

How to Solve a Cubic Equation

Part 4 –The 111 case

James F. Blinn
 Microsoft Research
 blinn@microsoft.com

Originally published in
IEEE Computer Graphics and Applications
 Jan/Feb 2007, pages ??

Our Story So Far

We have been looking for roots of the homogeneous cubic equation

$$f(x, w) = Ax^3 + 3Bx^2w + 3Cwx^2 + Dw^3 = 0$$

The first step is to calculate the coefficients of the Hessian quadratic and use them to find the discriminant, Δ

$$\begin{aligned}\delta_1 &= AC - B^2 \\ \delta_2 &= AD - BC \\ \delta_3 &= BD - C^2 \\ \Delta &= 4\delta_1\delta_3 - \delta_2^2\end{aligned}$$

We then perform a coordinate transform to “depress” the cubic, turning it into one that has an x^2w coefficient of zero. The simplest such transformation is a translation by the quantity B/A , but a more general transformation that does this can be expressed as

$$\begin{bmatrix} x & w \end{bmatrix} = \begin{bmatrix} \tilde{x} & \tilde{w} \end{bmatrix} \begin{bmatrix} t & u \\ s & v \end{bmatrix}$$

where

$$\begin{aligned}s &= -t^2B - 2tuC - u^2D = -\frac{1}{3}f_w(t, u) \\ v &= t^2A + 2tuB + u^2C = \frac{1}{3}f_x(t, u)\end{aligned}$$

Some algebra gives us three new coefficients $\tilde{A}, \tilde{C}, \tilde{D}$ that are polynomial functions of the (t, u) values that parameterize the transformation. So now we have

$$\tilde{A}(t, u)\tilde{x}^3 + 3\tilde{C}(t, u)\tilde{x}\tilde{w}^2 + \tilde{D}(t, u)\tilde{w}^3 = 0$$

We then found that the \tilde{C} and \tilde{D} polynomials contain the \tilde{A} polynomial as a factor. Dividing this out and, without loss of generality, setting $\tilde{w} = 1$ gives the simple polynomial to solve

$$\tilde{x}^3 + 3\bar{C}(t, u)\tilde{x} + \bar{D}(t, u) = 0$$

And we further found that these quantities always satisfy the identity

$$\bar{D}^2(t, u) + 4\bar{C}^3(t, u) = -\tilde{A}^2(t, u)\Delta$$

From now on I'll omit the (t, u) parameter, assuming that we've picked a (t, u) and calculated the quantities $\bar{A}, \bar{C}, \bar{D}$ as simple scalars. So the equation we must solve is

$$\tilde{x}^3 + 3\bar{C}\tilde{x} + \bar{D} = 0 \tag{0.1}$$

and all our quantities satisfy the identity

$$\bar{D}^2 + 4\bar{C}^3 = -\tilde{A}^2\Delta \tag{0.2}$$

I have, so far, finished the solution only for the cases where $\Delta = 0$ (double roots) and where $\Delta < 0$ (one real root and a complex conjugate pair). In this installment I will address the case where $\Delta > 0$, which will yield three distinct real roots. To get the three roots in terms of just the values of \bar{C}, \bar{D} we must perform a trick similar to the one we did to solve the $\Delta < 0$ case. So let's review that trick.

The $\Delta < 0$ case (one real root)

The main trick to solving equation (0.1) is to match it up with the identity

$$(p+q)^3 - 3pq(p+q) - (p^3 + q^3) = 0$$

The match up looks like

$$\underbrace{(p+q)^3}_{\tilde{x}} - \underbrace{3pq}_{3\bar{C}} \underbrace{(p+q)}_{\tilde{x}} - \underbrace{(p^3 + q^3)}_{\bar{D}} = 0$$

In other words, if we can find p and q that satisfy

$$\begin{aligned} -pq &= \bar{C} \\ -p^3 - q^3 &= \bar{D} \end{aligned} \tag{0.3}$$

then our desired answer is simply

$$\tilde{x} = p + q$$

Some work with equations (0.3) ultimately led to the formulas

$$\begin{aligned} p &= \sqrt[3]{\frac{-\bar{D} + \sqrt{\bar{D}^2 + 4\bar{C}^3}}{2}} \\ q &= \sqrt[3]{\frac{-\bar{D} - \sqrt{\bar{D}^2 + 4\bar{C}^3}}{2}} \end{aligned}$$

