

MATH 228 lecture notes for November 8, 10, and 12

Russell Milne

November 2021

1 November 8: Solving systems of linear, homogeneous first-order ODEs with constant coefficients

As with uncoupled ODEs, we'll start our journey into systems of ODEs with the simplest possible case, primarily because it's one of the few for which there exist methods of solving analytically. More specifically, we'll start with linear, homogeneous first-order ODEs with constant coefficients. Remember at the very beginning of this course that we determined that the solution to the ODE $x' = kx$ is $x(t) = Ce^{kt}$. Suppose that we take a similar system:

$$\begin{cases} \frac{dx_1}{dt} = ax_1 + bx_2 \\ \frac{dx_2}{dt} = cx_1 + dx_2 \end{cases} \quad (1)$$

The solutions for $x_1(t)$ and $x_2(t)$ will also contain exponential functions. However, because each ODE in this system has two terms, each affecting the growth or decay of x_1 or x_2 over time, a single exponential function might not be enough to describe the dynamics of the two state variables. Instead, the solutions for x_1 and x_2 will each be a linear combination of two different exponential functions. An exception to this is if one of the state variables is uncoupled, i.e. its rate of change only depends on itself (and not any of the other state variables). (In the above system, that would mean that $b = 0$ and/or $c = 0$.) In this case, the ODE for that state variable can be solved independently of the rest of the system.

To solve this system of differential equations, we will need to rewrite it as a matrix. We can construct a vector whose entries are the state variables of the system, like this:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (2)$$

Using this vector-valued function, as well as some linear algebra, we can rewrite our system of ODEs as a matrix-vector equation:

$$\mathbf{x}' = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \mathbf{x} = \mathbf{M}\mathbf{x} \quad (3)$$

We will now assume that our vector-valued function \mathbf{x} will consist of exponential functions. In other words, for some vector \mathbf{u} that we will eventually solve for the entries of, and some constant r that we will also eventually find, we will assume the following:

$$\mathbf{x} = \mathbf{u}e^{rt} \implies \mathbf{x}' = r\mathbf{u}e^{rt} \quad (4)$$

Plugging this into our original system, we get the following:

$$r\mathbf{u}e^{rt} = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \mathbf{u}e^{rt} = \mathbf{M}\mathbf{u}e^{rt} \quad (5)$$

Because e^{rt} is nonzero, we can divide both sides by it. We then can rearrange terms to get our equation into the following form:

$$(\mathbf{M} - r\mathbf{I})\mathbf{u} = 0 \quad (6)$$

This looks very similar to the setup of an eigenvalue problem so far, which isn't a coincidence. Continuing with the theme, we need $\mathbf{M} - r\mathbf{I}$ to have a nontrivial kernel in order to get some solutions other than $\mathbf{u} = 0$. This occurs when the matrix $\mathbf{M} - r\mathbf{I}$ has a determinant of zero, so we will need to solve for the values of r that make it so. In this case, finding the determinant of a 2×2 matrix is straightforward:

$$\det(\mathbf{M} - r\mathbf{I}) = \begin{vmatrix} a-r & b \\ c & d-r \end{vmatrix} = (a-r)(d-r) - bc = r^2 + (-a-d)r + (ad-bc) \quad (7)$$

This produces a polynomial in r that can be solved. This is indeed the characteristic polynomial for the matrix \mathbf{M} , meaning that we need to find the eigenvalues of \mathbf{M} as part of finding the solution to our system. In this specific case, since we have assumed a 2×2 system, we will have at most 2 distinct eigenvalues, and therefore at most 2 values of r . Continuing this process, we still need to solve for \mathbf{u} for each value of r , which is equivalent to finding the eigenvectors of \mathbf{M} . Once we do this, then as we assumed that $\mathbf{x} = \mathbf{u}e^{rt}$ earlier, we will have all of the information necessary to construct a solution \mathbf{x} . (Note how I say "a solution" instead of "the solution". This is because there will be multiple eigenvalues and eigenvectors and hence multiple solutions; I'll elaborate on this later.)

