

Brief Description of the DP Algorithm Design

Objective

To estimate the maximum expected profit in T days given the current storage condition.

Mathematical Model

The entire storage can be divide into m categories and the items belonging to the i th category are expiring in i days where $i \in [0, m - 1]$. Note the maximum number of items within each category is limited to k . We store the number of items in all the categories in a m dimensional vector where the first element represents the number of items to be expired today and the last element represents the number of items to be expired in $m - 1$ days. We denote the vector as X and X_t is the storage status of the t -th day.

Everyday in the evening, the shop manager has to clean up the items that has been expired on that day and each item has a disposal cost θ . Next he shall decide the depletion amount. For each item depleted, an amount of salvage denoted as s will be recovered and counted into benefits. And after the depletion, all the items remained in the storage will be charged a holding cost each which is denoted as h . And in the morning of the following day, the manager will decide the total amount of items to order, the order cost is c per item. During the opening hours, those items which are going to be expired will be sold based on the possion distribution, and the selling price is r . Therefore the net income for the entire period can be expressed as

$$\begin{aligned} Profit = & + \text{numDeplete} \cdot s \\ & - \text{totalNum} \cdot h \\ & - \text{numOrder} \cdot c \\ & + \min(\mathbf{D}, \text{totalNum} + \text{numOrder}) \cdot r \\ & - X_0 \cdot \theta \end{aligned}$$

where

- numDeplete is the number of items been depleted on the night before
- totalNum is the number of items remained in the storage after depletion but before the next ordering
- D is the random variable to denote the number of items sold in a day. And it follows the Possion distribution.

So in order to find out the maximum expected value, we have the following recursion relation:

$$V_t(X_t) = \max_{z \in Z, q \in Q} (\mathbf{E}[Profit(X_t, z, q, \mathbf{D})] + \alpha \mathbf{E}[V_{t+1}(S(X_t, z, q), z, q, \mathbf{D})])$$

where

- z represents the number of depletion with range Z , and q is the number of ordering with the limit Q
- S is the function to evaluate X_{t+1} given X_t, z, q

Therefore V_t is attainable once given V_{t+1} .

We shall start from the table V_T , in which all the entries equal to zero, to recursively update the table V_i ($i = T - 1, T - 2, \dots$) and go all the way down to $i = 0$, which is the table of today.

Approaches

The problem itself has some fine features leading to a simpler solution.

- **Inherited Relation of Depletion**

All the state values within the table shall be organized in a tree, where the optimal depletion number of a child node has some relation with its parent node.

For example, for a problem with $m = 3$, $k = 2$, the following structure can be obtained.

