

1 实验说明

1.1 实验目的

利用不同算法设计策略来求解组合优化问题。

1.2 实验内容

假设有 n 个物品和一个背包，每个物品重量为 $0 < w_i < 100$, $i = 1, 2, \dots, n$, 价值为 $0 < v_i < 100$, $i = 1, 2, \dots, n$ 。背包最大容量为 c 。请问该如何选择物品才能使装入背包中的物品总价值最大？最大价值是多少？请按照如下要求完成算法：

1. 请利用分治法来求解该问题，给出最优解值以及求解时间；
2. 请利用动态规划算法来求解该问题，给出得到的最优解值以及求解时间；
3. 请利用贪心算法来求解该问题，给出得到的最优解值以及求解时间；
4. 请利用回溯法来求解该问题，给出得到的最优解值以及求解时间；
5. 请利用分支限界法来求解该问题，给出得到的最优解值以及求解时间；
6. 请利用蒙特卡洛算法来求解该问题，给出得到的最优解值以及求解时间；
7. 请利用深度强化学习算法来求解该问题，给出最优解值以及求解时间；
8. 给定相同输入，比较上述算法得到的最优解值和求解时间。当 n 比较大的时候，上述算法运算时间可能很长，请在算法中增加终止条件以确保在有限时间内找到最优解的值。

2 问题描述及建模

0-1 背包问题 (0-1 Knapsack) 是一种 NP 完全的组合优化问题。定义 x_i 为第 i 个物品是否放入背包中，则其公式化描述如下：

$$\begin{aligned} \max \quad & \sum_{i=1}^n v_i x_i \\ \text{s.t.} \quad & \sum_{i=1}^n w_i x_i \leq c \\ & x_i \in \{0, 1\}, i = 1, 2, \dots, n \end{aligned}$$

解空间 $X = \{(x_1, x_2, \dots, x_n) \mid x_i \in \{0, 1\}, i = 1, 2, \dots, n\}$ 的大小为 2^n ，故用穷举法求解该问题的时间复杂度为 $O(2^n)$ 。

2.1 分治法 (Divide and Conquer)

定义子问题 $m(i, c)$ 为：从第 i 个至第 n 个物品中选取若干物品放入容量为 c 的背包中，使得总价值最大，即

$$m(i, c) = \max \left\{ \sum_{j=i}^n v_j x_j \mid \sum_{j=i}^n w_j x_j \leq c, x_j \in \{0, 1\}, j = i, i+1, \dots, n \right\}$$

则原问题为 $m(1, c)$ 。注意到 $m(1, c)$ 具有最优子结构性质：令 (x_1, x_2, \dots, x_n) 为 $m(1, c)$ 的一个最优解，则 (x_2, x_3, \dots, x_n) 是 $m(2, c - w_1 x_1)$ 的最优解。否则，假设 (y_2, y_3, \dots, y_n) 是 $m(2, c - w_1 x_1)$ 的最优解，有 $\sum_{j=2}^n v_j y_j > \sum_{j=2}^n v_j x_j$ 且 $\sum_{j=2}^n w_j y_j \leq c - w_1 x_1$ ，可以推出 $v_1 x_1 + \sum_{j=2}^n v_j y_j > v_1 x_1 + \sum_{j=2}^n v_j x_j$ 且 $w_1 x_1 + \sum_{j=2}^n w_j y_j \leq c$ ，与 (x_1, x_2, \dots, x_n) 是 $m(1, c)$ 的最优解矛盾。由此，有递归式

$$m(i, c) = \begin{cases} \max\{m(i+1, c), v_i + m(i+1, c - w_i)\}, & \text{if } w_i \leq c \\ m(i+1, c), & \text{otherwise} \end{cases} \quad i = 1, 2, \dots, n-1$$

边界条件为 $m(i, c) = v_i \mathbb{I}(w_i \leq c)$ 。

当 $c = 0$, $m(i, c) = 0$, v_i, w_i 均为整数时，递归树的大小为 $O(nc)$ ，故分治法求解该问题的时间复杂度为 $O(nc)$ 。伪代码如下

