

Speedup of Uniform Bicubic Spline Interpolation

Viliam Kačala, Lukáš Miňo, and Csaba Török

Institute of Computer Science, Faculty of Science, Pavol Jozef Šafárik University in Košice
 Jesenná 5, 040 01 Košice, Slovakia
 viliam.kacala@student.upjs.sk, lukas.mino@upjs.sk, csaba.torok@upjs.sk

Abstract: The goal of the paper is to introduce an efficient algorithm for computation of derivatives of bicubic spline surfaces over equispaced grids with C^2 class continuity. The algorithm is based on a recently proposed approach using a special approximation property between quartic and cubic polynomials. The proposed solution replaces the classical de Boor’s systems of equations with systems of reduced size and simple remainder explicit formulas. We will show that the proposed new algorithm is numerically equivalent to de Boor’s algorithm and the former is more than 50% faster.

1 Introduction and problem statement

Modeling of surfaces plays a key role in a wide variety of computer science fields such as graphics or CAD applications. The paper proposes an efficient computational algorithm for interpolating uniform bicubic spline surfaces of class C^2 . The standard way to accomplish that is an almost 60 years old classic algorithm designed by Carl de Boor [2] that uses tridiagonal systems of equations. We will refer to it as the *full* algorithm.

Even though the evaluation of such systems is quite straightforward in general, running in linear time, it can be still viewed as slow especially in real time computing scenarios. For this reason we present a new *reduced* algorithm for uniform bicubic spline surfaces that, while being in the same complexity class, is faster and needs lower memory requirements.

The structure of the paper is as follows. The remaining part of this section is devoted to the problem statement and introducing terminology. Section 2 presents the new reduced algorithm along with proof of its equality to the full algorithm. The third section provides experimental measurements of actual time savings of the new approach along with some words about its efficient implementation. To be self contained, we provide the equations of de Boor’s full algorithm for surfaces and the reduced algorithm for curves in Appendix, for easier comparison with ours.

Now let’s jump into the problem statement. Consider a uniform grid

$$[x_0, x_1, x_2, \dots, x_{I-1}] \times [y_0, y_1, y_2, \dots, y_{J-1}], \quad (1)$$

where

$$x_i = x_0 + ih_x, \quad i = 1, 2, 3, \dots, I - 1, \quad I > 1, h_x \in \mathbb{R}^+, \\ y_j = y_0 + ih_y, \quad j = 1, 2, 3, \dots, J - 1, \quad J > 1, h_y \in \mathbb{R}^+.$$

According to [2], a spline surface is defined by given values

$$z_{i,j}, \quad i = 0, 1, 2, \dots, I - 1, \quad j = 0, 1, 2, \dots, J - 1 \quad (2)$$

at the grid-points, and given first directional derivatives

$$d_{i,j}^x, \quad i = 0, I - 1, \quad j = 0, 1, 2, \dots, J - 1 \quad (3)$$

at the boundary verticals,

$$d_{i,j}^y, \quad i = 0, 1, 2, \dots, I - 1, \quad j = 0, J - 1 \quad (4)$$

at the boundary horizontals and cross derivatives

$$d_{i,j}^{x,y}, \quad i = 0, I - 1, \quad j = 0, J - 1 \quad (5)$$

at the four corners of the grid.

The task is to define a quadruple $[z_{i,j}, d_{i,j}^x, d_{i,j}^y, d_{i,j}^{x,y}]$ at every grid-point $[x_i, y_j]$, based on which a bicubic clamped spline surface S of class C^2 is constructed with properties

$$S(x_i, y_j) = z_{i,j}, \quad \frac{\partial S(x_i, y_j)}{\partial y} = d_{i,j}^y, \\ \frac{\partial S(x_i, y_j)}{\partial x} = d_{i,j}^x, \quad \frac{\partial^2 S(x_i, y_j)}{\partial x \partial y} = d_{i,j}^{x,y}.$$

For $I = 7$ and $J = 5$, we provide visual illustration of the input situation in Figure 1, where bold values represent (2) – (5) while the remaining non-bold values represent the unknown derivatives to compute.