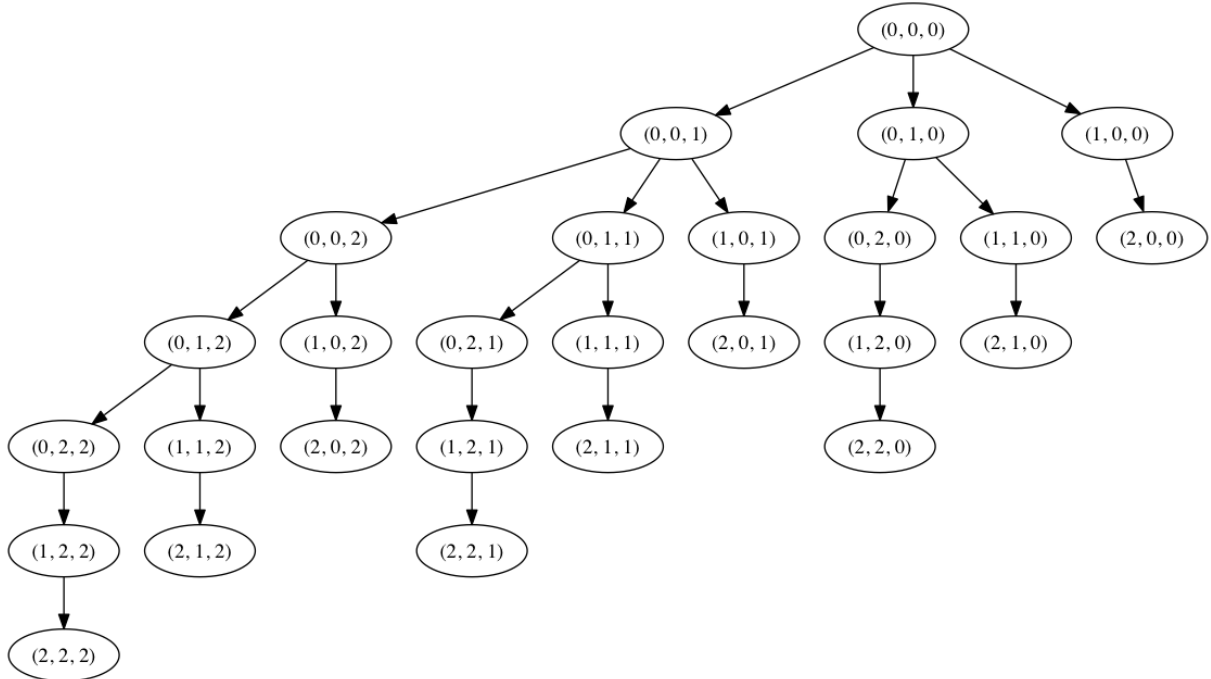


Figure 1: The tree structure for $m = 3$, $k = 2$

- **Limited Demands per day**

There is a upper limit for the number of customers that may visit the shop everyday. And maximum number should be rather small (hopefully less than a hundred). So instead of launching the kernels for thousands of times to get better estimation, we may simply calculate the result for every demand within the range and adjust the result by the weights given by the Possion distribution.

Based on the above analysis of the problem, a parallel process to deal with the problem can be described as follows. Take the problem discussed above as the example.