In order to illustrate this, I will work through a concrete example. Suppose we have the following system:

$$\begin{cases} \frac{dx_1}{dt} = x_1 + 2x_2 \\ \frac{dx_2}{dt} = 3x_1 + 2x_2 \end{cases} \quad (8)$$

We can turn this system into a matrix \mathbf{M} , and find its eigenvalues and eigenvectors:

$$\mathbf{M} = \begin{bmatrix} 1 & 2 \\ 3 & 2 \end{bmatrix} \implies \det(\mathbf{M} - r\mathbf{I}) = \begin{vmatrix} 1-r & 2 \\ 3 & 2-r \end{vmatrix} = r^2 - 3r - 4 \quad (9)$$

The characteristic polynomial factors into $(r-4)(r+1)$, which means that we have eigenvalues of $r = 4$ and $r = -1$. Now, we will plug in these values to $\mathbf{M} - r\mathbf{I}$ in order to find the eigenvectors of \mathbf{M} . Let's start with $r = 4$:

$$r = 4 \implies \begin{bmatrix} -3 & 2 \\ 3 & -2 \end{bmatrix} \mathbf{u} = \begin{bmatrix} -3 & 2 \\ 3 & -2 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = 0 \quad (10)$$

When solving this system for u_1 and u_2 , we get the relation that $3u_1 = 2u_2$. This points to $[2 \ 3]^T$ being an eigenvector with the associated eigenvalue of $r = 4$. Now, we will plug in $r = -1$:

$$r = -1 \implies \begin{bmatrix} 2 & 2 \\ 3 & 3 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = 0 \quad (11)$$

This brings us to $u_1 = -u_2$, and hence an eigenvector of $[-1 \ 1]^T$ (or any scalar multiple thereof, of course). We now have two solutions to the system, each one of the following form:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \mathbf{u}e^{rt} = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} e^{rt} \quad (12)$$

In the specific case that we are dealing with, we have the following two solutions:

$$\mathbf{x}^{(1)} = \begin{bmatrix} 2 \\ 3 \end{bmatrix} e^{4t}; \quad \mathbf{x}^{(2)} = \begin{bmatrix} -1 \\ 1 \end{bmatrix} e^{-t} \quad (13)$$

What can we do with these? Well, as with finding solutions of second-order differential equations, our general solution will just be a linear combination of $\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$. (This is referred to as the principle of superposition. It can be easily proved by taking the derivatives of $\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$, plugging them into the original system, and determining what cancels.) As a result of all this, we can get the following general solutions for x_1 and x_2 in the case that we're looking at. For C_1 and C_2 arbitrary constants, which can be obtained by applying initial conditions, we have:

$$\begin{cases} x_1 = 2C_1e^{4t} - C_2e^{-t} \\ x_2 = 3C_1e^{4t} + C_2e^{-t} \end{cases} \quad (14)$$

Note that we have two basis functions that make up our solution for each state variable, just like we did when we were solving second-order ODEs. This is not an accident. Remember that for an ODE of the form $y'' + ay' + by = 0$, we can make the substitution $z = y'$ and get the following system:

$$\begin{cases} y' = z \\ z' = -by - az \end{cases} \quad (15)$$

Here, the solutions for both y and z will consist of two basis functions each, which can be found using the same method as we went through above.

When we go through this process for finding solutions, it is important that the solutions that we come up with are linearly independent. Any set of solutions of a system of differential equations in which all of the individual solutions are linearly independent over some interval is called a fundamental set of solutions for that system over that interval. Determining that some solutions of a system of ODEs form a fundamental set of solutions can be done by evaluating the Wronskian of the functions. However, if the ODEs in the system are linear and homogeneous with constant coefficients, then the solutions will all be exponentials and hence evaluating the Wronskian becomes very easy.

In fact, for a given set of solutions to a linear, homogeneous system of ODEs over some interval, the Wronskian will either be identically zero over the entire interval or will never be zero at any point in the interval. This holds regardless of what the coefficients are; they need not be constants. This result was first proved by Niels Henrik Abel, and therefore is one of the many results that are referred to as “Abel’s theorem”.

