
EECS 16B Designing Information Devices and Systems II
 Fall 2021 Note 19: Linearization and Quadratic
 Approximation

Overview

So far, we have spent a great deal of time studying linear systems described by differential equations of the form:

$$\frac{d}{dt}\vec{x}(t) = A\vec{x}(t) + B\vec{u}(t) + \vec{w}(t) \quad (1)$$

or linear difference equations (aka linear recurrence relations) of the form:

$$\vec{x}[i + 1] = A\vec{x}[i] + B\vec{u}[i] + \vec{w}[i]. \quad (2)$$

Here, the linear difference equations came from either discretizing a continuous model or were learned from data.

In the real world, however, it is very rare for physical systems to perfectly obey linear differential equations. For example:

- (a) Up until now, we have considered transistors to be binary – turning off or on at some voltage differential. But in reality, transistors work continuously, and the relevant governing equations are highly nonlinear. For example, [the Wikipedia article on BJTs has some of them](#).
- (b) Robotics and control in general has highly nonlinear dynamics. The reason for this is because the system has a variety of geometric relationships that must hold, and so there are many trigonometric functions involved, which are all non-linear but have good linear approximations. An example of this analysis may be found [when applying control to stabilize an inverted pendulum on a cart](#).
- (c) Machine learning and optimization also have highly nonlinear systems. In particular, the concept of *gradient flow* sets up the differential equation $\frac{d}{dt}x(t) = -\frac{d}{dx(t)}f(x(t))$, and if one solves for $x(t)$ then one can characterize the whole process of optimization via *gradient descent*. In modern machine learning, such as in neural networks, the function f is nonlinear.

Don't worry, we don't expect you to know anything from these examples that hasn't been taught elsewhere in the course! They are just there for flavor.

Our systems will obey nonlinear differential equations of the form:

$$\frac{d}{dt}\vec{x} = \vec{f}(\vec{x}(t), \vec{u}(t)) + \vec{w}(t), \quad (3)$$

where \vec{f} is some non-linear function. Through linearization, we will see how to approximate general systems as locally linear systems, which will in turn allow us to apply all the techniques we have developed so far. We will then expand our theory to quadratic approximations as a straightforward generalization.

1 Taylor Expansion

First, we must recall how to approximate the simplest nonlinear functions — scalar functions of one variable — as linear functions. To do so, we will use the Taylor expansion we learned in calculus. This expansion tells us that for analytic functions $f(x)$ (functions that can be infinitely differentiated everywhere), we can write

$$f(x) = f(x^*) + f'(x^*)(x - x^*) + \frac{1}{2}f''(x^*)(x - x^*)^2 + \dots \quad (4)$$

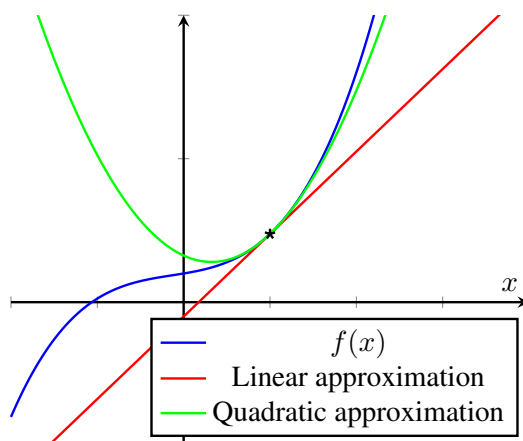
$$= f(x^*) + \frac{df}{dx}(x^*) \cdot (x - x^*) + \frac{1}{2} \frac{d^2f}{dx^2}(x^*) \cdot (x - x^*)^2 + \dots \quad (5)$$

near any particular point $x = x^*$, known as the *expansion point*.¹ By truncating this expansion after the first few terms, we can obtain a locally polynomial approximation for f of a desired degree d . For instance, we may take just the first two terms to linearly approximate $f(x)$ at x^* as

$$f(x) \approx f(x^*) + \frac{df}{dx}(x^*) \cdot (x - x^*). \quad (6)$$

It is further known that, when working in the “neighborhood” of $x = x^*$ (i.e. in a small interval around that point), the error of such a truncated approximation is bounded if the function is well behaved². The exact magnitude of this bound can be determined³ but is not very important for our purposes — the important part is that it is indeed bounded.

Before going further, let’s see what these approximations look like.