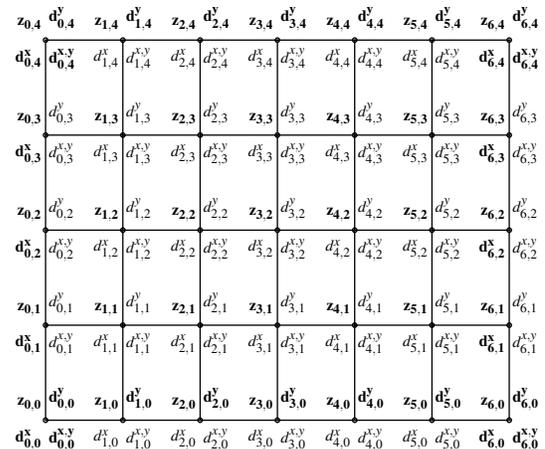


Figure 1: Input situation for a 7 × 5 uniform grid.

we get (6) from (10), (7) from (24) and the remainder formulas (9) from (26) that alongside with Lemma 3 finishes the proof. \square

2.2 Algorithm

Now we have sufficient theoretical foundations to propose an efficient algorithm to compute spline derivatives over a uniform grid of points. The reduced algorithm consists of four main steps, each evaluating systems of equations derived from (7) and remainder formulas (9). The algorithm as a whole takes the same input values and provides identical results as the original full algorithm in Appendix. Therefore it is intended as a drop-in replacement for the full algorithm.

Based on Lemma 1 the full systems (21) – (23) can be solved using analogical reduced systems and explicit formulas. Hence we get immediately the reduced algorithm for solving the unknown derivatives of a C^2 spline surface that appears to provide faster computations than de Boor's algorithm.

Theorem 1 (Reduced algorithm for surfaces). *If the values (2) – (5) over the uniform grid (1) are given, then the unknown values*

$$\begin{aligned} d_{i,j}^x, & \quad i = 1, 2, 3, \dots, I-2, \quad j = 0, 1, 2, \dots, J-1, \\ d_{i,j}^y, & \quad i = 0, 1, 2, \dots, I-1, \quad j = 1, 2, 3, \dots, J-2, \\ d_{i,j}^{x,y}, & \quad i = 1, 2, 3, \dots, I-2, \quad j = 0, J-1, \\ d_{i,j}^{x,y}, & \quad i = 0, 1, 2, \dots, I-1, \quad j = 1, 2, 3, \dots, J-2 \end{aligned}$$

for spline surface of class C^2 are uniquely determined by the following linear systems and formulas in four main steps:

Step 1a. Computation of d^x along the horizontals from equation systems for inner even-indexed grid-points.

For each $j = 0, 1, 2, \dots, J-1$,

$$\begin{aligned} & \text{solve_system}(\quad \\ & \quad d_{i+2,j}^x - 14d_{i,j}^x + d_{i-2,j}^x = \\ & \quad = \frac{3}{h_x}(z_{i+2,j} - z_{i-2,j}) - \frac{12}{h_x}(z_{i+1,j} - z_{i-1,j}), \quad (12) \\ & \quad \text{where } i = 2, 4, 6, \dots, I-3 \\ & \quad). \end{aligned}$$

Step 1b. Computation of d^x along the horizontals from explicit formulas for inner odd-indexed grid-points.

For each $i = 1, 3, 5, \dots, I-2$ and $j = 0, 1, 2, \dots, J-1$,

$$d_{i,j}^x = \frac{3}{4h_x}(z_{i+1,j} - z_{i-1,j}) - \frac{1}{4}(d_{i+1,j}^x + d_{i-1,j}^x). \quad (13)$$

Step 2a. Computation of d^y along the verticals from equation systems for inner even-indexed grid-points.

For each $i = 0, 1, 2, \dots, I-1$,

$$\begin{aligned} & \text{solve_system}(\quad \\ & \quad d_{i,j+2}^y - 14d_{i,j}^y + d_{i,j-2}^y = \\ & \quad = \frac{3}{h_y}(z_{i,j+2} - z_{i,j-2}) - \frac{12}{h_y}(z_{i,j+1} - z_{i,j-1}), \quad (14) \\ & \quad \text{where } j = 2, 4, 6, \dots, J-3 \\ & \quad). \end{aligned}$$

Step 2b. Computation of d^y along the verticals from explicit formulas for inner odd-indexed grid-points.

