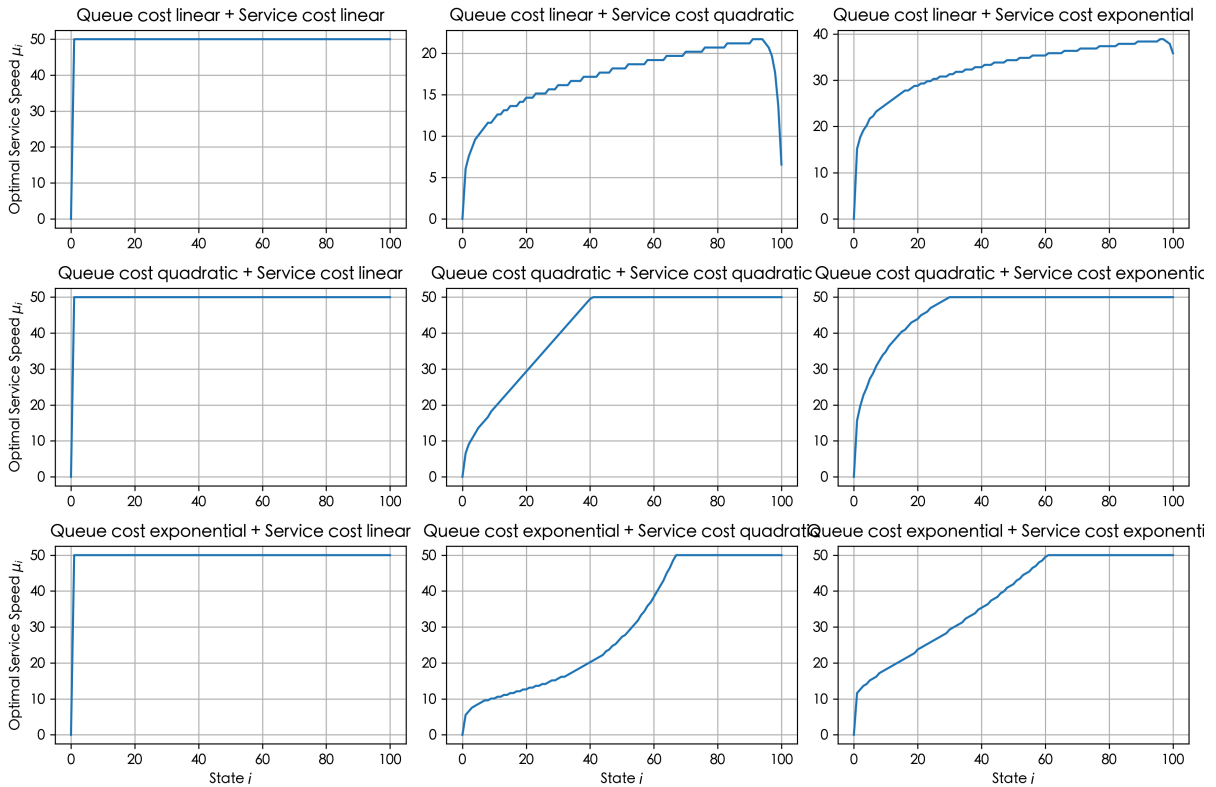


图 1 不同损失组合下的最优服务速度策略

3.5. $M/M/1$ 排队系统的仿真模拟

本节简述仿真模拟的逻辑。由于 CTMDP 的跳变总是发生在瞬间，除了跳变的时刻外其他时刻该过程的状态均在跳变间隔中恒定不变，因此可以采用离散事件法对该系统进行建模。具体地，将系统状态 i 表示为当前系统中顾客人数，将系统服务速度 μ_i 表示为当前系统服务速率，将系统到达速率 λ 表示为顾客到达系统的平均间隔时间的倒数，将系统最大容量 N 表示为系统最多可以容纳的顾客人数。

维护一个事件优先队列。仿真模拟开始时，系统中顾客人数为零，系统服务速率为零，并在队列中添加一个新的到达事件，事件距离系统当前时刻的间隔采样自到达过程的间隔分布。仿真系统处理完每个事件后，将会直接跳转至下一个事件的发生时刻。如果该事件是一个到达事件，系统中顾客数 $+1$ ，如果达前系统为空，系统开始服务该顾客，并在队列中添加一个服务完成时间的事件，即离开事件；如果到达前系统已满，则忽略该到达事件。如果该事件是离开事件，系