In the example that I showed you in this lecture, everything worked out well, since the matrix \mathbf{M} representing the specific system we were attempting to solve had eigenvalues that were real and distinct. However, not all matrices (or even all real-valued matrices) have this property. Of the ones that don’t, there are two main categories that they can fall into. One possibility for real-valued matrices is that we have a repeated eigenvalue, such as in the following system:

$$\begin{cases} \frac{dx_1}{dt} = x_1 - x_2 \\ \frac{dx_2}{dt} = x_1 + 3x_2 \end{cases} \quad (16)$$

The other is that the eigenvalues are complex. In real-valued matrices, this means that they will be complex conjugates, but if we have a complex-valued matrix, the eigenvalues can be any complex number. In this system, for instance, both eigenvalues are purely imaginary:

$$\begin{cases} \frac{dx_1}{dt} = -x_1 + 3x_2 \\ \frac{dx_2}{dt} = -2x_1 + x_2 \end{cases} \quad (17)$$

As it turns out, solving each of these cases requires some of the same techniques that you have previously applied when finding solutions to second-order ODEs with constant coefficients, because of our assumption (as always) that the solutions are exponentials of some kind. This means that these cases are actually easier to solve than you might think!

2 November 10: Solving systems of ODEs with repeated or complex eigenvalues, or with forcing terms

In the previous lecture, we looked at an easy example of a system of linear, homogeneous first-order ODEs with constant coefficients. However, in that case, everything worked out about as well as it possibly could have due to the eigenvalues being real and distinct. What if instead we have complex eigenvalues? For instance, consider the following system:

$$\begin{cases} \frac{dx_1}{dt} = -x_1 + 3x_2 \\ \frac{dx_2}{dt} = -2x_1 + x_2 \end{cases} \quad (18)$$

In order to solve this, we first need to find the eigenvalues and eigenvectors of the matrix corresponding to this system. We know how to do this already:

$$\mathbf{M} = \begin{bmatrix} -1 & 3 \\ -2 & 1 \end{bmatrix} \implies \det(\mathbf{M} - r\mathbf{I}) = (-1 - r)(1 - r) + 6 \quad (19)$$

A closer look at the characteristic polynomial of our matrix \mathbf{M} , or indeed the entries of \mathbf{M} , should reveal something important. Because the diagonal entries of \mathbf{M} are 1 and -1, we will get terms of $+r$ and $-r$, which cancel each other out. Additionally, the constant term in the characteristic polynomial will be positive. (From the previous lecture, we know the formula for all of the coefficients of the characteristic polynomial of a 2×2 matrix. The constant term is $m_{1,1}m_{2,2} - m_{1,2}m_{2,1}$, which in this case is $-1 + 6 = 5$.) This means that we can get our characteristic polynomial into the form $r^2 = -5$, which yields the purely imaginary roots of $r = \pm i\sqrt{5}$. As we're still assuming that our solutions will be exponentials of some kind, this means that we will get some linear combination of $e^{i\sqrt{5}t}$ and $e^{-i\sqrt{5}t}$ for solutions. These are, of course, both wavefunctions.

Since we have two eigenvalues, r_1 and r_2 , that are complex conjugates, the corresponding eigenvectors \mathbf{u}^1 and \mathbf{u}^2 will also be complex conjugates. We can see this by taking the complex conjugate of the equation $(\mathbf{M} - r_1\mathbf{I})\mathbf{u}^{(1)} = 0$, which we know is satisfied because r_1 is an eigenvalue with associated eigenvector \mathbf{u}^1 . This results in the following relation:

$$\overline{(\mathbf{M} - r_1\mathbf{I})\mathbf{u}^{(1)}} = 0 \implies (\mathbf{M} - \bar{r}_1\mathbf{I})\overline{\mathbf{u}^{(1)}} = 0 \quad (20)$$

This is true because \mathbf{M} and \mathbf{I} are real-valued (by our assumption), so they are their own complex conjugates. Hence, because $\bar{r}_1 = r_2$, the eigenvector associated with r_2 is the complex conjugate of that associated with r_1 . This makes the process of finding two linearly independent solutions much easier.