For each $i = 0, 1, 2, \dots, I-1$ and $j = 1, 3, 5, \dots, J-2$,

$$d_{i,j}^y = \frac{3}{4h_y}(z_{i,j+1} - z_{i,j-1}) - \frac{1}{4}(d_{i,j+1}^y + d_{i,j-1}^y). \quad (15)$$

Step 3a. Computation of $d^{x,y}$ along the horizontals from equation systems for inner even-indexed grid-points.

For each $j = 0, J-1$,

$$\begin{aligned} & \text{solve_system}(\quad \\ & \quad d_{i+2,j}^{x,y} - 14d_{i,j}^{x,y} + d_{i-2,j}^{x,y} = \\ & \quad = \frac{3}{h_x}(d_{i+2,j}^x - d_{i-2,j}^x) - \frac{12}{h_x}(d_{i+1,j}^x - d_{i-1,j}^x), \quad (16) \\ & \quad \text{where } i = 2, 4, 6, \dots, I-3 \\ & \quad). \end{aligned}$$

Step 3b. Computation of $d^{x,y}$ along the horizontals from explicit formulas for inner odd-indexed grid-points.

For each $i = 1, 3, 5, \dots, I-2$ and $j = 0, J-1$,

$$d_{i,j}^{x,y} = \frac{3}{4h_x}(d_{i+1,j}^x - d_{i-1,j}^x) - \frac{1}{4}(d_{i+1,j}^{x,y} + d_{i-1,j}^{x,y}). \quad (17)$$

Step 4a. Computation of $d^{x,y}$ along the verticals from equation systems for inner even-indexed grid-points.

For each $i = 0, 1, 2, \dots, I-1$,

$$\begin{aligned} & \text{solve_system}(\quad \\ & \quad d_{i,j+2}^{x,y} - 14d_{i,j}^{x,y} + d_{i,j-2}^{x,y} = \\ & \quad = \frac{3}{h_y}(d_{i,j+2}^y - d_{i,j-2}^y) - \frac{12}{h_y}(d_{i,j+1}^y - d_{i,j-1}^y), \quad (18) \\ & \quad \text{where } j = 2, 4, 6, \dots, J-3 \\ & \quad). \end{aligned}$$

Step 4b. Computation of $d^{x,y}$ along the verticals from explicit formulas for inner odd-indexed grid-points.

for each $i = 0, 1, 2, \dots, I-1$ and $i = 1, 3, 5, \dots, I-2$,

$$d_{i,j}^{x,y} = \frac{3}{4h_y}(d_{i,j+1}^y - d_{i,j-1}^y) - \frac{1}{4}(d_{i,j+1}^{x,y} + d_{i,j-1}^{x,y}). \quad (19)$$

If I is even, then the last model equation in steps (12) and (16) needs to be accordingly replaced by the model equation derived according to (8). Analogically, if J is even, the same applies to steps (14) and (18).

We underline that the systems (12), (14), (16) and (18) are diagonally dominant and use the same tridiagonal matrix with constant main diagonal and upper, lower diagonals.

In the next section it will be shown that the reduced algorithm is faster than the full one and additionally it needs lower memory requirements.

The main difference between original full algorithm and the reduced one is that the latter computes only half of the unknown values using LU factorization where the remaining half of the unknowns are solved by simple remainder formulas. Because of this, the reduced algorithm is supposed to be more efficient in terms of execution speed and memory requirements as well. We will discuss and measure it's speed in Section 3.

3 Performance evaluation

To be worth of actual implementation, the reduced algorithm should be significantly faster than full one. The elementary task in both algorithms is computation of tridiagonal systems of equations. As the full and the reduced systems are diagonally dominant we use a modified LU factorization method as the internal solver for such systems [4]. Our optimizations of the LU factorization relies on a fact, that the left-hand side matrices comprise only constant values.