We have plotted a nonlinear function $f(x)$ and its linear approximation constructed around $x^* = 0.5$. Notice that though the approximation deviates greatly from $f(x)$ over the whole domain, in the close neighborhood

¹Note that the two expressions above are just two different notations for the same thing!

²For example, it has a bounded derivative beyond the level at which we truncate in the neighborhood of interest.

³Remember from your calculus courses — this is an application of the mean value theorem. You can take the maximum absolute value for the (presumably bounded) next derivative and use it in the expansion to get an explicit upper bound. The fact that the factorial $(k + 1)!$ in the denominator is very rapidly growing is what makes this bound get very small as we add terms in the expansion. But for us, the more important fact is that the bound includes $(x - x^*)^{k+1}$. This is what makes the bound get very small in a tight enough neighborhood around x^* . This tells us that we can use approximations that are not that high order (like linear) and still have a good local approximation. If the second derivative is bounded by 2, then in a neighborhood of distance 0.1 from x^* , the linear approximation will be good to a precision of 0.01.

of x^* , the error of the approximation is very small. Including one more term of our Taylor expansion, we obtain the quadratic approximation. Qualitatively, this approximation matches our function better in the neighborhood of our expansion point. In practice, we rarely go beyond the quadratic terms when trying to find a polynomial approximation of a function in a local region.

The ability to approximate is generally good because it lets us attack a problem that we don't know to solve by approximating it using another problem that we do know how to solve. But this only makes engineering sense if we have some way of dealing with the impact of the approximation error. Control provides a very useful example for this kind of thinking. If we approximate the f in a nonlinear differential equation by a linear function, why can we trust that everything will work out when we use this approximation to control? The answer is provided by the disturbance term $w(t)$ — our feedback control design for the linearized system should be able to reject bounded disturbances. In the neighborhood of the expansion point, our approximation differs from the true function by a small bounded amount. This approximation error can therefore be absorbed into the existing disturbance term — we just make the disturbance a little bigger. As long as the feedback controlled system stays within the region that the approximation is valid, we are fine. In effect, we pass the buck and let the feedback control deal with the approximation error. We'll see in a later note on classification how iteration provides another way to deal with the impact of approximation errors.

2 Multivariate Scalar-Valued Functions

Unfortunately, we can't directly apply the Taylor expansions that we learned in our calculus courses to solve our problem of linearizing a nonlinear system. Why is that? Well, so far we only know how to linearize scalar functions of one variable. But any nontrivial control system has at least *two* variables: the state and the control input! Even with just a scalar state and one-dimensional control input, our system would look like something of the form

$$\frac{d}{dt}x(t) = f(x(t), u(t)). \quad (7)$$

So we need to figure out how to deal with $f(x, u)$. Ideally, we'd be able to pick an expansion point (x^*, u^*) , around which we could linearize f as

$$f(x^* + \delta x, u^* + \delta u) \approx f(x^*, u^*) + a \cdot \delta x + b \cdot \delta u, \quad (8)$$

where a and b depend only on the function f and the x^* and u^* . It makes sense that a should somehow represent the “sensitivity” of f to variations in x , and b the sensitivity of f with respect to variations in u .⁴

To get a better understanding of how we can compute these two quantities, it will help to consider an example. Let $f(x, u) = x^3 u^2$. Then we see that, near an expansion point (x^*, u^*) , we can rewrite

$$f(x^* + \delta x, u^* + \delta u) = (x^* + \delta x)^3 (u^* + \delta u)^2 \quad (9)$$

$$= (x^{*3} + 3x^{*2} \cdot \delta x + 3x^* \cdot (\delta x)^2 + (\delta x)^3)(u^{*2} + 2u^* \cdot \delta u + (\delta u)^2) \quad (10)$$

$$= x^{*3} u^{*2} + 3x^{*2} u^{*2} \cdot \delta x + 2x^{*3} u^* \cdot \delta u + \dots \quad (11)$$

Why did we drop a lot of terms when making our approximation?

Imagine that we are very near this expansion point, so $|\delta x|, |\delta u| \ll x^*, u^*$. In other words, we can treat δx and δu as two comparably tiny quantities, that are both much smaller than x^* or u^* . Then which terms in

⁴Here note that δx and δu are one variable each, not the product of two variables – in general, δz usually means “a small quantity that's of the same type as z ”. This notation comes from calculus, where we reason with ϵ and δ in the following way: “for every ϵ there exists a δ such that...”. Usually δ is used in the argument of the function, in these contexts.

the above product (when expanded out fully) are most significant?