Let's find the eigenvector associated with $r_1 = i\sqrt{5}$. Plugging this into our system yields the following:

$$\begin{bmatrix} -1 - i\sqrt{5} & 3 \\ -2 & 1 - i\sqrt{5} \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = 0 \quad (21)$$

We therefore get the following relations for u_1 and u_2 :

$$\begin{cases} (1 + i\sqrt{5})u_1 = 3u_2 \\ 2u_1 = (1 - i\sqrt{5})u_2 \end{cases} \quad (22)$$

From the second of these equations, we can obtain the eigenvector $\mathbf{u}^{(1)} = \left[\frac{1}{2}(1 - i\sqrt{5}) \ 1\right]^T$, which has the associated eigenvalue $r = i\sqrt{5}$. If instead we take $r = -i\sqrt{5}$, then we will get a vector for $\mathbf{u}^{(2)}$ in which every entry is the complex conjugate of the corresponding entry in $\mathbf{u}^{(1)}$. Specifically, this is $\left[\frac{1}{2}(1 + i\sqrt{5}) \ 1\right]^T$. Now that we know $r_1, r_2, \mathbf{u}^{(1)}$ and $\mathbf{u}^{(2)}$, we can get a general solution for this problem in the same way as if the eigenvalues were both real.

As a matter of fact, if the eigenvalues are complex conjugates, we can even get real-valued solutions. Suppose that we have some system of ODEs where the characteristic polynomial has complex roots, so that the equation $(\mathbf{M} - r\mathbf{I})\mathbf{u} = 0$ will feature complex values for r and the entries of \mathbf{u} . In other words, if we have $\mathbf{x} = \mathbf{u}e^{rt}$, we can assume that $\mathbf{x} = (\mathbf{v} + i\mathbf{w})e^{(\alpha+i\beta)t}$, where \mathbf{v} and \mathbf{w} have real elements, and α and β are both real constants. After expanding everything out, we therefore get the following solution:

$$\mathbf{x} = e^{\alpha t}(\mathbf{v} \cos \beta t - \mathbf{w} \sin \beta t) + ie^{\alpha t}(\mathbf{v} \sin \beta t + \mathbf{w} \cos \beta t) \quad (23)$$

However, this is just one solution. There will be two, and the second one will have values for \mathbf{u} and r that are the complex conjugates of what we found above. As a matter of fact, by taking scalar multiples of both of these solutions and adding them together, we get that both of the following are solutions for \mathbf{x} :

$$\mathbf{x} = e^{\alpha t}(\mathbf{v} \cos \beta t - \mathbf{w} \sin \beta t) \quad (24)$$

$$\mathbf{x} = e^{\alpha t}(\mathbf{v} \sin \beta t + \mathbf{w} \cos \beta t) \quad (25)$$

This, once again, is similar to what we found for second-order linear ODEs with constant coefficients. The functions above are linearly independent (this can be checked using the Wronskian), and we can use them as a basis for the general solution. (This basis is actually preferred, since the functions involved are real-valued.)

So, we have solved the case of a 2×2 system of ODEs with real, constant coefficients where the eigenvalues are complex. What about higher-dimensional systems? Well, in a 3×3 system of ODEs (once again, with real, constant coefficients), you could potentially have a matrix where one of the eigenvalues is real and the other two are complex conjugates of each other. In this case, you can find the solution associated with the real eigenvalue as you normally would, and the two solutions associated with the complex eigenvalues using the method above. For 4×4 systems, and those whose dimensionality is even higher,

you could have multiple different complex conjugate pairs, and hence you would have to do the above method multiple times.

What about if the matrix that we form out of our system has repeated eigenvalues? In some cases, that may not matter. For instance, in any $n \times n$ matrix that has repeated eigenvalues but is real-valued and symmetric, we will have n linearly independent eigenvectors regardless. (A symmetric matrix is one that is equal to its own transpose. In other words, if \mathbf{M} is the matrix, then $m_{ij} = m_{ji} \forall i, j$.) One example of this is the identity matrix \mathbf{I} , which we saw previously has an eigenvalue of multiplicity 2 but two linearly independent eigenvectors.