While the LU factorization for tridiagonal linear systems is a quite efficient way to solve the equations which are interdependent, it has some drawbacks in the sense of modern CPU's capabilities like usage of vectorized SIMD instructions especially due to the structure of tridiagonal matrices. By breaking down the computational task of the full approach into reduced systems and simple explicit formulas, half of the equations are independent and therefore more versatile to manual or automatic compiler optimizations.

Memory requirements For the sake of completeness some words about applied data structures and memory requirements should be given.

The input grid (1) of size $I \times J$ needs $I + J$ memory to store x and y coordinates of the total IJ grid-points. Each grid-point however requires $4IJ$ space for $z_{i,j}$, $d_{i,j}^x$, $d_{i,j}^y$, $d_{i,j}^{x,y}$ values, where most of them will be computed by the full or reduced algorithm. This gives us $4IJ + I + J$ space requirement to store input/output values.

The needs of the full and reduced algorithms are quite low regarding the size of the grid. The full tridiagonal systems in Lemma 2 require $2 \cdot \max(I, J)$ space to store the right-hand side vector and an auxiliary buffer vector used for the LU factorization. In case of the reduced algorithm, only half number of the equations form the tridiagonal system, therefore they require only $\max(I, J)$ space for the LU part.

Data structures Since the grid may contain tens of thousands or more grid-points, the most effective representation is the jagged array structure for each of the z , d^x , d^y , $d^{x,y}$ values. Each equation system from any of the two algorithms always depends on one row of a jagged array, therefore entire rows of the jagged structure can be effectively cached under the assumption, that the size of the row is not very large. Notice that the iterations for computing the $d_{i,j}^y$ and most of the $d_{i,j}^{x,y}$ values in both algorithms have interchanged indices compared to the iteration throughout the $d_{i,j}^x$ values. We mention that an efficient implementation needs to setup the jagged arrays in accordance with how we want to iterate the data [7].

3.1 Measurements

Now let us compare the implementations of both algorithms. We implemented a benchmark in C++ 17 compiled with MSVC 2017 using `-O2` optimization level and individual code generation for each tested CPU, i.e. CPUs with AVX2 support were running binaries compiled to AVX2 instruction set whereas older CPUs received for instance only SSE2 compiled binaries. Testing environments comprised several computers with various CPUs ranging from rather obsolete Nehalem to recent Skylake microarchitecture. All systems had 8 – 32 GB of RAM, SSD and Windows 10 installed. The tests were conducted on freshly booted PCs after 2 minutes of idle time without running any non-essential services or processes like web browsers, database engines, etc.

On all testing computers the tests were conducted on two datasets, a small one on a grid of size $I, J = 100$, and a large one, where grid dimensions were $I, J = 1000$. Both datasets comprised the grid $[x_0, x_1, \dots, x_{I-1}] \times [y_0, y_1, \dots, y_{J-1}]$, where $x_0 = -20$, $x_{I-1} = 20$, $y_0 = -20$, $y_{J-1} = 20$ and the values $z_{i,j}$, $d_{i,j}^x$, $d_{i,j}^y$, $d_{i,j}^{x,y}$, see (2) – (5), were computed from function $\sin \sqrt{x^2 + y^2}$ at each grid-point. These datasets were chosen arbitrarily and the behaviour of the algorithms was correct for any input values. The speedup values were gained averaging 50 measurements of each algorithm.

Tables 1 and 2 contain results for both datasets and consist of four columns. The first column contains the tested CPUs ordered by their release date. Columns two through four contain measured execution times in microseconds for both algorithms and their speed ratios.

To provide better comparison for more datasets, Table 3 provides measurements on more different input sizes, however for the sake of readability it contains measurements only from Intel i7 6700K as the fastest testing CPU.

Let us review the measured performance improvement of the reduced algorithm. Results of Table 1 say that the reduced algorithm is more than fifty percent faster than the full one as long as the grid dimensions remains relatively small to enable the entire rows of the grid as well as the auxiliary LU vectors to be cached.

