

# Divide and Conquer

Robert Kirkby\*

July 16, 2024

## Abstract

Gordon & Qiu (2018) introduce a 'binary monotonicity' algorithm that exploits monotonicity to solve value function problems faster. This brief pdf describes the adaption of 'two-level monotonicity' for the GPU that is implemented in VFI Toolkit.

**Keywords:** Value Function Iteration, Divide and Conquer, VFI Toolkit.

---

\*Thanks to AAA. Thanks to Rāpoi at Victoria University of Wellington for the use of their computing facilities. Kirkby: Victoria University of Wellington. Please address all correspondence about this article to Robert Kirkby at <robertdkirkby@gmail.com>, [robertdkirkby.com](http://robertdkirkby.com).

This document is going to largely assume you already read Gordon and Qiu (2018). The basic idea is that for almost all value function problems, we know that the solution will be 'monotone', in the sense that  $aprime(a)$  is a (weakly) increasing function of  $a$ ; where  $aprime$  is next period endogenous state and  $a$  is this period endogenous state. We can exploit this monotonicity to check less points (as many would fail to be monotone and so cannot be correct) and so can find the solution faster and using less memory.

To start, consider the naive pure discretization approach that is default in VFI Toolkit. Say we put 12 points on the endogenous state. Then we will check all 12-by-12 points, and find the solution (which is to say that for each grid point of  $a$  we find the point in  $aprime$  that maximizes the value function problem). Figure 1 illustrates this idea, the pale blue indicates the points we have evaluated, and the dark blue indicates the solution. So pure discretization is going to have to try all 144 (12 times 12) points. Because we can use the GPU this is workable as they can all be evaluated in parallel, but it would be very slow on a CPU.

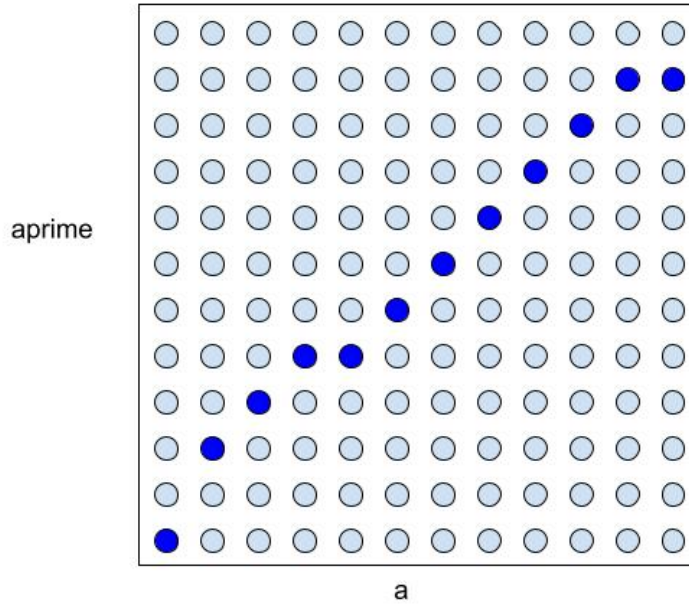


Figure 1: Pure Discretization

'Simple monotonicity' takes advantage of that fact that we know that  $aprime$  is weakly increasing in  $a$ . We start with the first grid point of  $a$ , and look at all the possible  $aprime$  to find the solution for that first  $a$  grid point. Since we know that the second grid point of  $a$  must have an  $aprime$  solution for a larger value of  $aprime$  we only need to check 'larger'  $aprime$  values. Then for the third grid point of  $a$ , we only need to consider values of  $aprime$  larger than the  $aprime$  that was the solution for the second grid point. Figure 2 illustrates this idea, the pale blue indicates the points we have evaluated, and the dark blue indicates the solution. Notice that 'simply monotonicity' has

substantially reduced the number of points we have to evaluate.