统中顾客数 -1 ，如果离开后系统中顾客人数不为零，将在事件序列中添加新的离开事件；否则将系统服务速率降低至 0。在系统处理每一个事件时，都会计算自上一个事件以来的代价函数。如果事件队列中最近的事件发生时间超过了提前设定的仿真时长，系统终止，并根据仿真时长和总损失计算平均损失。

4. 模型参数和折扣参数对最优策略的影响

本节中固定排队损失为线性，服务损失为二次函数，研究系统参数中到达速率 λ 、系统最大容量 N 、折扣参数 β 对最优策略的影响。考虑 $\lambda \in \{5, 10, 20\}$ ， $N \in \{100, 1000, 2000\}$ ， $\beta \in \{0.001, 0.01, 0.1\}$ ，下图显示不同组合下的最优服务速度策略。我们发现一个有趣的情况。当 λ 较大或 β 较小时，系统的最优策略下服务速度随着系统内顾客数量增加而下降。

但 N 较大时，即使 λ 较大并不会产生上述情形，这说明系统的容量会显著影响最优策略。另外注意到当 β 较大时，即使增大系统容量，系统最优策略在大顾客量时下降的现象依然出现，这说明折扣参数会显著影响系统策略。

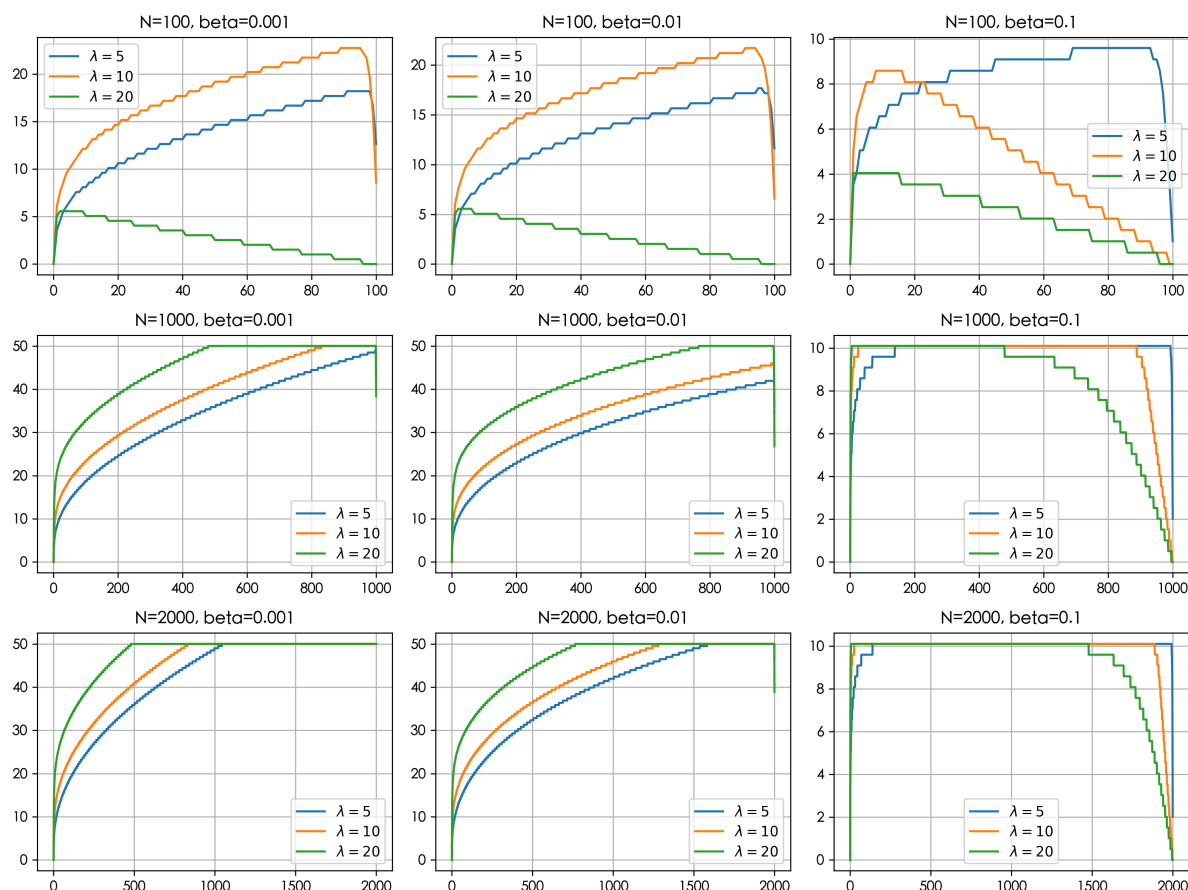


图 2 不同损失组合下的最优服务速度策略

5. 代码附录

5.1. $M/M/1$ 仿真模拟

Python

```

1  # mml.py
2  import heapq
3  import os
4  import numpy as np
5  from collections import namedtuple
6  from typing import Callable, List, Tuple, Optional
7  import matplotlib.pyplot as plt
8  import matplotlib as mpl
9  from tqdm import tqdm
10
11 plt.rcParams['font.family'] = list({"SimHei", "Heiti TC"} & set(f.name for f in
    mpl.font_manager.fontManager.ttflist))[0]
12
13 Event = namedtuple('Event', ['time', 'type', 'data'])
14 # type: 'arrival', 'departure', 'rate_change'
15
16 class ControlledMM1Queue:
17
18     def __init__(self,
19                 arrival_rate: float,
20                 service_rate_policy: Callable[[int], float],
21                 service_cost_func: Callable[[float], float],
22                 queue_cost_func: Callable[[int], float],
23                 max_customers: int = None):
24         self.arrival_rate = arrival_rate
25         self.service_rate_policy = service_rate_policy
26         self.service_cost_func = service_cost_func
27         self.queue_cost_func = queue_cost_func
28         self.max_customers = max_customers
29
30         # system state
31         self.current_time = 0.0
32         self.num_customers = 0
33         self.current_service_rate = 0.0
34
35         # event priority queue
36         self.event_queue: List[Event] = []