CPU	$I, J = 100$		
	Full	Reduced	Speedup
Intel i5 M430	783	438	1.79
AMD A6 3650M	922	567	1.64
Intel i7 4790	482	250	1.93
Intel i7 6700K	355	190	1.86
AMD X4 845	648	347	1.86

Table 1: Comparison of full and reduced algorithms on small dataset. Times are in microseconds.

CPU	$I, J = 1000$		
	Full	Reduced	Speedup
Intel i5 M430	116 856	90 541	1.29
AMD A6 3650M	116 390	78 892	1.48
Intel i7 4790	54 719	33 584	1.63
Intel i7 6700K	37 029	22 571	1.64
AMD X4 845	78 860	48 955	1.61

Table 2: Comparison of full and reduced algorithms on large dataset. Times are in microseconds.

The speedup of the reduced algorithm in Table 2 with a larger dataset is not so high due to difficulties of keeping data in fast but small L1 cache, resulting in higher ratio of cache misses. In the case of even larger dataset, let's say in the order of billions of grid-points, the performances of both algorithms are similar with the reduced one gaining only small advantage.

In our experiments we observed that the reduced algorithm was always faster than the full one, however exact speedup depends on the size of the input grid. The general rule is: the larger the grid the smaller the speedup.

The highest speedup does not depend on the grid size if it is from range of 50 to 200.

CPU	Full	Reduced	Speedup
$I, J = 50$	92	50	1.84
$I, J = 100$	355	190	1.86
$I, J = 200$	1 417	778	1.82
$I, J = 300$	3 203	1 797	1.78
$I, J = 400$	5 791	3 234	1.79
$I, J = 1000$	37 029	22 571	1.64
$I, J = 1500$	88 236	54 141	1.63
$I, J = 2000$	178 144	115 674	1.54

Table 3: Multiple dataset comparison of full and reduced algorithms tested on i7 6700K. Times are in microseconds.

4 Discussion

Let us briefly discuss the new algorithm from the numerical and experimental point of view. While the classic full algorithm is composed of four series of tridiagonal systems of equations, the reduced algorithm breaks down each equation system into a reduced one approximately half the size of the original one and simple mutually independent explicit formulas.

In addition, from the numerical point of view, the reduced tridiagonal subsystems are diagonally dominant and therefore computationally stable [1], similarly to the full systems. The remainder explicit formulas (9) are simple and thus do not present an issue.

The maximal numerical difference in our C++ implementation was in the order of 10^{-16} on several different datasets so the reduced algorithm yields numerically accurate results.

Since the algorithm consists of many independent systems of linear equations, it can be also effectively parallelized for both CPU and GPU architectures.

There is also a future work for further reduction of the number of equations in the tridiagonal systems.

5 Conclusion

The paper introduced a new algorithm to compute the unknown derivatives for uniform bicubic spline surfaces of class C^2 . The algorithm reduces the size of the tridiagonal systems of equations by half and computes the remaining unknown derivatives using simple explicit formulas. A substantial decrease of execution time of derivatives at grid-points has been achieved with lower memory space requirements at the cost of a slightly more complex implementation where the measured speedup ranges from 1.3 to 1.9 depending on the grid size and CPU architecture.

Acknowledgements

This work was partially supported by projects Technicom ITMS 26220220182 and APVV-15-0091 Effective algorithms, automata and data structures.

Appendix

Carl de Boor in [2] proposed an algorithm for computation of the unknown derivatives of spline surface of class C^2 over any grid with four types of full systems of linear equations. The following lemma reformulates this algorithm for uniform grid.

Lemma 2 (Full algorithm for surfaces). *If the values (2)–(5) over the uniform grid (1) are given, then the unknown*

- [9] Szabó I.: Approximation Algorithms for 3D Data Analysis. PhD Thesis, P. J. Šafárik University in Košice, Slovakia (2016)
- [10] Török Cs.: On reduction of equations' number for cubic splines. *Matematicheskoe modelirovanie*, Vol. 26 , No. 11, ISSN 0234–0879, (2014), 33–36
- [11] Török, Cs.: Speed-up of Interpolating Spline Construction, (to appear)