What if the matrix is not symmetric? Then, we need to find another linearly independent solution somehow. I will illustrate this with an example. Suppose we have the following system:

$$\begin{cases} \frac{dx_1}{dt} = x_1 - x_2 \\ \frac{dx_2}{dt} = x_1 + 3x_2 \end{cases} \quad (26)$$

Using what we already know, we can get $r = 2$ and $\mathbf{u}^{(1)} = [1 \ -1]^T$, but not a second solution. Our first guess for another one will be to assume something of the form $\mathbf{u}^{(2)}te^{2t}$ (keeping r the same), as we did for second-order ODEs when this challenge arose. In order for this to be a solution, it must satisfy $\mathbf{x}' = \mathbf{M}\mathbf{x}$, or alternatively the following:

$$2\mathbf{u}^{(2)}te^{2t} + \mathbf{u}^{(2)}e^{2t} = \begin{bmatrix} 1 & -1 \\ 1 & 3 \end{bmatrix} \mathbf{u}^{(2)}te^{2t} \quad (27)$$

However, this equation cannot be true unless $\mathbf{u}^{(2)} = 0$, since we have a term of e^{2t} on the left-hand side that does not show up anywhere on the right-hand side. This means that we have to assume instead that our second solution takes the form of $\mathbf{v}te^{2t} + \mathbf{w}e^{2t}$, so that we can properly balance terms. We then get the following:

$$2\mathbf{v}te^{2t} + (\mathbf{v} + 2\mathbf{w})e^{2t} = \mathbf{M}(\mathbf{v}te^{2t} + \mathbf{w}e^{2t}) = \begin{bmatrix} 1 & -1 \\ 1 & 3 \end{bmatrix} (\mathbf{v}te^{2t} + \mathbf{w}e^{2t}) \quad (28)$$

If we set the coefficients on both sides equal to each other, we get $2\mathbf{v} = \mathbf{M}\mathbf{v}$ and $(\mathbf{v} + 2\mathbf{w}) = \mathbf{M}\mathbf{w}$. The first of these equations just states that \mathbf{v} needs to be an eigenvector of \mathbf{M} . As we previously found $[1 \ -1]^T$ to be such an eigenvector, we will use that for \mathbf{v} . For the second equation, we need that $(\mathbf{M} - 2\mathbf{I})\mathbf{w} = \mathbf{v} = [1 \ -1]^T$. Solving this is only slightly different than solving an eigenvalue problem, as we will still get a relation between w_1 and w_2 . Specifically, we get the following:

$$\begin{bmatrix} -1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \quad (29)$$

It therefore follows that if $w_1 = k$, $w_2 = -k - 1$. From this, we can obtain \mathbf{w} :

$$\mathbf{w} = \begin{bmatrix} 0 \\ -1 \end{bmatrix} + k \begin{bmatrix} 1 \\ -1 \end{bmatrix} \quad (30)$$

We therefore have everything we need to construct our second solution, which is the following:

$$\mathbf{x} = \mathbf{v}te^{2t} + \mathbf{w}e^{2t} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} te^{2t} + \begin{bmatrix} 0 \\ -1 \end{bmatrix} e^{2t} + k \begin{bmatrix} 1 \\ -1 \end{bmatrix} e^{2t} \quad (31)$$

The last of these terms has the same form as what we got for the first solution and hence adds nothing new. Therefore, we can ignore it by assuming that $k = 0$. After this, taking a linear combination of the above solution with the first one will yield the general solution of the system.

What if we have eigenvalues with multiplicity more than 2? This can come up if we have a system of ODEs that is larger than just 2×2 , for instance an eigenvalue of multiplicity 3 in a 3×3 system. In that case, you will just need extra powers of t to keep generating basis functions that are linearly independent. So, if your first basis function is some vector times e^{rt} for the repeated eigenvalue r , and your second is $\mathbf{v}te^{rt} + \mathbf{w}e^{rt}$ for vectors \mathbf{v} and \mathbf{w} (as we went through in the above example), then a hypothetical third basis vector would be $\mathbf{a}t^2e^{rt} + \mathbf{b}te^{rt} + \mathbf{c}e^{rt}$ for vectors \mathbf{a} , \mathbf{b} and \mathbf{c} , a fourth one would also include a term with t^3e^{rt} , and so on.