The $\Delta > 0$ case (three real roots)

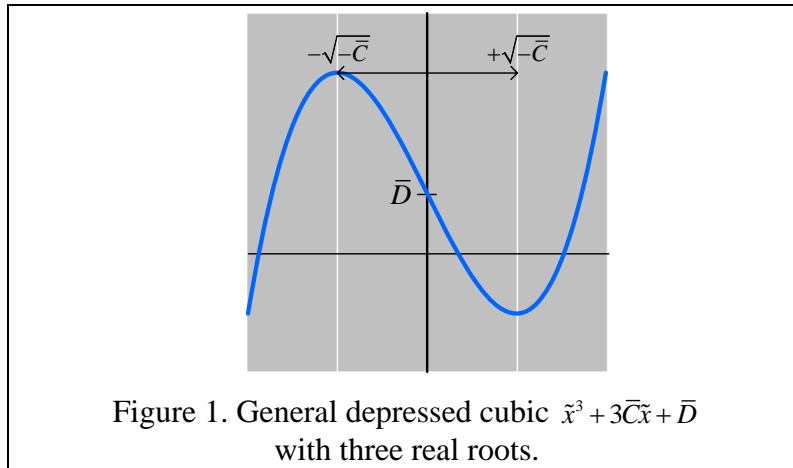
In this case we are going to match up equation (0.1) with a different identity. Nickalls [1] gives the nicest description I've found on how to do this; I've used that as the inspiration for my discussion here. In explaining this I am going to go back and forth between various algebraic and geometric derivations, sometimes pointing out relationships in more than one way. This is just to enhance our intuition and understanding of the situation so that we can be able to generate a range of solution formulations. This will help us when we start to investigate numeric stability problems.

The first root

The first thing to notice is that if a cubic has three real roots it will look something like Figure 1. Furthermore, a depressed cubic like equation (0.1) will have its second derivative be zero at $\tilde{x} = 0$ (since $\tilde{B} = 0$). This puts its inflection point on the vertical axis, with the local maximum and minimum symmetrically on either side at $\tilde{x} = \pm\sqrt{-\bar{C}}$, as figure 1 also shows. For the square root to make sense we would like to be reassured that $\bar{C} < 0$. Well, the derivative of equation (0.1), evaluated at $\tilde{x} = 0$, is $3\bar{C}$, and

figure 1 shows that this is negative. We can also see that \bar{C} must be negative when $\Delta > 0$ by rearranging the identity of equation (0.2) as $4\bar{C}^3 = -\bar{D}^2 - \bar{A}^2\Delta$.

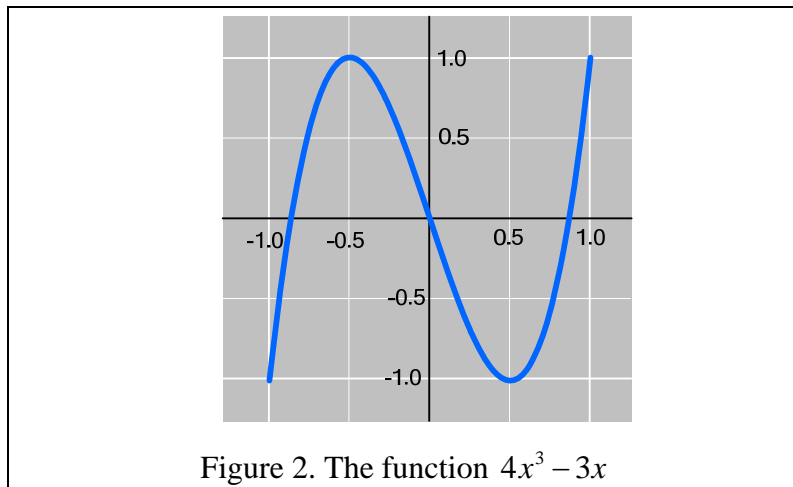
Another property of a depressed cubic is that its roots sum to zero; a property not visually obvious from Figure 1, but that we will show shortly in a more geometrically obvious way.



Now comes the trick. We are going to try to make formula (0.1) look like the trigonometric identity

$$4\cos^3\theta - 3\cos\theta - \cos3\theta = 0$$

and then associate the quantity $\cos\theta$ with x . To see what it will take to do this, let's look at a plot of the function $4x^3 - 3x$ in figure 2.