37
38         # history
39         self.history = {
40             'time': [0.0],
41             'num_customers': [0],
42             'service_rate': [0.0],

```

```

43         'total_cost': [0.0]
44     }
45     self.total_cost = 0.0
46     self.last_event_time = 0.0
47
48     def reset(self):
49         self.current_time = 0.0
50         self.num_customers = 0
51         self.current_service_rate = 0.0
52         self.event_queue = []
53         self.history = {
54             'time': [0.0],
55             'num_customers': [0],
56             'service_rate': [0.0],
57             'total_cost': [0.0]
58         }
59         self.total_cost = 0.0
60         self.last_event_time = 0.0
61
62     def _add_event(self, time: float, event_type: str, data=None):
63         heapq.heappush(self.event_queue, Event(time, event_type, data))
64
65     def _exponential_sample(self, rate: float) -> float:
66         if rate <= 0:
67             return float('inf')
68         return np.random.exponential(1.0 / rate)
69
70     def _schedule_arrival(self):
71         inter_arrival = self._exponential_sample(self.arrival_rate)
72         next_arrival_time = self.current_time + inter_arrival
73         self._add_event(next_arrival_time, 'arrival')
74
75     def _schedule_departure(self):
76         if self.num_customers > 0:
77             # schedule next departure time
78             service_rate = self.service_rate_policy(self.num_customers)
79             service_time = self._exponential_sample(service_rate)
80             next_departure_time = self.current_time + service_time
81             self._add_event(next_departure_time, 'departure')
82             self.current_service_rate = service_rate
83         else:
84             # system is empty
85             self.current_service_rate = 0.0
86
87     def _update_cost(self):
88         dt = self.current_time - self.last_event_time