Well, any term involving δx or δu will immediately be made small, so the most significant term is $x^{*3}u^{*2}$. This term is just a constant, and isn't very interesting since it doesn't vary at all as we move around the neighborhood of our expansion point. So what can we get next? Well, we want as few δx and δu 's as possible in our terms. We've already got the only term with zero of them, so the next most significant terms will be those involving exactly one tiny quantity: $3x^{*2}u^{*2} \cdot \delta x$ and $2u^*x^{*3} \cdot \delta u$. Every other term will have products of tiny terms and we know that products of tiny numbers are even tinier numbers. If you will recall, this is the same type of reasoning that you saw in your calculus course when you first encountered derivatives.

Thus, using only these most significant terms, we obtain the approximation

$$x^3u^2 \approx x^{*3}u^{*2} + 3x^{*2}u^{*2} \cdot \delta x + 2x^{*3}u^* \cdot \delta u. \quad (12)$$

The $x^{*3}u^{*2}$ term corresponds to the $f(x^*, u^*)$ term in our desired linearization. In a similar manner, $3x^{*2}u^{*2}$ corresponds to the coefficient a and $2u^*x^{*3}$ to the coefficient b .

How could we have obtained these coefficients directly from f ?

2.1 Partial Derivatives

The key insight is the following. Imagine that u was in fact a constant, so our function was just the scalar valued function of a single scalar variable $f(x) = x^3u^2$. Then what is its derivative? Well, since u^2 is being treated as a constant, its derivative is $3x^2u^2$. Similarly, by treating x as a constant and writing our nonlinear function as $f(u) = x^3u^2$, we can differentiate to obtain $2ux^3$. These quantities are known as the *partial derivatives* of f , and are denoted as $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial u}$ in analogy with the single-variable notation $\frac{df}{dx}$, or alternatively the more compact notation f_x, f_u in analogy with the single-variable notation f' .

Now when we evaluate the partial derivatives at our expansion point (x^*, u^*) , we exactly get the terms from above. In other words, we obtain the desired coefficients by treating all the variables except the one of interest as a constant, and then differentiating f with respect to the remaining variable. Thus, we get

$$a = \frac{\partial f}{\partial x}(x^*, u^*) = f_x(x^*, u^*) \quad \text{and} \quad b = \frac{\partial f}{\partial u}(x^*, u^*) = f_u(x^*, u^*). \quad (13)$$

Then we can linearize any well-behaved scalar-valued function of 2 variables about an expansion point (x^*, u^*) as

$$f(x, u) \approx f(x^*, u^*) + \frac{\partial f}{\partial x}(x^*, u^*) \cdot (x - x^*) + \frac{\partial f}{\partial u}(x^*, u^*) \cdot (u - u^*) \quad (14)$$

If we had more than 2 variables as the input, and instead a vector of variables, we can generalize this to get the linearization of $f(\vec{x})$ around an expansion point \vec{x}^* :

$$f(\vec{x}) \approx f(\vec{x}^*) + \sum_{i=1}^n \frac{\partial f}{\partial x_i}(\vec{x}^*) \cdot (x_i - x_i^*) \quad (15)$$

3 Vector-Valued Functions

With partial derivatives in hand, we can return to our original problem of linearizing a vector-valued function $\vec{f}(\vec{x}, \vec{u})$.

First, we should ask ourselves what the desired form of the linearization could be. Recalling our experience with linear systems, we'd probably like our linearization near an expansion point (\vec{x}^*, \vec{u}^*) to look something like

$$\vec{f}(\vec{x}, \vec{u}) \approx \vec{f}(\vec{x}^*, \vec{u}^*) + A(\vec{x} - \vec{x}^*) + B(\vec{u} - \vec{u}^*), \quad (16)$$

where A and B are some matrices that presumably depend on \vec{f} and \vec{x}^* and \vec{u}^* .