Algorithm 1: Divide and Conquer

```
Input:  $n, c, w[1..n], v[1..n]$ .
Output:  $\text{KNAPSACK}(1, c)$ .
function  $\text{KNAPSACK}(i, c)$ :
  1. if  $i = n$  then
    2. return  $v_n \mathbb{I}(w_n \leq c)$ 
  3. if  $w_i > c$  then
    4. return  $\text{KNAPSACK}(i+1, c)$ 
  5. return  $\max\{\text{KNAPSACK}(i+1, c), v_i + \text{KNAPSACK}(i+1, c - w_i)\}$ 
```

2.2 动态规划 (Dynamic Programming)

上述分治法可能会重复求解一些问题。为了避免这种情况，可以使用带备忘录的递归算法 (也称带备忘录的动态规划)，以空间换取时间；或使用自底向上的动态规划。伪代码如下

Algorithm 2: Dynamic Programming

```
Input:  $n, c, w[1..n], v[1..n]$ .
Output:  $m[1, c]$ .
1. Initialize an  $2D$  array  $m[1..n, 0..c]$  with  $0$ .
2. for  $j = w_n$  to  $c$  do
3.    $m[n, j] \leftarrow v_n$ 
4. for  $i = n-1$  downto  $1$  do
5.   for  $j = 0$  to  $w_i - 1$  do
6.      $m[i, j] \leftarrow m[i+1, j]$ 
7.   for  $j = w_i$  to  $c$  do
8.      $m[i, j] \leftarrow \max\{m[i+1, j], v_i + m[i+1, j - w_i]\}$ 
9. return  $m[1, c]$ 
```

应当声明的是，上述算法逐一求解子问题，但并非所有子问题的解均被使用，这意味着可能存在多余冗余计算。此外，该算法只适用于 c 与 w_i , $i = 1, 2, \dots, n$ 均为整数的情形。

2.3 贪心算法 (Greedy Algorithm)

考虑 0-1 背包问题的特殊情况。如果 $w_1 \leq w_2 \leq \dots \leq w_n$ 且 $v_1 \geq v_2 \geq \dots \geq v_n$ ，则问题具有贪心解性质：每次求解 $m(i, c)$ ，只要 $w_i \leq c$ ，就选择 $x_i = 1$ ，进而求解唯一子问题 $m(i+1, c - w_i)$ ；否则，选择 $x_i = 0$ ，进而求解唯一子问题 $m(i+1, c)$ 。换言之，我们逐一选择单位重量价值 $\frac{v_i}{w_i}$ 最大的物品 i ，直至背包装满。下证，该贪心解最优。

令 (x_1, x_2, \dots, x_n) 为问题 $m(1, c)$ 的一个最优解。(1) 若 $x_1 = 1$ ，由可行性，必有 $w_1 \leq c$ 。显然，该贪心符合贪心策略。(2) 若 $x_1 = 0$ 且 $w_1 > c$ ，该解同样符合贪心策略。(3) 若 $x_1 = 0$ 且 $w_1 \leq c$ ，令 $x_k = 1$ 为最优解中标最小的非零分量，则必有 $w_k = w_1$, $v_k = v_1$ 。从而 $(1, x_2, \dots, x_{k-1}, 0, x_{k+1}, \dots, x_n)$ 同样为最优解且符合贪心策略；否则，若 $w_k > w_1$, $v_k < v_1$ ，易知 $(1, x_2, \dots, x_{k-1}, 0, x_{k+1}, \dots, x_n)$ 仍为可行解且总价值更大，与 (x_1, x_2, \dots, x_n) 为最优解矛盾。综上所述，贪心策略可以得到一个最优解。

对于一般情况，首先将物品按单位重量价值排序，使得 $\frac{v_1}{w_1} \geq \frac{v_2}{w_2} \geq \dots \geq \frac{v_n}{w_n}$ 。若此时满足 $w_1 \leq w_2 \leq \dots \leq w_n$ 且 $v_1 \geq v_2 \geq \dots \geq v_n$ ，则贪心算法可以得到最优解。伪代码如下