The two differences between figures 1 and 2 are a horizontal scale and a vertical displacement. The horizontal scale needs to map the quantity $\sqrt{-\bar{C}}$ to the value $1/2$, so we define the transformation

$$\tilde{x} = 2\sqrt{-\bar{C}} \tilde{\tilde{x}} \tag{0.4}$$

How to Solve a Cubic Equation – Part 4

I realize that the stack of hats on the x begins to make it look like reference [2], but it frees us to use other decorations such as subscripts to distinguish between the roots. Anyway, plugging equation (0.4) into equation (0.1) gives us

$$\begin{aligned} (2\sqrt{-\bar{C}}\tilde{x})^3 + 3\bar{C}(2\sqrt{-\bar{C}}\tilde{x}) + \bar{D} &= \\ -8\bar{C}\sqrt{-\bar{C}}\tilde{x}^3 + 6\bar{C}\sqrt{-\bar{C}}\tilde{x} + \bar{D} &= 0 \end{aligned}$$

and dividing out the (positive) quantity $-2\bar{C}\sqrt{-\bar{C}} = \sqrt{-4\bar{C}^3}$ gives us

$$4\tilde{x}^3 - 3\tilde{x} + \frac{\bar{D}}{-2\bar{C}\sqrt{-\bar{C}}} = 0 \quad (0.5)$$

Compare this with the identity

$$4\cos^3\theta - 3\cos\theta - \cos 3\theta = 0$$

and we see that

$$\begin{aligned} \tilde{x} &= \cos\theta \\ \cos 3\theta &= \frac{\bar{D}}{2\bar{C}\sqrt{-\bar{C}}} \end{aligned}$$

In other words, one of the roots of the scaled depressed cubic is

$$\tilde{x}_1 = \cos\left(\frac{1}{3}\cos^{-1}\left(\frac{\bar{D}}{2\bar{C}\sqrt{-\bar{C}}}\right)\right)$$

Now let's say a few words about the vertical displacement value for the cubic of equation (0.5). Since it equals the quantity $-\cos 3\theta$ we need to reassure ourselves that it is always in the range $-1\dots+1$. Our first reassurance comes from looking at figure 2 and seeing that if we displace vertically by more than ± 1 we will no longer have a cubic with three real roots. Our second reassurance comes from rewriting equation (0.2) as

$$\bar{D}^2 + \tilde{A}^2\Delta = -4\bar{C}^3 \quad (0.6)$$

so of course

$$\bar{D}^2 < -4\bar{C}^3$$

Since the right side is positive we can divide by it to give

$$\frac{\bar{D}^2}{-4\bar{C}^3} = \left(\frac{\bar{D}}{\sqrt{-4\bar{C}^3}}\right)^2 < 1$$

As a more geometric demonstration, look at equation (0.6) and think of it as an instance of the Pythagorean Theorem. This gives us the triangle in figure 3, which illustrates the geometrical relation between 3θ and \bar{C}, \bar{D} .