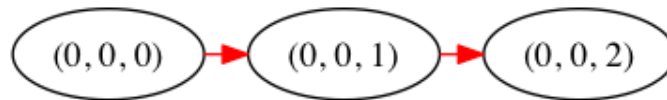


Figure 2: step 1

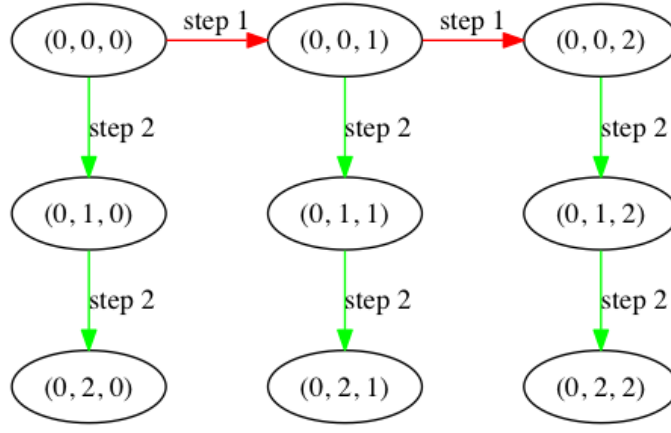


Figure 3: step 2

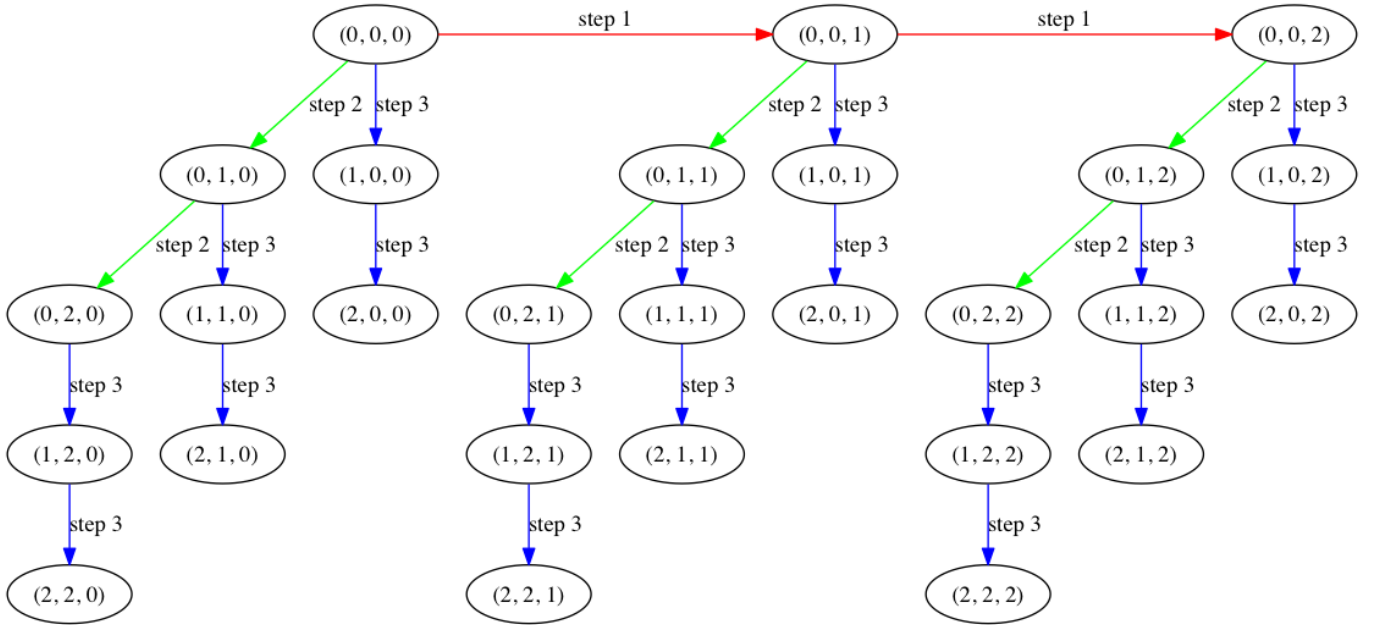


Figure 4: step 3

So in total m iterations are needed and within each iteration k^{i-1} threads operate in parallel where i is the index of the iteration.

Algorithm Implementation

There are three tables kept in the GPU memory all the time, the V_{t+1} , the V_t , and a table holding the all the depletion amount for each state of V_{t+1} . And what we are going to do here is to update the table of V_t based on V_{t+1} and the depletion information.

As has been discussed above, the value table will be updated in m iterations (by updating V_t and V_{t+1} alternatively) and within each iteration, k^{i-1} threads will be used to update a total of k^i entries in the V_{t+1} table. And for each thread, it takes care of k nodes in the V_{t+1} table which are considered as the “descents” of the corresponding entry in the V_t table as has been discussed above.

So what we are going to do here is to implement an algorithm which contains two for loops, the outer one is deployed on the CPU and it takes care of the m iterations and the inner loop is deployed on the specific GPU threads and each takes care of a node family.

Program Description

Functions:

- *checkCudaErrors*: This function is to detect any errors while invoking CUDA functions. And this function should be called whenever a CUDA function is called.
- *cudaHostAlloc*: This function is to allocated some memory on host (the CPU memory and map this memory to the CUDA memory). So if some memory is allocated in this way, those memory is visible to the CUDA SMs.
- *cudaSetDeviceFlags(cudaDeviceMapHost)*: To enable the function mentioned above.
- *cudaThreadSynchronize*: This function is for synchronize the threads within a kernel call. And the function should be called whenever the kernel tries to do some sequential modification on the global memory.
- *cudaDeviceSynchronize*: This function is to synchronize the entire block and is rarely used.
- *threadIdx* & *blockIdx*: In a CUDA device, the computation units are divided into different blocks, each block is divided into 1024 threads. All the threads can operate simultaneously. Though there is this “wrap” thing, we don’t need to care too much about it on this problem. And the *blockIdx* is a system parameter which stores the index of the current block in the global scope. And note since the blocks are managed in 3D, the *blockIdx* is a structure contains three attributes, x , y , and z . Similar to the *blockIdx*, the *threadIdx* stores the thread index of the current thread within the block, and it’s also a structure with three attributes.
- *global* and *device*: the global keyword denotes the kernel function which is visible to the host function, and the device keyword denotes the kernel function which is only visible to another kernel function. And the keyword *inline* should be used while writing small device function for the best efficiency.

Some Left Overs

Currently the bottleneck of this algorithm is the data transmitting speed between the host and the device. And the problem size is limited by the GPU memory. We are still looking for solutions to break the limit and implement the SLI configuration.