Let our state \vec{x} be n -dimensional and our control \vec{u} be k -dimensional. Since $\frac{d}{dt}\vec{x} = \vec{f}(\vec{x}, \vec{u})$, the output of \vec{f} must also be n -dimensional. Thus, by the rules of matrix multiplication, we know that A must be an $n \times n$ matrix and B an $n \times k$ matrix, for the dimensions to all work out. But what are the contents of these two matrices? To determine this, we will look at our vectors elementwise, expressing the function we wish to linearize as

$$\vec{f}(\vec{x}, \vec{u}) = \begin{bmatrix} f_1(\vec{x}, \vec{u}) \\ f_2(\vec{x}, \vec{u}) \\ \vdots \\ f_n(\vec{x}, \vec{u}) \end{bmatrix}, \quad (17)$$

where the f_i are each scalar-valued functions. Let's now consider a single one of these f_i , and try to linearize it about an expansion point. Note that though we write each f_i as functions of the vectors \vec{x} and \vec{u} , they are really simply functions of the $n + k$ individual scalars $x_1, \dots, x_n, u_1, \dots, u_k$ that form the components of the two vector-valued arguments. Thus, we can use the linearization techniques from the earlier section to approximate

$$f_i(\vec{x}, \vec{u}) \approx f_i(\vec{x}^*, \vec{u}^*) + \sum_{j=1}^n \frac{\partial f_i}{\partial x_j}(\vec{x}^*, \vec{u}^*) \cdot (x_j - x_j^*) + \sum_{j=1}^k \frac{\partial f_i}{\partial u_j}(\vec{x}^*, \vec{u}^*) \cdot (u_j - u_j^*). \quad (18)$$

We can rearrange the above linearization as the matrix multiplication

$$f_i(\vec{x}, \vec{u}) \approx f_i(\vec{x}^*, \vec{u}^*) \quad (19)$$

$$+ \begin{bmatrix} \frac{\partial f_i}{\partial x_1}(\vec{x}^*, \vec{u}^*) & \dots & \frac{\partial f_i}{\partial x_n}(\vec{x}^*, \vec{u}^*) \end{bmatrix} \cdot (\vec{x} - \vec{x}^*) \quad (20)$$

$$+ \begin{bmatrix} \frac{\partial f_i}{\partial u_1}(\vec{x}^*, \vec{u}^*) & \dots & \frac{\partial f_i}{\partial u_k}(\vec{x}^*, \vec{u}^*) \end{bmatrix} \cdot (\vec{u} - \vec{u}^*). \quad (21)$$

The two row vectors above can be thought of as the derivatives⁵ of $f_i(\vec{x}, \vec{u})$ with respect to \vec{x} and \vec{u} , so we denote these rows by $\frac{\partial f_i}{\partial \vec{x}}$ and $\frac{\partial f_i}{\partial \vec{u}}$. These are row-vector-valued functions of column-vectors \vec{x}^*, \vec{u}^* , in the same way that a derivative of a scalar-valued function of scalars is itself a scalar-valued function. So,

$$f_i(\vec{x}, \vec{u}) \approx f_i(\vec{x}^*, \vec{u}^*) + \frac{\partial f_i}{\partial \vec{x}}(\vec{x}^*, \vec{u}^*) \cdot (\vec{x} - \vec{x}^*) + \frac{\partial f_i}{\partial \vec{u}}(\vec{x}^*, \vec{u}^*) \cdot (\vec{u} - \vec{u}^*) \quad (22)$$

This was the linearization for just one element of \vec{f} , so we can now stack the above linearization for all f_i to obtain a linearization for the original vector-valued $\vec{f}(\vec{x}, \vec{u})$:

$$\vec{f}(\vec{x}, \vec{u}) \approx \vec{f}(\vec{x}^*, \vec{u}^*) \quad (23)$$

⁵In general, the derivative of a scalar valued function of a vector has to be a row vector. This is because it needs to act on an change in a column vector and return a scalar that represents the change to the function. Rows multiply columns to give scalars. For those of you who might have taken a course in multivariable calculus, it is vital that you not confuse the gradient ∇f with the derivative. The gradient is a column vector and is the transpose of the derivative of the function f .

$$+ \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(\vec{x}^*, \vec{u}^*) & \cdots & \frac{\partial f_1}{\partial x_n}(\vec{x}^*, \vec{u}^*) \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1}(\vec{x}^*, \vec{u}^*) & \cdots & \frac{\partial f_n}{\partial x_n}(\vec{x}^*, \vec{u}^*) \end{bmatrix} \cdot (\vec{x} - \vec{x}^*) \quad (24)$$

$$+ \begin{bmatrix} \frac{\partial f_1}{\partial u_1}(\vec{x}^*, \vec{u}^*) & \cdots & \frac{\partial f_1}{\partial u_k}(\vec{x}^*, \vec{u}^*) \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial u_1}(\vec{x}^*, \vec{u}^*) & \cdots & \frac{\partial f_n}{\partial u_k}(\vec{x}^*, \vec{u}^*) \end{bmatrix} \cdot (\vec{u} - \vec{u}^*), \quad (25)$$

We call these two matrices of partial derivatives the *Jacobian(s)* of \vec{f} with respect to \vec{x} and \vec{u} .