So far, we've just looked at homogeneous systems. What happens when we include a forcing term? In that case, it's actually quite easy to get solutions, so long as the matrix \mathbf{M} corresponding to the system is diagonalizable. (This is equivalent to it being of full rank, or in other words, having n linearly independent eigenvectors if it is an $n \times n$ matrix.) So, suppose we have the problem $\mathbf{x}' = \mathbf{M}\mathbf{x} + \mathbf{g}(t)$, where \mathbf{g} is a vector of time-dependent functions (i.e. our forcing terms for each ODE). Now, let \mathbf{T} be a matrix whose columns are the eigenvectors of \mathbf{M} , and let $\mathbf{x} = \mathbf{T}\mathbf{y}$ for some vector \mathbf{y} . We can substitute this into our equation $\mathbf{x}' = \mathbf{M}\mathbf{x} + \mathbf{g}(t)$ to get the following:

$$\mathbf{T}\mathbf{y}' = \mathbf{M}\mathbf{T}\mathbf{y} + \mathbf{g}(t) \quad (32)$$

The derivative of $\mathbf{T}\mathbf{y}$ is just $\mathbf{T}\mathbf{y}'$, as all of the entries in \mathbf{T} are constants. When we left-multiply everything in this equation by the inverse of \mathbf{T} (namely \mathbf{T}^{-1}) to isolate \mathbf{y}' , we get the following:

$$\mathbf{y}' = (\mathbf{T}^{-1}\mathbf{M}\mathbf{T})\mathbf{y} + \mathbf{T}^{-1}\mathbf{g}(t) \quad (33)$$

However, because of our choice for \mathbf{T} , the matrix $\mathbf{T}^{-1}\mathbf{M}\mathbf{T}$ is a diagonal matrix whose entries are the eigenvalues of \mathbf{M} , in the same columns as the corresponding eigenvectors of \mathbf{M} are in \mathbf{T} . Additionally, $\mathbf{T}^{-1}\mathbf{g}$ is just another vector of time-dependent functions (which we will call \mathbf{h}). Because the matrix $\mathbf{T}^{-1}\mathbf{M}\mathbf{T}$ is diagonal and $\mathbf{T}^{-1}\mathbf{g}$ contains no variables other than t , we now have

an uncoupled system, in which we can solve each equation separately. Specifically, we will get the following, for $i = 1, \dots, n$ and r_i the eigenvalues of \mathbf{M} :

$$y_i'(t) = r_i y_i + h_i(t) \quad (34)$$

This means that we can solve for each y_i using techniques that we already know. After constructing our vector \mathbf{y} , we just need to left-multiply it by \mathbf{T} to get the solution \mathbf{x} , since we defined $\mathbf{x} = \mathbf{T}\mathbf{y}$ earlier.

3 November 12: Introduction to numerical integration, nonlinear dynamics, and mathematical modelling

Throughout this course, we have focused mostly on differential equations (or systems thereof) which are easy to solve by hand. Hence, you've mainly seen solutions that are composed of exponentials, trigonometric functions (which are just complex exponentials), and polynomials (often in the form of Taylor series). This is because the theory behind differential equations first arose during the 1800s, when calculations were typically done by hand. However, at the same time, physicists and other scientists were coming up with differential equations that described natural phenomena, of which most were (and are) quite difficult to find an analytical solution for. Bessel's equation, for instance, takes the following form:

$$x^2 y'' + xy' + (x^2 - \alpha^2)y = 0 \quad (35)$$

This is still a linear differential equation, and its singular point at $x = 0$ is regular, so we at least have the tools to solve it by hand if we wanted to. However, we can't get a closed-form solution, and are instead left with an infinite polynomial series (here α is as defined above):

$$J_\alpha(x) = \sum_{n=0}^{\infty} \frac{(-1)^n}{n! \cdot \Gamma(n + \alpha + 1)} \left(\frac{x}{2}\right)^{2n+\alpha} \quad (36)$$

Despite the fact that Bessel's equation looks relatively "nice" compared to other ODEs that we could come up with, finding the solution by hand is very challenging, and even requires us to use the Gamma function. If we restricted ourselves to solving by hand, this would leave us with a dilemma: we can use ODEs to describe processes in physics, chemistry, biology, and other fields that are definitely worth knowing, but solving them by hand is difficult. This is actually part of what prompted the theory behind numerical integration: the Runge-Kutta method, which is still one of the most popular methods for numerically integrating ODEs today, was developed at the beginning of the 1900s (after more complicated differential equations had proliferated) by Carl Runge and Wilhelm Kutta.