Algorithm 3: Greedy Algorithm

```
Input:  $n, c, w[1..n], v[1..n]$ .
Output:  $m$ .
1. Let  $k_i$  be the sorted index of  $\frac{v[i]}{w[i]}$  such that  $\frac{v[k_1]}{w[k_1]} \leq \frac{v[k_2]}{w[k_2]} \leq \dots \leq \frac{v[k_n]}{w[k_n]}$ .
2.  $m \leftarrow 0$ 
3. for  $i = 1$  to  $n$  do
4.   if  $w[k_i] \leq c$  then
5.      $m \leftarrow m + v[k_i]$ 
6.      $c \leftarrow c - w[k_i]$ 
```

上述算法的时间复杂度取决于排序策略。若不存在排列使得 $w_1 \leq w_2 \leq \dots \leq w_n$ 且 $v_1 \geq v_2 \geq \dots \geq v_n$ 同时满足，则单位重量价值贪心策略不一定能得到最优解。同样，重量贪心策略 (按 w_i 排序) 与价值贪心策略 (按 v_i 排序) 也不一定能得到最优解。

2.4 回溯法 (Backtracking)

在分治法中，我们实际上以递归的方式对解空间树进行了深度优先搜索 (Depth-First Search)。然而，每次调用函数 $\text{KNAPSACK}(i, c)$ ，我们均分配新的内存用于记录当前参数 i, c 与函数返回值 $m(i, c)$ ，造成大量空间浪费。

事实上，若将当前待选物品 i 与背包剩余容量 c 视为问题的状态，以全局变量的形式存储，每次递归调用结束后，我们仍能以 $O(1)$ 时间回溯至上一个状态，以避免比较函数返回值。我们选择不回溯的过程中构造最优解，而是自顶向下地记录搜索路径和累计价值 cv ，更新最优值 m 。

经过上述调整，我们得到了基本的回溯算法。可行性约束函数 (Constraint) 缩小了解空间树的规模。经过调整，我们发现可以用如下限界函数 (Bound) 进一步将解空间树剪枝：若 $\sum_{i=1}^n v_i x_j + \sum_{j=i+1}^n v_j \leq m$ ，即当前累计价值与子问题的价值上界之和不超过当前最优解，则不再搜索该子树。

我们用数组 $r[1..n]$ 记录预处理得到子问题的价值上界。伪代码如下

Algorithm 4: Backtracking

```
Input:  $n, c, w[1..n], v[1..n]$ .
Output:  $m$ .
1.  $i \leftarrow 1; m \leftarrow 0; cv \leftarrow 0; r[n+1] \leftarrow 0$ 
2. for  $j = n$  downto  $1$  do
3.    $r[j] \leftarrow r[j+1] + v[j]$ 
4.    $\text{KNAPSACK}()$ 

function  $\text{KNAPSACK}()$ :
  1. if  $i > n$  then
    2. return
  3. if  $cv + r[i] \leq m$  then // Bound
    4. return
  5. if  $w[i] \leq c$  then // Constraint
    6.  $c \leftarrow c - w[i]; cv \leftarrow cv + v[i]$ 
    7.  $m \leftarrow \max\{m, cv\}$ 
    8.  $i \leftarrow i+1$ 
    9.    $\text{KNAPSACK}()$ 
    10.  $i \leftarrow i-1$ 
    11.  $c \leftarrow c + w[i]; cv \leftarrow cv - v[i]$ 
    12.  $i \leftarrow i+1$ 
    13.  $\text{KNAPSACK}()$ 
    14.  $i \leftarrow i-1$ 
  15. return
```

2.5 分支限界法 (Branch and Bound)

与回溯法不同，分支限界法通常采用广度优先 (Breadth-First) 或最佳优先 (Best-First) 的方式搜索解空间树。在 0-1 背包问题中，可用如下启发式函数 (Heuristic) 定义节点的优先级：

$$b = \sum_{j=1}^{i-1} v_j x_j + m_{\text{greedy}}(i, c) \quad (1)$$

其中， c 为当前背包剩余容量， $m_{\text{greedy}}(i, c)$ 为将子问题放宽为小数背包问题后的最优值，可由贪心算法得到。显然， $m_{\text{greedy}}(i, c)$ 是子问题最优值的一个上界，因此 (1) 式实际上给出了一个限界函数。