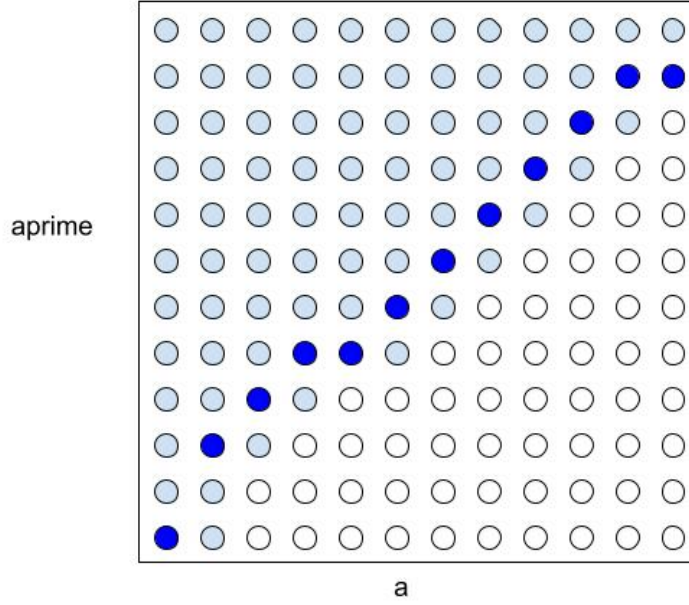


Figure 2: Simple Monotonicity

'Binary monotonicity' introduced by Gordon and Qiu (2018) improves on this further. We first look at the smallest and largest  $a$  grid points, check all  $a_{prime}$ , and find their solutions. We then turn to the middle grid point of  $a$ , and only have to check above (or equal to) the  $a_{prime}$  from the first  $a$  grid point, and below the  $a_{prime}$  of the last  $a$  grid point. Next, the lower quarter grid point of  $a$ , where we only have to check above  $a_{prime}$  of the smallest  $a$  grid point, and below  $a_{prime}$  of the middle  $a$  grid point. Similarly, the upper quarter grid point of  $a$ , where we only have to check above  $a_{prime}$  of the middle grid point, and below  $a_{prime}$  of the largest  $a$  grid point. Keep subdividing each interval in 2 (this description is a bit rough, see Gordon and Qiu (2018) for precise formulation; hopefully I have communicated the rough idea). Figure 3 illustrates this idea, the pale blue indicates the points we have evaluated, and the dark blue indicates the solution. You can see just how few points binary monotonicity has to evaluate compared to the alternatives.

On the CPU, this is about the end of the story. Fewer points is better. But on the GPU there is more to it. Yes, less points to evaluate will be faster on the GPU than more points to evaluate. But serial on the GPU will be much slower to evaluate than parallel. And binary monotonicity is *very* serial.

So now we come to the idea of 'two-level monotonicity', which is what this note introduces and is implemented in VFI Toolkit. We choose ' $n$ ' points for the first level, let's make it 3 for our example. In the first step we solve the 'first level': we choose  $n = 3$  equally spaced points on the  $a$

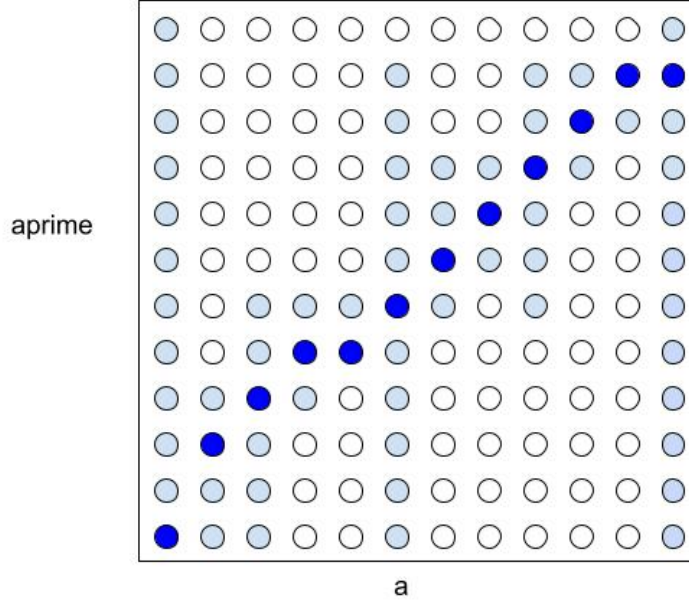


Figure 3: Binary Monotonicity

grid (these  $n$  will include both endpoints)<sup>1</sup> and for each of these we evaluate all  $a_{prime}$  to find the solution. For the second level we use a for-loop counting  $ii$  over  $n - 1$  iterations, and solve all  $a$  grid points between the  $ii$ -th and  $ii + 1$ -th  $a$  grid points from level one, and allow a range for  $a_{prime}$  from the solutions corresponding to these two points from the first level. Figure 4 illustrates this idea, the pale blue indicates the points we have evaluated, and the dark blue indicates the solution.

Obviously, 'two-level monotonicity' involves evaluating more points than 'binary monotocity'. So on a CPU 'binary monotocity' will be faster. But on a GPU, 'two-level monotonicity' provides a reasonable compromise between evaluating less points, and being more parallel/less serial.