```

```

89         if dt > 0:
90             # add period cost
91             queue_cost = self.queue_cost_func(self.num_customers) * dt
92             service_cost = self.service_cost_func(self.current_service_rate) *
                dt
93             self.total_cost += queue_cost + service_cost
94
95         # update history
96         self.history['time'].append(self.current_time)
97         self.history['num_customers'].append(self.num_customers)
98         self.history['service_rate'].append(self.current_service_rate)
99         self.history['total_cost'].append(self.total_cost)
100
101         self.last_event_time = self.current_time
102
103     def _handle_arrival(self):
104         self._update_cost()
105
106         # maximum capacity check
107         if self.max_customers is None or self.num_customers <
            self.max_customers:
108             self.num_customers += 1
109
110         # start service when system turns to non-empty
111         if self.num_customers == 1:
112             self._schedule_departure()
113
114         # schedule next arrival
115         self._schedule_arrival()
116
117     def _handle_departure(self):
118         self._update_cost()
119
120         if self.num_customers > 0:
121             self.num_customers -= 1
122
123         # serve next customer if system is non-empty
124         if self.num_customers > 0:
125             self._schedule_departure()
126         else:
127             self.current_service_rate = 0.0
128
129     def run(self, T: float) -> dict:
130         self.reset()
131
132         self._schedule_arrival()

```



```
133
134     while self.event_queue:
135         # check to the nearest event
136         event = heapq.heappop(self.event_queue)
137
138         # check time limit
139         if event.time > T:
140             break
141
142         # jump to the nearest event
143         self.current_time = event.time
144
145         # handle event
146         if event.type == 'arrival':
147             self._handle_arrival()
148         elif event.type == 'departure':
149             self._handle_departure()
150
151         print(f"Simulating... {self.current_time/T*100:.2f}", end='\r')
152
153         # update current time and cost
154         self.current_time = T
155         self._update_cost()
156
157         return self.history
158
159 def get_average_cost(self, T: float = None) -> float:
160     if T is None:
161         T = self.history['time'][-1]
162     return self.total_cost / T
```

5.2. 策略迭代和价值迭代

```

1  # solve_bellman.py
2
3  import numpy as np
4  from typing import Callable, List
5  import matplotlib.pyplot as plt
6  import matplotlib as mpl
7  from tqdm import tqdm
8  from scipy.sparse import diags
9  from scipy.sparse.linalg import spsolve
10
11  plt.rcParams['font.family'] = list({"SimHei", "Heiti TC"} & set(f.name for f in
    mpl.font_manager.fontManager.ttflist))[0]
12
13
14  class CTMDPControlledQueue:
15
16      def __init__(self,
17                  lambda_rate: float,
18                  max_state: int,
19                  c_func: Callable[[int], float],
20                  q_func: Callable[[float], float],
21                  mu_space: List[float],
22                  beta: float = 0.1):
23          self.lambda_rate = lambda_rate
24          self.max_state = max_state
25          self.c_func = c_func
26          self.q_func = q_func
27          self.mu_space = np.array(sorted(mu_space))
28          self.beta = beta
29
30          # nu upper bound
31          self.V = lambda_rate + np.max(self.mu_space)
32
33          assert np.any(np.isclose(self.mu_space, 0.0))
34
35          self._c_vec = np.array([c_func(i) for i in range(max_state + 1)])
36          self._q_vec = np.array([q_func(mu) for mu in self.mu_space])
37
38      def _bellman_rhs(self, i: int, mu: float, J: np.ndarray) -> float:
39          # calculate the RHS of Bellman equation
40          if i == 0:
41              # constraint at 0
42              return self._c_vec[i] + (self.V - self.lambda_rate) * J[i] +
                  self.lambda_rate * J[i+1]
43          elif i == self.max_state:

```

Python

```

44         # constraint at N
45         return self._c_vec[i] + self.q_func(mu) + mu * J[i-1] + (self.V -
46             mu) * J[i]
47     else:
48         # other cases
49         return (self._c_vec[i] + self.q_func(mu) +
50             mu * J[i-1] + (self.V - self.lambda_rate - mu) * J[i] +
51             self.lambda_rate * J[i+1])
52
53 def value_iteration(self, tolerance: float = 1e-6, max_iter: int = 20000) -
54 > tuple:
55     J = np.zeros(self.max_state + 1)
56     policy = np.zeros(self.max_state + 1)
57
58     pb = tqdm(range(max_iter))
59     for iteration in pb:
60         J_new = np.zeros_like(J)
61
62         for i in range(self.max_state + 1):
63             # iterate w.r.t. Bellman equation definition
64             if i == 0:
65                 J_new[i] = self._bellman_rhs(i, 0.0, J) / (self.beta +
66                     self.V)
67                 policy[i] = 0.0
68             else:
69                 rhs_values = np.array([
70                     self._bellman_rhs(i, mu, J) for mu in self.mu_space
71                 ])
72                 best_idx = np.argmin(rhs_values)
73                 policy[i] = self.mu_space[best_idx]
74                 J_new[i] = rhs_values[best_idx] / (self.beta + self.V)
75
76         # check convergence
77         diff = np.max(np.abs(J_new - J))
78         if diff < tolerance:
79             print(f"Value iteration converges at {iteration}th iteration
80                 with J difference {diff:.2e}")
81             break
82
83         pb.set_postfix({"Diff": f"{diff:.3e}"})
84         J = J_new
85
86         if iteration == max_iter - 1:
87             print("Reached maximum iterations")
88
89     return J, policy

```

```

86
87     def policy_iteration(self, max_iter: int = 500) -> tuple:
88         num_states = self.max_state + 1
89         policy = np.full(num_states, self.mu_space[0])
90         J = np.zeros(num_states)
91
92         for iteration in range(max_iter):
93             # setting up linear system AJ = b
94             main_diag = np.full(num_states, self.beta + self.V)
95             upper_diag = np.full(num_states - 1, -self.lambda_rate)
96             lower_diag = np.full(num_states - 1, 0.0)
97             b = np.zeros(num_states)
98
99             # case 0
100            b[0] = self._c_vec[0]
101            main_diag[0] = self.beta + self.lambda_rate
102            # upper_diag[0] = -self.lambda_rate 已设置
103
104            # case 1..N-1
105            for i in range(1, self.max_state):
106                mu_i = policy[i]
107                lower_diag[i-1] = -mu_i
108                #  $(\beta+V)J_i - \mu J_{i-1} - (V-\lambda-\mu)J_i - \lambda J_{i+1} = c_i + q(\mu)$ 
109                #  $\Rightarrow -(\mu)J_{i-1} + (\beta+\lambda+\mu)J_i - \lambda J_{i+1} = c_i + q(\mu)$ 
110                main_diag[i] = self.beta + self.lambda_rate + mu_i
111                b[i] = self._c_vec[i] + self.q_func(mu_i)
112
113            # case N
114            mu_N = policy[self.max_state]
115            lower_diag[self.max_state - 1] = -mu_N
116            main_diag[self.max_state] = self.beta + mu_N
117            b[self.max_state] = self._c_vec[self.max_state] + self.q_func(mu_N)
118
119            # solve for AJ = b
120            A = diags([lower_diag, main_diag, upper_diag], [-1, 0, 1],
121                    format='csc')
122            J_new = spsolve(A, b)
123
124            # evaluate bellman equation
125            if iteration == 0:
126                for i in range(min(3, num_states)):
127                    if i == 0:
128                        lhs = (self.beta + self.lambda_rate) * J_new[i] -
129                            self.lambda_rate * J_new[i+1]
130                        rhs = b[i]
131                    elif i == self.max_state:

```

```

130             mu = policy[i]
131             lhs = -mu * J_new[i-1] + (self.beta + mu) * J_new[i]
132             rhs = b[i]
133         else:
134             mu = policy[i]
135             lhs = -mu * J_new[i-1] + (self.beta + self.lambda_rate
136             + mu) * J_new[i] - self.lambda_rate * J_new[i+1]
137             rhs = b[i]
138
139         # policy improvement
140         new_policy = np.zeros(num_states)
141         new_policy[0] = 0.0
142
143         for i in range(1, num_states):
144             if i == self.max_state:
145                 # edge case
146                 rhs_values = np.array([
147                     self._c_vec[i] + self.q_func(mu) + mu * J_new[i-1] +
148                     (self.V - mu) * J_new[i]
149                     for mu in self.mu_space
150                 ])
151             else:
152                 rhs_values = np.array([
153                     self._c_vec[i] + self.q_func(mu) + mu * J_new[i-1] +
154                     (self.V - self.lambda_rate - mu) * J_new[i] +
155                     self.lambda_rate * J_new[i+1]
156                     for mu in self.mu_space
157                 ])
158             best_idx = np.argmin(rhs_values)
159             new_policy[i] = self.mu_space[best_idx]
160
161         # convergence check
162         policy_diff = np.max(np.abs(new_policy - policy))
163         J_diff = np.max(np.abs(J_new - J))
164
165         if policy_diff < 1e-8 and J_diff < 1e-8:
166             residual = self.compute_bellman_residual(J_new, new_policy)
167             break
168
169         policy = new_policy
170         J = J_new.copy()
171
172     return J, policy
173
174 def compute_bellman_residual(self, J: np.ndarray, policy: np.ndarray) ->
175 float:

```

```

172     """Calculate Bellman Residual"""
173     residual = 0.0
174
175     for i in tqdm(range(self.max_state + 1), desc="Evaluating Bellman
176 residual"):
177         rhs = self._bellman_rhs(i, policy[i], J)
178         T_J = rhs / (self.beta + self.V)
179         residual = max(residual, abs(J[i] - T_J))
180
181     return residual
182
183 def compare_policies(self, J_vi: np.ndarray, policy_vi: np.ndarray,
184                      J_pi: np.ndarray, policy_pi: np.ndarray,
185                      num_states: int = 20):
186
187     print(f"J diff: {np.max(np.abs(J_vi - J_pi)):.10f}")
188     print(f"Policy diff: {np.max(np.abs(policy_vi - policy_pi)):.10f}")
189
190     print("\n State | J_VI(i) | J_PI(i) | ΔJ_VI | ΔJ_PI | μ*_VI |
191 μ*_PI")
192     print("-"*80)
193
194     mismatch_count = 0
195     for i in range(min(num_states, self.max_state) + 1):
196         delta_J_vi = J_vi[i] - J_vi[i-1] if i > 0 else 0.0
197         delta_J_pi = J_pi[i] - J_pi[i-1] if i > 0 else 0.0
198
199         print(f" {i:3d} | {J_vi[i]:8.4f} | {J_pi[i]:8.4f} |
200 {delta_J_vi:6.2f} | "
201               f"{delta_J_pi:6.2f} | {policy_vi[i]:4.1f} |
202 {policy_pi[i]:4.1f} |")
203
204     if mismatch_count == 0:
205         print("Test Passed")
206     else:
207         print(f"Testing Failed with {mismatch_count} different states out
208 of {num_states}")
209
210     residual_vi = self.compute_bellman_residual(J_vi, policy_vi)
211     residual_pi = self.compute_bellman_residual(J_pi, policy_pi)
212     print(f"\nBellman Residual VI: {residual_vi:.2e}, PI:
213 {residual_pi:.2e}")
214
215 def compute_average_cost_steady_state(solver, policy):
216     lambda_rate = solver.lambda_rate
217     max_state = solver.max_state

```

```

213     # construct Q matrix
214     Q = np.zeros((max_state+1, max_state+1))
215
216     for i in range(max_state+1):
217         mu = policy[i]
218
219         if i < max_state:
220             Q[i, i+1] = lambda_rate # arrival
221
222         if i > 0:
223             Q[i, i-1] = mu # service
224
225         # departure
226         Q[i, i] = -(lambda_rate + mu)
227
228     # solve  $\pi Q = 0, \sum \pi = 1$ 
229     eigenvals, eigenvecs = np.linalg.eig(Q.T)
230     zero_idx = np.argmin(np.abs(eigenvals))
231     pi = np.abs(eigenvecs[:, zero_idx].real)
232     pi /= np.sum(pi)
233
234     # average cost
235     costs = np.array([
236         solver.c_func(i) + solver.q_func(policy[i])
237         for i in range(max_state+1)
238     ])
239     average_cost = np.dot(pi, costs)
240     return average_cost
241
242 if __name__ == "__main__":
243     # LAMBDA = 10.0
244     # MAX_STATE = 100
245     # BETA = 0.01
246
247     def linear(x):
248         return x
249
250     def quadratic(x):
251         return 0.5 * x ** 2
252
253     def exponential(x):
254         return np.exp(0.1 * x)
255
256
257     mu_space = np.linspace(0, 50, 100)
258