Additionally, in this course, we have only briefly touched upon systems of differential equations, otherwise known as dynamical systems. This is for the same reasons: once you get beyond linear systems with constant coefficients, analytical solutions are hard or impossible to find. This is despite the fact that a great many problems arise from the need to understand how different things interact with each other, which necessitates solving differential equations that are coupled together in a system rather than being independent from one another. I will illustrate this with an example. Suppose a biologist is studying two species, of which one is a predator and one is its prey. The biologist won't know the specific functions that describe the population sizes of the prey ($N(t)$) and predators ($P(t)$), but might be able to observe their rates of change in the field. Therefore, to predict how the populations of the two species will change over time, the biologist might construct two ODEs (one for each species) using the information regarding the rates of change. This process of building differential equations (and dynamical systems) that describe things in real life based on their observed or predicted rates of change is called mathematical modelling.

Let's see if we can put together this predator-prey system ourselves. We'll start with the ODE for the prey, which is $\frac{dN}{dt}$. We know that with plentiful food and adequate living space, a population will tend to grow exponentially. The reason for this is because we can model the growth of the population by assuming that over one unit of time, each individual in the population will produce r offspring, for some constant r . This is the birth rate of the population, which can be determined from field data. If we made the admittedly unrealistic assumption that our population could grow forever, we would arrive at this ODE:

$$\frac{dN}{dt} = rN \tag{37}$$

However, just as organisms are born, they can also die. We assumed that there was a predator that would eat the prey, so let's consider how that would affect the rate of change of the prey population. Obviously, if there are more predators, the odds of a prey organism encountering one (and being eaten by it) becomes greater. Likewise, if the number of predators is constant, then they will have a greater chance of running into prey if there is more prey. This leads us to a term that scales with both predator and prey population sizes that describes how often the prey get eaten by predators. We can add it to our ODE describing the prey population:

$$\frac{dN}{dt} = rN - \alpha NP \tag{38}$$

What about the predators? What rates govern how they are born and die? For this, we will consider the life cycle of a predator. They gain energy by eating prey, and use that energy to carry out all of their vital processes (including reproduction). So, we can assume that how many predators are born depends not only on how many predators there already are, but how many prey

organisms there are as well. The predator death rate is simpler: since we have assumed that the predators don't have any predators themselves, we can just say that some constant proportion of the predator population will die over a given unit of time. This brings us to the following ODE for the predators:

$$\frac{dP}{dt} = \beta NP - mP \quad (39)$$

Combining the two into a 2×2 system gives us the following:

$$\begin{cases} \frac{dN}{dt} = rN - \alpha NP \\ \frac{dP}{dt} = \beta NP - mP \end{cases} \quad (40)$$

Since the populations of the two species can be directly measured in the field, we can also specify an initial condition:

$$\begin{cases} N(t=0) = N_0 \\ P(t=0) = P_0 \end{cases} \quad (41)$$

This is called the Lotka-Volterra model, which was independently derived by Alfred Lotka and Vito Volterra in the 1920s. One of its most famous applications is to explain the dynamics of lynx and hare populations in northern Canada; the data collected by the Hudson's Bay Company on the number of lynx and hare pelts its trappers obtained fits the model very well, even over the very long timescale over which these statistics were tracked.

Notice that both the predator and prey ODEs depend on both N and P , and that this dependence is nonlinear (due to the NP term) in both equations. This doesn't look like anything that we know how to solve. Indeed, it's all but impossible to do this analytically. (I recommend plugging the input $x' = x - xy$, $y' = xy - y$ into Wolfram Alpha to see what it gives as an output.) Therefore, we need to find out the properties of this function some other way.