Formally, the Jacobian is the partial derivative of a vector-valued function \vec{f} with respect to a vector input \vec{x} . Thus, in analogy with our previous notations, we may write the Jacobian as $\frac{\partial \vec{f}}{\partial \vec{x}}$. Other notations include, for example, $J_{\vec{x}} \vec{f}$. These are matrix-valued functions, just as the derivative of a scalar function by a scalar variable is a scalar function.

Here they are evaluated at (\vec{x}^*, \vec{u}^*) . This gives us

$$\vec{f}(\vec{x}, \vec{u}) \approx \vec{f}(\vec{x}^*, \vec{u}^*) + \frac{\partial \vec{f}}{\partial \vec{x}}(\vec{x}^*, \vec{u}^*) \cdot (\vec{x} - \vec{x}^*) + \frac{\partial \vec{f}}{\partial \vec{u}}(\vec{x}^*, \vec{u}^*) \cdot (\vec{u} - \vec{u}^*) \quad (26)$$

so we have solved for $A = \frac{\partial \vec{f}}{\partial \vec{x}}(\vec{x}^*, \vec{u}^*)$ and $B = \frac{\partial \vec{f}}{\partial \vec{u}}(\vec{x}^*, \vec{u}^*)$ in our desired linear model.

Note that this formula is exactly how our other linearization formulas look.

One great development is that modern software frameworks like PyTorch and TensorFlow have automatic support for computing derivatives of functions that are expressed naturally using code. This means that all these ideas around linearization are now quite easy to implement — you don't have to compute all the derivatives by hand in practice.

4 Equilibrium Points

Now, let's see how this linearization plugs back into our original differential equation. We see that

$$\frac{d}{dt} \vec{x}(t) = f(\vec{x}, \vec{u}) = \vec{f}(\vec{x}^*, \vec{u}^*) + \frac{\partial \vec{f}}{\partial \vec{x}}(\vec{x}^*, \vec{u}^*) \cdot (\vec{x}(t) - \vec{x}^*) + \frac{\partial \vec{f}}{\partial \vec{u}}(\vec{x}^*, \vec{u}^*) \cdot (\vec{u}(t) - \vec{u}^*) + \vec{w}(t). \quad (27)$$

Notice that even though we are using an approximation, we have used an equals sign $=$ above. This is because we have folded the approximation error into the disturbance term $\vec{w}(t)$.

So have we gotten this into a form that we know how to deal with? The first apparent problem is that we have $\vec{x}(t) - \vec{x}^*$ on the right, but $\vec{x}(t)$ on the left. To fix this, we will define new state variables $\Delta \vec{x}(t) = \vec{x}(t) - \vec{x}^*$ and new control inputs $\Delta \vec{u}(t) = \vec{u}(t) - \vec{u}^*$.

By rewriting $\vec{x}(t) = \Delta \vec{x}(t) + \vec{x}^*$, we have

$$\frac{d}{dt} \vec{x}(t) = \frac{d}{dt} (\vec{x}^* + \Delta \vec{x}(t)) = \frac{d}{dt} \vec{x}^* + \frac{d}{dt} \Delta \vec{x}(t). \quad (28)$$

Since \vec{x}^* is a constant for all time, we know that $\frac{d}{dt}\vec{x}^* = \vec{0}$. Hence

$$\frac{d}{dt}\vec{x}(t) = \frac{d}{dt}\Delta\vec{x}(t). \quad (29)$$

Expanding out the derivative, we obtain

$$\frac{d}{dt}\Delta\vec{x}(t) = \frac{d}{dt}\vec{x}(t) \quad (30)$$

$$= \vec{f}(\vec{x}^*, \vec{u}^*) + \frac{\partial \vec{f}}{\partial \vec{x}}(\vec{x}^*, \vec{u}^*) \cdot (\vec{x}(t) - \vec{x}^*) + \frac{\partial \vec{f}}{\partial \vec{u}}(\vec{x}^*, \vec{u}^*) \cdot (\vec{u}(t) - \vec{u}^*) + \vec{w}(t) \quad (31)$$

$$= \vec{f}(\vec{x}^*, \vec{u}^*) + \frac{\partial \vec{f}}{\partial \vec{x}}(\vec{x}^*, \vec{u}^*) \cdot \Delta\vec{x}(t) + \frac{\partial \vec{f}}{\partial \vec{u}}(\vec{x}^*, \vec{u}^*) \cdot \Delta\vec{u}(t) + \vec{w}(t). \quad (32)$$