我们利用启发式函数 b 同时定义了分支策略与限界策略，得到优先队列式分支限界算法。伪代码如下

Algorithm 5: Branch and Bound

```
Input:  $n, c, w[1..n], v[1..n]$ .
Output:  $m$ .
1. Let  $k_i$  be the sorted index of  $\frac{v[i]}{w[i]}$  such that  $\frac{v[k_1]}{w[k_1]} \leq \frac{v[k_2]}{w[k_2]} \leq \dots \leq \frac{v[k_n]}{w[k_n]}$ .
2. Let  $q$  be a priority queue (max-heap) of node states in the form of  $(b, i, c, cv)$ .
3.  $m \leftarrow 0$ ;
4.  $b \leftarrow \text{BOUND}(1, c, 0)$ 
5. ENQUEUE( $q, (b, i, c, 0)$ )
6. loop
7.    $(b, i, c, cv) \leftarrow \text{DEQUEUE}(q)$ 
8.   if  $i > n$  or  $c = 0$  then
9.     break
10.    if  $b \leq m$  then
11.      continue
12.    ENQUEUE( $q, (\text{BOUND}(i+1, c, cv), i+1, c, cv)$ )
13.    if  $w[i] \leq c$  then
14.       $m \leftarrow \max\{m, cv + v[i]\}$ 
15.      ENQUEUE( $q, (b, i+1, c - w[i], cv + v[i])$ )
```

function $\text{BOUND}(i, c, cv)$:

```
1. for  $j \leftarrow i$  to  $n$  do
2.   if  $w[k_j] \leq c$  then
3.      $c \leftarrow c - w[k_j]$ 
4.      $cv \leftarrow cv + v[k_j]$ 
5.   else
6.     break
7. if  $j \leq n$  then
8.    $cv \leftarrow cv + \frac{v[k_j]}{w[k_j]} \cdot c$ 
9. return  $cv$ 
```

受限于搜索顺序，我们需要同时维护优先队列中所有节点的状态，因此空间开销远大于回溯法。上述算法搜索到一个可行解后即终止，由于启发式函数的良好性质，该解必定最优。

2.6 蒙特卡洛算法 (Monte Carlo Algorithm)

蒙特卡洛算法泛指一类随机算法，其输出解以一定的概率接近于问题的真实解，且随着算法重复次数的增加，输出解的正确概率逐渐提高。通常，蒙特卡洛算法依赖于重复的随机抽样来近似问题的解。

为了解决使用蒙特卡洛算法解决 0-1 背包问题，我们首先定义解空间上的概率分布 p ，使得最优解为该分布的众数，从而通过随机抽样得到的可能性最大。一个常见的选择是玻尔兹曼 (Boltzmann) 分布：

$$p(x) \propto \exp\left(-\frac{E(x)}{T}\right) \quad (2)$$

其中， $E(x) = -\sum_{i=1}^n v_i x_i$ 为状态 x 的势能 (Potential Energy)，由该状态下背包总价值取负得到 (若 x 不可行，则定义 $E(x) = +\infty$)； $T > 0$ 为温度 (Temperature)，可以预先选定或随着算法的迭代逐渐降低。由 (2) 式可知，当温度 x 越接近最优解，则背包价值越大，能量越小，该状态被抽样的概率 $p(x)$ 越大，符合我们对 $p(x)$ 的要求。

此外，观察发现，温度 T 越低，能量差对应的概率差越大，算法越倾向于抽样能量较小的状态；相反，当温度 T 趋近于正无穷时，算法对所有解中均匀抽样。因此，超参数 T 控制了算法对探索与利用 (Exploration-Exploitation) 的权衡，对求解速度与解的质量有重要影响。