```

```

259     ls = []
260
261     # for c_func in [linear, quadratic, exponential]:
262     #     for q_func in [linear, quadratic, exponential]:
263
264     #         solver = CTMDPControlledQueue(
265     #             lambda_rate=LAMBDA,
266     #             max_state=MAX_STATE,
267     #             c_func=c_func,
268     #             q_func=q_func,
269     #             mu_space=mu_space,
270     #             beta=BETA
271     #         )
272
273     #         J_pi, policy_pi = solver.policy_iteration()
274     #         ls.append((c_func.__name__, q_func.__name__, J_pi, policy_pi))
275
276     #         from main import ControlledMM1Queue
277
278     #         def service_rate_policy(a):
279     #             return policy_pi[a]
280
281     #         queue = ControlledMM1Queue(
282     #             arrival_rate=LAMBDA,
283     #             service_rate_policy=service_rate_policy,
284     #             service_cost_func=q_func,
285     #             queue_cost_func=c_func,
286     #             max_customers=MAX_STATE
287     #         )
288
289     #         queue.run(1e5)
290     #         avg_cost = compute_average_cost_steady_state(solver, policy_pi)
291     #         print(f"Stationary average cost (steady state): {avg_cost:.4f}")
292     #         print(f"Stationary average cost (simulation):
293     #             {queue.get_average_cost():.4f}")
294
295     #         plt.figure(figsize=(12, 8), dpi=300)
296
297     #         for i, (c_name, q_name, J_pi, policy_pi) in enumerate(ls):
298     #             plt.subplot(3, 3, i+1)
299     #             plt.plot(policy_pi)
300     #             plt.title(f"Queue cost {c_name} + Service cost {q_name}")
301     #             plt.grid()
302
303     #             if i > 5:
304     #                 plt.xlabel("State $i$")

```



```

304
305     #     if i % 3 == 0:
306     #         plt.ylabel("Optimal Service Speed  $\mu_i$ ")
307
308     # plt.tight_layout()
309     # plt.savefig("ctmdp_value_iteration.png")
310
311     c_func = linear
312     q_func = quadratic
313
314     for N in [100, 1000, 2000]:
315         for beta in [0.001, 0.01, 0.1]:
316             for LAMBDA in [5, 10, 20]:
317                 solver = CTMDPControlledQueue(
318                     lambda_rate=LAMBDA,
319                     max_state=N,
320                     c_func=c_func,
321                     q_func=q_func,
322                     mu_space=mu_space,
323                     beta=beta
324                 )
325
326                 J_pi, policy_pi = solver.policy_iteration()
327                 ls.append((N, beta, LAMBDA, J_pi, policy_pi))
328
329     fig, axs = plt.subplots(3, 3, figsize=(12, 9), dpi=300)
330     axs = axs.flatten()
331
332     for i, (N, beta, LAMBDA, J_pi, policy_pi) in enumerate(ls):
333         print(i // 3, N, beta, LAMBDA, policy_pi.shape)
334         axs[i//3].plot(policy_pi, label=f" $\lambda = \{LAMBDA\}$ ")
335         axs[i//3].set_title(f"N={N}, beta={beta}")
336         axs[i//3].grid()
337         axs[i//3].legend()
338
339     fig.tight_layout()
340     fig.savefig("ctmdp_policy_iteration_1.png")
341
342     # print("\n" + "="*80)
343     # print("Policy iteration results")
344     # print("="*80)
345     # solver.compare_policies(J_vi, policy_vi, J_pi, policy_pi)
346
347     # #
348     # # np.savez('ctmdp_converged_results.npz',
349     # #         J_vi=J_vi, policy_vi=policy_vi,

```

```
350     # #           J_pi=J_pi, policy_pi=policy_pi,
351     # #           lambda_rate=LAMBDA, beta=BETA, max_state=MAX_STATE)
352
353     # avg_cost = compute_average_cost_steady_state(solver, policy_pi)
354     # print(f"Stationary average cost: {avg_cost:.4f}")
355
356
```