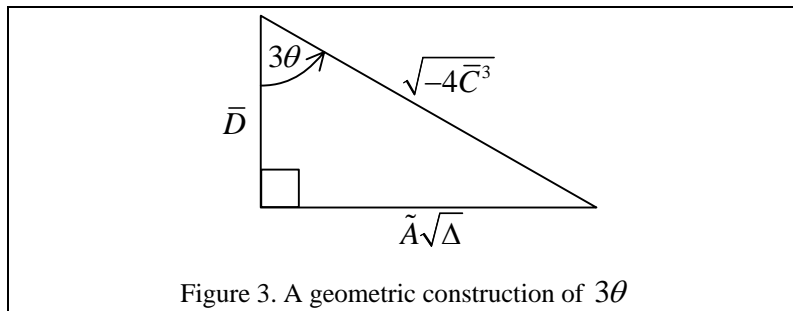


Figure 3. A geometric construction of 3θ

The other roots

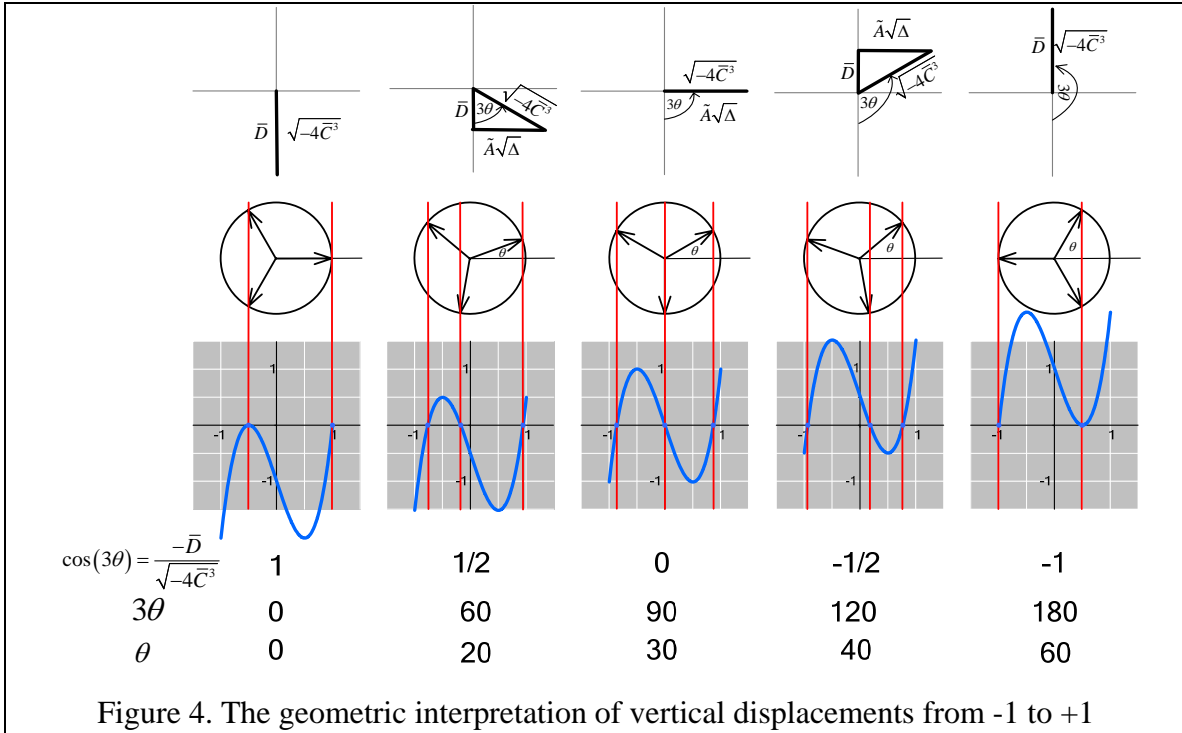
We now know one root. How do we get the other two? The mathematicians would tell us to divide the expression $(\tilde{x} - \tilde{x}_1)$ out of equation (0.5) and solve the resulting quadratic. But there is a simpler way. First notice in figure 2 that for any vertical displacement between -1 and $+1$ the three roots will satisfy

$$\begin{aligned}\tilde{x}_1 &\in \{+0.5, +1.0\} \\ \tilde{x}_2 &\in \{-0.5, +0.5\} \\ \tilde{x}_3 &\in \{-1.0, -0.5\}\end{aligned}$$

Now look at figure 4 where I've shown the geometric situation for several vertical displacements. You can see that the value of θ ranges in value from 0 to 60 degrees. The unit circles above each plot contain three vectors spaced 120 degrees apart. The vector sum of these three vectors is zero, and in particular the x component of that vector sum is zero. This is the promised geometric demonstration that the sum of the roots is zero. Figure 4 thus shows that the three roots of the scaled depressed cubic are

$$\begin{aligned}\tilde{x}_1 &= \cos \theta \\ \tilde{x}_2 &= \cos(\theta - 120^\circ) \\ \tilde{x}_3 &= \cos(\theta + 120^\circ)\end{aligned}$$

As a check, if you triple any of these three arguments to the cosine you will get the same result, 3θ , so they all satisfy the cubic.