由于无法得知从概率 $p(x)$ 到 x 的逆映射，直接按 (2) 式中独立同分布地抽样是困难且低效的。马尔可夫链蒙特卡罗 (Markov Chain Monte Carlo, MCMC) 算法提供了一种解决方案：从某个初始状态 $x^{(0)}$ 出发，每次抽样结果 $x^{(t)}$ 均依赖于上一次抽样结果 $x^{(t-1)}$ ，从而构造一个马尔可夫链 $\{x^{(t)}\}_{t=0}^{\infty}$ 。通过设计抽样机制，使得该马尔可夫链的稳态分布为 $p(x)$ ，最终以状态出现的频率近似概率 $p(x)$ 。

Metropolis-Hastings 是一种常用的 MCMC 算法，其抽样机制如下：

1. 从某一提议 (Proposal) 分布 $q(x'|x^{(t-1)})$ 中抽样得到 x' ；
 2. 计算接受率 (Acceptance Ratio) $\alpha(x', x^{(t-1)}) = \min\left\{1, \frac{p(x')q(x^{(t-1)}|x')}{p(x^{(t-1)})q(x'|x^{(t-1)})}\right\}$ ；
 3. 以概率 $\alpha(x', x^{(t-1)})$ ，令 $x^{(t)} = x'$ ；否则拒绝 x' ，保留 $x^{(t-1)}$ 。
- 由以上步骤，有状态转移概率 $p(x'|x^{(t-1)}) = q(x'|x^{(t-1)})\alpha(x', x^{(t-1)})$ 。容易证明，该抽样机制满足细致平衡 (Detailed Balance) 条件 $p(x)p(x'|x) = p(x')p(x|x')$ ，因此马尔可夫链的稳态分布为 $p(x)$ 。

在 0-1 背包问题中，我们选取提议分布 $q(x'|x)$ 为 x 的邻域中的均匀分布，即以等概率将状态 x 中的任一位置取反 (对应放入背包或拿出背包)，因而 $q(x'|x) = q(x|x') = \frac{1}{n}$ 。此时，接受率 $\alpha(x', x) = \min\left\{1, \frac{p(x')}{p(x)}\right\} = \min\left\{1, \exp\left(\frac{E(x)-E(x')}{T}\right)\right\}$ 。可见，当 x' 为不可行解时，能量为正无穷，算法必然拒绝 x' ，从而保证解的可行性；当 x' 为可行解，且价值较 x 更大时，能量差 $E(x') - E(x) < 0$ ，算法必然接受 x' ，从而实现价值函数的优化；当 x' 为可行解，且价值较 x 更小时，算法以一定概率接受 x' ，从而避免对解空间的探索。

为了确保在有限时间预算内得到较优解，兼顾解的质量与求解速度，我们按照一定的退火方案 (Annealing Schedule) 逐渐减小温度 T ，从而使求解法在前期更多地探索解空间，在后期更多地利用已有信息。考虑到实现的简单性，此处使用线性退火方案 $T = T_{\text{max}} - \frac{\text{elapsed}}{\text{budget}}(T_{\text{max}} - T_{\text{min}})$ ，其中 T_{max} 与 T_{min} 分别为初始温度与最终温度， elapsed 为算法已用时间， budget 为算法总时间预算。

综上，我们得到了模拟退火 (Simulated Annealing) 算法的完整流程，伪代码如下

Algorithm 6: Monte Carlo (Simulated Annealing)

```
Input:  $n, c, w[1..n], v[1..n], T_{\text{max}}, T_{\text{min}}, \text{annealing schedule}$ .
Output:  $m$ .
1.  $T \leftarrow T_{\text{max}}$ ;
2. Initialize  $x[1..n]$  as zero vector.
3. while  $T > T_{\text{min}}$  do
4.   Uniformly sample  $i \in \{1, 2, \dots, n\}$ .
5.   if  $x[i] = 0$  and  $\sum_{j=1}^n w[j]x[j] + x[i] > c$  then
6.     continue
7.    $\Delta E \leftarrow (2x[i] - 1)v[i]$ 
8.   if  $\Delta E < 0$  then
9.      $x[i] \leftarrow 1 - x[i]$ 
10.     $m \leftarrow \max\{m, \sum_{i=1}^n v[i]x[i]\}$ 
11.   else
12.     Uniformly sample  $u \in [0, 1]$ .
13.     if  $u < \exp(-\Delta E/T)$  then
14.        $x[i] \leftarrow 1 - x[i]$ 
15.   Update  $T$  according to annealing schedule.
```