Then the final term to remove to get a linear model is the $\vec{f}(\vec{x}^*, \vec{u}^*)$ on the right side of the equation.⁶ Notice that this term depends only on the choice of expansion point. Thus, when we linearize a system, it is our job to pick expansion points such that

$$\vec{f}(\vec{x}^*, \vec{u}^*) = \vec{0}. \quad (33)$$

Such points (\vec{x}^*, \vec{u}^*) are called *equilibrium points* of \vec{f} . When we linearize a dynamical system around such an equilibrium, the approximation is now a linear system:

$$\frac{d}{dt}\Delta\vec{x}(t) = \frac{\partial \vec{f}}{\partial \vec{x}}(\vec{x}^*, \vec{u}^*) \cdot \Delta\vec{x}(t) + \frac{\partial \vec{f}}{\partial \vec{u}}(\vec{x}^*, \vec{u}^*) \cdot \Delta\vec{u}(t) + \vec{w}(t) \quad (34)$$

$$= A \cdot \Delta\vec{x}(t) + B \cdot \Delta\vec{u}(t) + \vec{w}(t) \quad (35)$$

where now $A = \frac{\partial \vec{f}}{\partial \vec{x}}(\vec{x}^*, \vec{u}^*)$ and $B = \frac{\partial \vec{f}}{\partial \vec{u}}(\vec{x}^*, \vec{u}^*)$.

The system's stability can be checked by looking at the eigenvalues of A and the local controllability can be checked by looking at A, B together. At this point, if the system is controllable, we can use feedback control to place the eigenvalues of $A + BK$ where we want and thereby get a feedback control law that will stabilize the linearized system about the chosen equilibrium point. Note that this feedback control will drive $\Delta\vec{x} \rightarrow 0$ which means it drives $\vec{x} \rightarrow \vec{x}^*$. In this way, we can think of \vec{x}^* as our desired value of $\vec{x}(t)$. It is also important to remember that the actual control applied will be $\vec{u}(t) = \vec{u}^* + K(\vec{x}(t) - \vec{x}^*)$ — because the K has been calculated for the linearized system and we need to apply a control in the original nonlinear system.

In general, there are many potential solutions for an equilibrium point. Which one you choose depends on the context or design goal. For example, one might want to find a feedback control law that stabilizes a nonlinear system about a particular state \vec{x}^* . In that case, you will want to solve for the corresponding \vec{u}^* that satisfies (33) with that particular state \vec{x}^* . It can be thought of as the control that would keep the original nonlinear system held at that particular point in the absence of all disturbances.

⁶Note that linear functions are of the form $\vec{g}(\Delta\vec{x}, \Delta\vec{u}) = A(\Delta\vec{x}) + B(\Delta\vec{u}) = \begin{bmatrix} A & B \end{bmatrix} \begin{bmatrix} \Delta\vec{x} \\ \Delta\vec{u} \end{bmatrix}$, *not* of the form $A(\Delta\vec{x}) + B(\Delta\vec{u}) + \vec{c} = \begin{bmatrix} A & B \end{bmatrix} \begin{bmatrix} \Delta\vec{x} \\ \Delta\vec{u} \end{bmatrix} + \vec{c}$. The latter functions are called *affine* functions, and do *not* obey the linear property $\vec{g}(a_1\vec{y}_1 + a_2\vec{y}_2) = a_1\vec{g}(\vec{y}_1) + a_2\vec{g}(\vec{y}_2)$, for $\vec{y} = \begin{bmatrix} \Delta\vec{x} \\ \Delta\vec{u} \end{bmatrix}$. We have lots of theory for linear \vec{g} , but not affine \vec{g} . Modifying the system to “fold” the affine offset \vec{c} into A or B , thus making the system linear, would change properties of the system, such as controllability or stability. We like dealing with purely linear systems so as to not deal with that.

The ability to linearize is what makes linear control practical since many (if not most) real-world systems that we encounter are fundamentally nonlinear. The linearized system of differential equations can also be discretized to give rise to a discrete-time linear control system. Here, the critical choice is that we must choose a discretization time interval length that is small enough so that the approximation remains valid. The system must not move outside of the region of approximation validity (i.e. the neighborhood where the approximation error induced by linearization is within the desired bound) during the interval between taking samples. As long as we sample fast enough, we are free to design our control laws in discrete-time.