My rough tests suggest that for GPU 'two-level monotonicity' is a good balance. In practice, I often seem to see little gain from using 11 or 15 points in the first level, rather than just 5. This suggests that the costs of serialization are outweighing the benefits of evaluating less points fairly early on.

## Decision Variables and Exogenous Shocks

The obvious approach is to apply this idea 'conditional on' any decision variables and exogenous shocks. But on GPU this would involve lots of loops as the number of  $a_{prime}$  points to be evaluated would be (potentially and likely often) different in each case.

<sup>1</sup>In Matlab code we choose the indexes  $level1ii$  as  $level1ii = round(linspace(1, n_a, vfoptions.level1n));$ , where  $n_a$  is number of grid points for  $a$ , and  $vfoptions.level1n$  is the number  $n$  of grid points to be used in the first level.

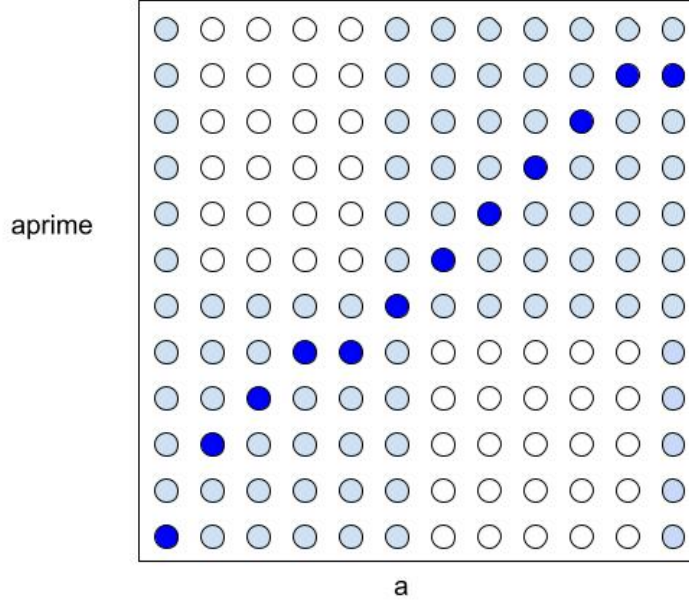


Figure 4: Two-level Monotonicity

My first idea just use the  $\min()$  and  $\max()$  across all decision variables and exogenous shocks of the relevant lower and upper  $a_{prime}$  range for each given  $a$  grid point in the first level. This will involve evaluating more points, but helps keep more parallel and less serial. A very quick and dirty test I did suggested this was better than doing loops using different ranges for each value of the decision variables and exogenous shocks.

This was later improved to a method that (most of the time) required evaluating less grid points. Think of an the range conditional on  $d$  for which we have to check  $a_{prime}$  as  $alower(d)$  to  $aupper(d)$ . The idea was to use the acutal lower  $a_{prime}$  for each  $d$  as the lower bounds, but then have the range of points being  $max_d[aupper(d) - alower(d)]$ . This reduces the range of  $a_{prime}$  points to evalute from  $max_d[aupper(d)] - min_d[alower(d)]$  in the first idea to  $max_d[aupper(d) - alower(d)]$  in the current one. Note that while  $alower(d)$  is in principle the lower bound of this range, I overwrite this with  $n_a - range$  whenever it would otherwise lead to  $alower(d)$  plus the range going off the top of the grid on  $a$ .

This has the added advantage of loosening further the assumptions under which this algorithm will give the correct solution. I think of it as just requiring  $a_{prime}(a)$  is weakly increasing, conditional on fixing any decision variables and exogenous shocks. In fact the required assumption is slightly weaker because of taking the max range over decision variables and exogenous shocks (of first level results, to decide ranges of  $a_{prime}$  considered in second level iterations); although I am not sure if there are any economic applications in which the first conditions fail but the second hold, so you can probably just stick to thinking about the first conditions. Gordon and Qiu (2018)

have a discussion in their paper about when monotonicity holds and how you might prove it for a model.

Because it is very easy in VFI Toolkit to switch between the default pure discretization and the two-level monotocity solutions, the easiest thing in practice is probably just to run both and double-check they give the same answer (this is how I tested them :).

## References

- Grey Gordon and Shi Qiu. A divide and conquer algorithm for exploiting policy function monotonicity. Quantitative Economics, 9(2), 2018. doi: <https://doi.org/10.3982/QE640>.
- Robert Kirkby. VFI toolkit, v2. Zenodo, 2022. doi: <https://doi.org/10.5281/zenodo.8136790>.