Alternate calculations

There are various alternate ways to do this calculation. We can apply some trig identities and come up with

$$\begin{aligned}\tilde{x}_1 &= \cos \theta \\ \tilde{x}_2 &= -\frac{1}{2} \cos \theta + \frac{\sqrt{3}}{2} \sin \theta \\ \tilde{x}_3 &= -\frac{1}{2} \cos \theta - \frac{\sqrt{3}}{2} \sin \theta\end{aligned}$$

Or, since the roots sum to zero we could use

$$\begin{aligned}\tilde{x}_1 &= \cos \theta \\ \tilde{x}_2 &= -\frac{1}{2} \cos \theta + \frac{\sqrt{3}}{2} \sin \theta \\ \tilde{x}_3 &= -\tilde{x}_1 - \tilde{x}_2\end{aligned}$$

Slightly more exotically, we can use the symmetry of the cubic in figure 2 and do the calculation as follows;

$$\begin{aligned}\tilde{x}_1 &= +\cos\left(\frac{1}{3}\cos^{-1}\left(+\frac{\bar{D}}{2\bar{C}\sqrt{-\bar{C}}}\right)\right) \\ \tilde{x}_3 &= -\cos\left(\frac{1}{3}\cos^{-1}\left(-\frac{\bar{D}}{2\bar{C}\sqrt{-\bar{C}}}\right)\right) \\ \tilde{x}_2 &= -\tilde{x}_1 - \tilde{x}_3\end{aligned}$$

This seems like a lot of work, but has some numerical advantages. Also the roots \tilde{x}_1 and \tilde{x}_3 can be calculated in parallel on modern GPU shaders so the extra calculation is basically free.

Finally, I have had the experience that cubics that are very close to having a double root can sometimes, due to numerical noise, generate a value of $\bar{D}/2\bar{C}\sqrt{-\bar{C}}$ that is slightly greater than +1 or slightly less than -1 causing the arccosine function to croak. Figure 3 gives us a way out of this by showing us how to get 3θ by using a somewhat more bulletproof arctangent instead of an arc cosine. The two-parameter version is what we want here and looks like:

$$\theta = \frac{1}{3} \operatorname{atan2}(\tilde{A}\sqrt{\Delta}, -\bar{D})$$

All of these variants give us a toolkit that we can use to improve the numerical properties of our algorithm.

Some Observations

It should not be surprising that in solving cubic equations, sooner or later you are going to have to take the cube root of something. But it might not be immediately obvious what cubic equations have to do with arc cosines. Well, in the general case of complex coefficients and roots we would need to take the cube root of a complex number. How do you do this? You first express it in polar coordinates as a magnitude and a unit vector. You take cube root of the (real) magnitude, and combine it with a unit vector that is at 1/3 of the angle of the original unit vector. In our case, where coefficients are real, we have effectively arranged things so that we need to do only one of these operations. In the single-root case we only need to take the cube root of a real number; in the three-root case we only need to do the 1/3 angle calculation.

But we now have a solution to cubic polynomials (and also quadratic polynomials from earlier articles) in closed form... or do we? The only reason that we could express the roots in closed form is that we could include some exotic functions in the expression. All of the formulas involved some transcendental function T and required evaluation of

$$T\left(\frac{1}{n}T^{-1}(x)\right) \tag{0.7}$$

In the case of quadratic equations and single-root cubic equations, T is the exponential function. In the case of three-root cubics, T is the cosine function. (This makes sense since, in the complex domain, the exponential and the cosine are basically the same function.)

But whadaya think, transcendental functions grow on trees? They must be approximated by some sort of power series, or perhaps by some iterative scheme. So hidden inside of our closed-form solutions might be some iterative calculations. Our entire process of depressing and scaling the polynomials was basically to map the polynomial into a space where it matches a standard transcendental function. But that's OK. These functions have been instantiated in fairly efficient hardware algorithms. In fact to solve a quadratic where $n=2$ and $T=\exp$ the entire three step calculation of equation (0.7) has been encapsulated into a single hardware operation—it's called a square root. Perhaps future hardware can do the same thing with cube root, and the $\cos\left(\frac{1}{3}\cos^{-1}(x)\right)$ function.

We're not done yet

There is still a bit of a problem with this algorithm. In the real world of limited precision floating point it doesn't always get the right answer. And, as in human affairs, a significant cause of problems is depression. In my next and final installment of this topic, I'll exercise our collection of root finding techniques to generate an algorithm that compensates for this and gets the right answer (pretty much) all the time.

References

[1] Nickalls, R. W. D., "A new approach to solving the cubic: Cardan's solution revealed", *The Mathematical Gazette*, November 1993, vol. 77, pp 354-358.

Available online at <http://www.m-a.org.uk/docs/library/2059.pdf>

[2] Seuss, Dr., *The 500 hats of Bartholomew Cubbins*, Vanguard Press, 1938.