But what if we start in discrete-time? In general, it is not valid to linearize in discrete-time, since the sampling interval is fixed in discrete time and is probably larger than the interval in which our linearization is a useful approximation. Thus we generally need to start from continuous-time systems and linearize to discrete-time systems, or else go from a discrete-time system to a continuous-time system by taking the limit as the sampling interval goes to 0, linearize, and convert back by discretization.

5 Finding Equilibrium Points

To solve the nonlinear system of equations (33) and find equilibrium points, we can use graphical⁷ methods. We can also use numerical methods to solve the system by computer or in closed form.

So, how do you solve nonlinear systems of equations like (33)? This can be challenging at times but fortunately software packages exist for you to use. For example, SciPy has `scipy.optimize.fsolve` for you. Alternatively, you can use the Gauss-Newton method.

Specifically, in the homework, you will see a problem that shows you how iterative methods can be used to solve a general $\vec{g}(\vec{\theta}) = \vec{0}$. These begin with a guess $\vec{\theta}^*$ for a solution, and then proceed iteratively. One intuitive way (that you will explore in the homework) is to linearize the function \vec{g} around $\vec{\theta}$. This gives rise to a linear system of equations $\vec{g}(\vec{\theta}^*) + \vec{g}'(\vec{\theta}^*)(\vec{\theta}^* + \vec{\delta\theta}) = \vec{0}$. Here, $\vec{g}'(\vec{\theta}^*)$ denotes the derivative of \vec{g} with respect to $\vec{\theta}$, i.e. the matrix $\frac{\partial \vec{g}}{\partial \vec{\theta}}$. We can solve this linear system of equations using what you know for a candidate update $\vec{\delta\theta}$. At this point, we can update our guess and iterate. We could update our new guess $\vec{\theta}_{i+1}$ to be $\vec{\theta}_i + \vec{\delta\theta}$ itself. Or we could be more conservative (under the belief that the linearization was only valid in a small neighborhood) and just take a little step in that direction to $\vec{\theta}_i + \eta \vec{\delta\theta}$ where η is some small constant that is less than 1. In general, by iterating such an algorithm with a small enough step-size η , we find our way to a root. This is also the approach that we will take when we apply linearization ideas to the problem of classification in the next note.

6 Tracking Trajectories

You've seen in the homework that we can do linear feedback control to keep a linear system stably following a desired chosen open-loop trajectory. Looking at the equations earlier, we can similarly consider what happens when \vec{x}^* is no longer constant, and so we want $\vec{x}(t)$ to follow some desired trajectory $\vec{x}^*(t)$. Then to get the same linear system as (34), we will want points such that

$$\frac{d}{dt} \vec{x}^*(t) = \vec{f}(\vec{x}^*(t), \vec{u}^*(t)). \quad (36)$$

How do we find such nominal solutions? That is a bit out of the scope of this course but it turns out that the

⁷By graphical, we literally mean methods that involve drawing the graph of the function and seeing where it crosses zero or related intuitive approaches.

same kinds of iterative numerical approaches that work for finding equilibrium points can be leveraged here as well. We linearize locally, solve, and then relinearize and resolve repeatedly.

However, there is one more challenge that occurs when we use linearization around a nominal trajectory — the resulting A and B matrices are now generically functions of time. So does that mean we are doomed? Not at all! While one has to be careful, often we can succeed by taking the same modeling philosophy that we have taken so far. If we assume that we are sampling the system fast enough, we can hope that the A and B matrices don't change too much from one step to the next. As a result, we can approximate the A matrices as being locally constant over many samples and just fold the approximation error into the disturbances. Here, we need to be a bit more careful because any error in A (due to time-variation) multiplies the deviation of the state and we need the product to stay small. If the state is small and the Frobenius norm of the change in A is small, the product will also be small. Similarly for B — but here we must ensure that the control that we apply times the change in B stays small. This means that the feedback controls should be gentle enough. You will learn how to navigate the resulting tension in more advanced control courses.