2.7 深度强化学习 (Deep Reinforcement Learning)

近年来，深度强化学习在搜索 NP 难的组合优化问题的近似最优解方面取得了巨大的成功。现有工作可大致分为两类：(1) 利用深度强化学习模型直接构造问题的最优解；(2) 利用深度强化学习模型改进传统的局部搜索算法。此处，我们实现一种基于 REINFORCE 算法的模型，通过学习得到的策略，直接构造物品序列。该部分主要参考如下文献：

1. Kool, W., Hoof, H. van, & Welling, M. (2019). Attention, Learn to Solve Routing Problems! *ICLR 2019*.
 2. Kwon, Y.-D., Choo, J., Kim, B., Yoon, I., Gwon, Y., & Min, S. (2020). POMO: Policy Optimization with Multiple Optima for Reinforcement Learning. *NeurIPS 2020*.
- 类似于贪心算法的排序过程，我们选取 $n \in [5, 200]$ 进行多组实验。考虑到背包容量 c 和物品数量 n 的取值可能会影响问题的求解难度，我们固定 $c = 25n$ 。对于每组 n 和 c ，我们将背包剩余容量 $s^{(0)} = c - \sum_{j=1}^{t-1} w_{x_j}$ 用于问题的求解过程类似于选择排序，假定第 i 步时已选择了物品 $\{k_1, k_2, \dots, k_{i-1}\}$ ，则行动 $a^{(0)}$ 的目标是从 $\{1, 2, \dots, n\} \setminus \{k_1, k_2, \dots, k_{i-1}\}$ 中选择下一个物品 k_i 。可见，该问题的最优策略依赖于过去所有行动，即

$$\pi_{\theta}^{(i)} = \begin{cases} p_{\theta}(a^{(1)}|s^{(1)}) & \text{if } i = 1 \\ p_{\theta}(a^{(i)}|s^{(i)}, a^{(1:i-1)}) & \text{if } i = 2, 3, \dots, n \end{cases}$$

其中， $\pi_{\theta}^{(i)}$, $i = 1, 2, \dots, n$ 为策略函数， θ 为可学习参数。最终的解为所有行动的轨迹 (Trajectory) $\tau = (a^{(1)}, a^{(2)}, \dots, a^{(n)})$ ，其概率为

$$p_{\theta}(\tau|s^{(1)}) = p_{\theta}(a^{(1)}|s^{(1)}) \prod_{i=2}^n p_{\theta}(a^{(i)}|s^{(i)}, a^{(1:i-1)})$$

我们的目标是最大化期望累积奖励 $J(\theta, s) = \mathbb{E}_{p_{\theta}(\tau|s)}[R(\tau)]$ 。REINFORCE 是一种基于策略梯度的算法，该算法将 $J(\theta, s)$ 的梯度定义为

$$\nabla_{\theta} J(\theta, s) = \mathbb{E}_{p_{\theta}(\tau|s)}[(R(\tau) - b(s))\nabla_{\theta} \log p_{\theta}(\tau|s)]$$

其中， $b(s)$ 为基线函数 (Baseline)， $A(\tau, s) = R(\tau) - b(s)$ 称作优势函数 (Advantage)。

在 Policy Optimization with Multiple Optima (POMO) 中，Kwon et al. (2020) 中，模型每次从 N 个不同初始行动 $a^{(1)}, a^{(2)}, \dots, a^{(N)}$ 出发，同时推断出 N 个轨迹 $(\tau_1, \tau_2, \dots, \tau_N)$ ，从而将基线函数与策略梯度分别近似为

$$b(s) \approx \frac{1}{N} \sum_{i=1}^N R(\tau_i) \\ \nabla_{\theta} J(\theta, s) \approx \frac{1}{N} \sum_{i=1}^N (R(\tau_i) - b(s))\nabla_{\theta} \log p_{\theta}(\tau_i|s)$$