Here's another example of how we can derive a mathematical model from first principles. Suppose that a chemical reaction is going on inside a cell, in which an enzyme is converting a substrate molecule into a reaction product. The substrate can freely bind to and disengage from the enzyme, but the enzyme-catalyzed reaction converting the substrate to the product is one-way. This means that we have the following reaction schematic:



There are four different interacting objects in our system, namely the enzyme (E), the substrate (S), the enzyme-substrate complex (ES), and the product (P). Therefore, in order to model it, we'll need four ODEs, which leads to a larger system than anything we've seen before. Fortunately, we have everything we need to do so in the above reaction schematic. We can assume that the enzyme and substrate bind together at a rate k_1 and unbind at a rate k_{-1} , and that once the complex is formed, the reaction proceeds at a rate k_2 . Additionally, as with predators and prey, we can assume that how often the enzyme

and substrate encounter each other (and hence bind to form the complex) is proportional to how much of both of them there is within the cell. Using this information, we can write a system of ODEs that characterizes the chemical reaction, or specifically how the concentrations of the different molecules change over time. Using the notation $e = [E]$, $s = [S]$, $c = [ES]$ (“c” for “complex”), and $p = [P]$, we have the following:

$$\begin{cases} \frac{de}{dt} = -k_1se + k_{-1}c + k_2c \\ \frac{ds}{dt} = -k_1se + k_{-1}c \\ \frac{dc}{dt} = k_1se - k_{-1}c - k_2c \\ \frac{dp}{dt} = k_2c \end{cases} \quad (43)$$

If we are initializing the model at the beginning of the reaction, we will have the following initial conditions:

$$\begin{cases} e(t=0) = e_0 \\ s(t=0) = s_0 \\ c(t=0) = 0 \\ p(t=0) = 0 \end{cases} \quad (44)$$

One thing that jumps out from the system of 4 ODEs is the fact that all of the terms in $\frac{de}{dt}$ are the opposites of the terms in $\frac{dc}{dt}$. In other words, we can see that $\frac{de}{dt} + \frac{dc}{dt} = 0$, or equivalently $e + c = e_0 + 0 = e_0$ given our initial conditions. This makes sense biologically, as the total amount of enzyme present (including enzyme that is both bound and unbound to the substrate) will remain constant throughout the reaction. This means that $e + c$ is a conserved quantity, which is important. It means that we can simplify our system of ODEs by taking $e(t) = e_0 - c(t)$, reducing the dimensionality by 1 and giving us the following:

$$\begin{cases} \frac{ds}{dt} = -k_1s(e_0 - c) + k_{-1}c \\ \frac{dc}{dt} = k_1s(e_0 - c) - (k_{-1} + k_2)c \\ \frac{dp}{dt} = k_2c \end{cases} \quad (45)$$

This is still a nonlinear system, and none of the tools we have learned so far can be used to solve it. (It’s still doable, though, so long as you make a few assumptions beforehand.)

As you can see from what we did above, it’s pretty easy to generate a dynamical system if you already know something like a reaction diagram, which shows in graphical form all of the processes where some molecules are produced from other molecules. This means that even for large, complicated biochemical processes involving many different molecules (like the one we saw during lecture), we can construct a dynamical system that describes the process just by following the arrows in the reaction diagram. This concept can be generalized to other fields in which you might want to model a large network of interacting objects. To create a “reaction diagram”, you can just draw each object in the system as its own box, and then draw arrows between boxes that interact with

each other. Note that some of these reactions may be positive (i.e. the presence of Object A causes more of Object B to be produced), and some may be negative (i.e. Object A instead causes Object B to be consumed). Additionally, the functions that describe the interactions between the different boxes may be different than just the linear terms we've mostly seen in this course. For example, in the enzyme kinetics discussed above, some terms were linear and some terms were not, depending on the order of the reaction that each term represented. Other common terms that you'll see in mathematical models include logistic growth:

$$\frac{dx}{dt} = rx \left(1 - \frac{x}{K}\right) \quad (46)$$

as well as the saturation function:

$$\frac{dx}{dt} = \frac{x}{k + x} \quad (47)$$

and a generalized version of the saturation function, which is referred to in biology as the Hill function:

$$\frac{dx}{dt} = \frac{x^n}{k + x^n} \quad (48)$$

Of course, when you are creating a model, the most important part is making sure that your mathematical terms match what is observed in real life. This means that any arrow that you draw between boxes (and the functional form that you choose to represent the boxes' interaction) needs to have some rationale based on the problem at hand.