遵循 Kool et al. 2019 的思路，受 Transformer 启发，我们将策略网络设计为如下两部分：(1) Encoder，利用多头注意力机制将物品 i 的特征向量 (w_i, v_i) 嵌入到隐空间；(2) Decoder，将物品嵌入作为键和值，上下文嵌入 (此处，直接使用背包剩余容量 c) 作为查询，利用注意力机制计算单步行动的概率分布。具体的模型架构详见附件内的 TensorFlow 实现。

伪代码如下

Algorithm 7: Deep Reinforcement Learning (POMO)

```
function TRAINING(training set  $D$ , number of initial states  $N$ , number of epochs  $T$ , batch size  $B$ , learning rate  $\alpha$ ):
  1. Initialize policy network parameters  $\theta$ .
  2. for  $i \leftarrow 1$  to  $T$  do
  3.   Sample a batch of  $B$  problems  $\{P_1, P_2, \dots, P_B\}$  from  $D$ .
  4.   Encode item data into embeddings for each problem in the batch.
  5.   Select  $N$  initial actions  $\{a_{i,1}^{(0)}, a_{i,2}^{(0)}, \dots, a_{i,N}^{(0)}\}_{i=1}^B$  for each problem  $P_i$  in the batch.
  6.   Decode the policy and sample rollouts  $\{(\tau_{i,1}, \tau_{i,2}, \dots, \tau_{i,N})\}_{i=1}^B$  from  $N$  initial actions.
  7.   Compute baselines  $b_i \leftarrow \frac{1}{N} \sum_{j=1}^N R(\tau_{i,j})$ ,  $i = 1, 2, \dots, B$ .
  8.   Compute gradients  $\nabla_{\theta} J(\theta, s) \leftarrow \frac{1}{BN} \sum_{i=1}^B \sum_{j=1}^N (R(\tau_{i,j}) - b_i)\nabla_{\theta} \log p_{\theta}(\tau_{i,j}|s)$ .
  9.   Update policy network parameters  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta, s)$ .
  10. return  $\theta$ 
```

function INFERENCE($n, c, w[1..n], v[1..n]$):

```
1. Encode item data  $(w[1..n], v[1..n])$  into embeddings.
2. Select  $N$  initial actions  $\{a_{i,1}^{(0)}, a_{i,2}^{(0)}, \dots, a_{i,N}^{(0)}\}_{i=1}^B$ .
3. Decode the policy and select the best actions to form  $N$  solutions  $(\tau_1, \tau_2, \dots, \tau_N)$ .
4. return  $\tau^* = \arg \max_{\tau_j} R(\tau_j)$ .
```

3 经典算法性能比较

我们使用 C++ 分别实现 0-1 背包问题的分治法、动态规划算法、贪心算法、回溯法、分支限界法和蒙特卡洛算法。实验结果表明，对于不同的问题规模和应用需求，上述算法各有优劣。我们还实现了一个基于注意力机制的深度强化学习模型，并将其应用于 0-1 背包问题。该模型能够在较短时间求得较好的近似解，并且具有较好的泛化能力。该模型的性能超越了贪心策略，并且满足时间预算的要求。

4 深度强化学习算法性能评估

我们使用 TensorFlow 完整实现 Policy Optimization with Multiple Optima 算法。网络超参数如下：物品嵌入维数为 32，解码器层的键的维数为 16，注意力头数为 4，前馈隐藏层维数为 64，编码器层数为 2，解码器层对数几率的裁剪范围为 $[-10, 10]$ 。我们在 GTX 1080 Ti GPU 上训练 200 个 epoch。其中，每个 epoch 包含 5 个 batch，每个 batch 包含 64 个问题实例，每个问题均从 $n = 50$ 个不同初始物品分别出发得到 50 个行动轨迹，并据此计算梯度更新参数。为了提升模型的泛化能力，我们使用随机浮点数值作为 c, w, v ，并适当限制了 w 的下限，放宽了 w 的上限。我们使用 Adam 优化器，采用固定的学习率 8×10^{-5} 以及固定的权重衰减 1×10^{-7} 。训练过程花费约 7 分钟，损失函数变化如图 2 所示。注意下图中，训练集上的损失函数由随机抽样的行动轨迹计算得到，而验证集上的损失函数由概率最大的行动轨迹计算得到。

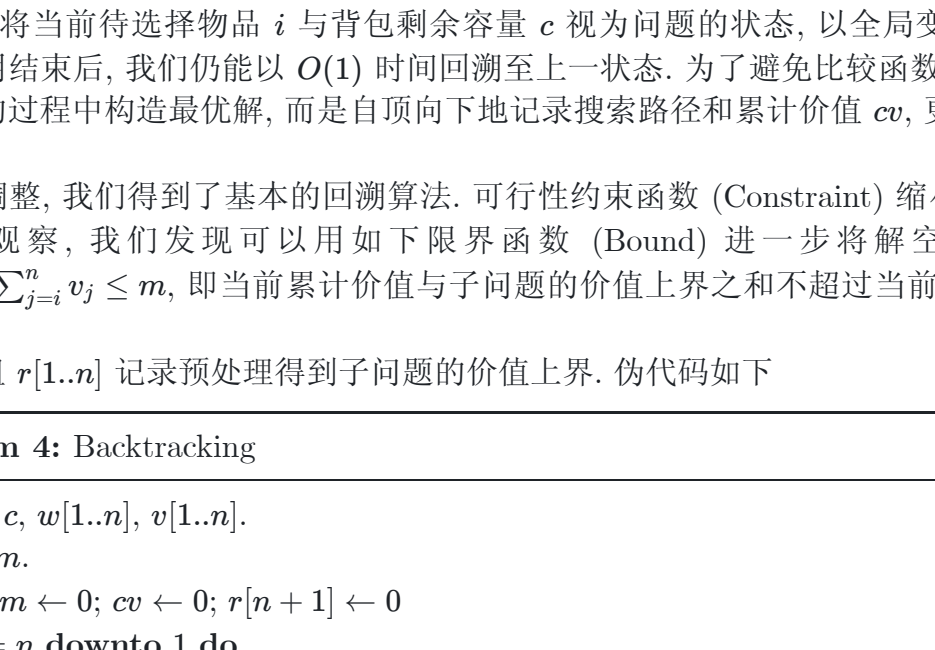


Figure 2. Learning curve of the REINFORCE agent.

我们将模型的性能与按单位重量价值排序的贪心策略以及模拟退火算法进行比较。每次测试，固定 $n = 50$, $c = 750$ ，每个方法均使用相同输入，模拟退火的时间预算限制为 200 毫秒。应当声明的是，我们的模型原生地支持如下特性：从 N 个不同初始物品输入上得到 N 个行动轨迹，选取最优的解作为最终结果。然而，该行为可能会导致性能评估的不公平性。因此，我们作出限制，只允许模型推断 1 个行动轨迹。实验结果如表 3 所示。

Table 3. Performance comparison between deep reinforcement learning method and classical baselines.

Method	Time (ms)	Optimality Gap
Greedy Strategy	0.1011	0.003
Simulated Annealing	201.3374	0.0252
Policy Optimization w/ Multiple Optima	166.8119	0.0022

可见，我们的深度强化学习模型超越了贪心策略的性能，并拥有较短的推理时间，展示出其优越性。由于该模型架构内在的特征，我们可以轻松地将其应用于任意问题规模 n 和 c ，并调整推理的超参数 N 达到最优性和效率的平衡。由于时间有限，此处不作多次测试。

5 结论

本实验中，我们实现了求解 0-1 背包问题的分治法、动态规划算法、贪心算法、回溯法、分支限界法和蒙特卡洛算法。实验结果表明，对于不同的问题规模和应用需求，上述算法各有优劣。我们还实现了一个基于注意力机制的深度强化学习模型，并将其应用于 0-1 背包问题。该模型能够在较短时间求得较好的近似解，并且具有较好的泛化能力。该模型的性能超越了贪心策略，并且满足时间预算的要求。