Finally, what about planning a nominal trajectory? Here, it turns out that the minimum-norm approach that you have learned can be iterated to get plans. The reason is that there was nothing in the approach that required the A and B matrices to stay constant. So it is possible to guess a set of controls to get to a destination, linearize, find a minimum norm solution, and then iterate this process to get a plan. Not all starting guesses will work, and so sometimes, more brute force guessing or smarter heuristics have to be used to speed this process up. This is also a topic for more advanced robotics courses.

7 Second-Order Approximation

There is one final loose end to tie up; namely, what happens if we take more than just the linear term in our truncated Taylor series approximation? It turns out that going more than two terms, even for scalar-valued functions of vector-arguments, gets very messy. We would need more advanced versions of the tools we have already covered, along with a deeper understanding of abstract linear algebra and the geometry of vector spaces, to efficiently tackle such problems, which form the backbone of so-called *tensor calculus* (tensors being essentially n -dimensional collections of numbers, in the sense that vectors are when $n = 1$ and matrices are when $n = 2$). However, there is one setting where going beyond linearization makes sense to us right now: when we are approximating a scalar function of a vector input by using up to two derivatives.

Namely, suppose $f(\vec{x})$ is a scalar-valued function of a vector $\vec{x} \in \mathbb{R}^n$. We would like to approximate f around a point \vec{x}^* using the first and second derivatives of f .

The first derivative of f , $\frac{\partial f}{\partial \vec{x}}$, is something we already know how to calculate from earlier in this note. But what about the second derivative? As before, let us retreat to the single-variable case and then do a concrete example to gain intuition, redoing our previous derivation with the quadratic term in mind.

First, let us consider the scalar case. Suppose f is a scalar function of a scalar variable x . Then Taylor's theorem tells us

$$f(x) = f(x^*) + \frac{df}{dx}(x^*) \cdot (x - x^*) + \frac{1}{2} \frac{d^2f}{dx^2}(x^*) \cdot (x - x^*)^2 + \dots \quad (37)$$

$$\approx f(x^*) + \frac{df}{dx}(x^*) \cdot (x - x^*) + \frac{1}{2} \frac{d^2f}{dx^2}(x^*) \cdot (x - x^*)^2. \quad (38)$$

Note that the second derivative has a coefficient of $\frac{1}{2}$ in front of it.

Now we do an example. Let us suppose we have the system

$$\frac{d}{dt}x(t) = f(x(t), u(t)) \quad (39)$$

where $f(x, u) = x^3u^2$ as before. Again, we perturb both $x = x^* + \delta x$ and $u = u^* + \delta u$. We write

$$f(x^* + \delta x, u^* + \delta u) = (x^* + \delta x)^3(u^* + \delta u)^2 \quad (40)$$

$$\approx x^{*3}u^{*2} + 3x^{*2}u^{*2} \cdot \delta x + 2x^{*3}u^* \cdot \delta u \quad (41)$$

$$+ 3x^*u^{*2} \cdot (\delta x)^2 + x^{*3} \cdot (\delta u)^2 + 6x^{*2}u^* \cdot (\delta x)(\delta u). \quad (42)$$

Here, in order to carry through a quadratic approximation in our variables δx , δu , we kept all the linear and quadratic terms in the expansion, and not just the linear terms like we did for the linear approximation. Here, note that quadratic doesn't necessarily mean that *one* variable has to have a power of two – just that the sum of powers of the variables for that term is equal to two, as seen in the last term where the powers on δx and δu are both 1. This definition is consistent with other contexts in which we find multivariate polynomials.

The key insight to continue from here is that we can extend our concept of partial derivatives to get second-order partial derivatives, and indeed n^{th} order partial derivatives. Indeed, one can write

$$\frac{\partial^2 f}{\partial x^2} = \frac{\partial}{\partial x} \left(\frac{\partial f}{\partial x} \right) \quad \frac{\partial^2 f}{\partial u^2} = \frac{\partial}{\partial u} \left(\frac{\partial f}{\partial u} \right) \quad \frac{\partial^2 f}{\partial x \partial u} = \frac{\partial}{\partial u} \left(\frac{\partial f}{\partial x} \right), \quad (43)$$

and so on. It is a fact of multivariable calculus that under mild conditions, such as the second derivatives being continuous at the point of evaluation, one can interchange the order of taking partial derivatives, i.e., $\frac{\partial^2 f}{\partial x \partial u} = \frac{\partial^2 f}{\partial u \partial x}$. The precise result is called **Clairaut's theorem**.

Rewriting our earlier approximation in terms of the first and second partial derivatives, we can write

$$f(x^* + \delta x, u^* + \delta u) \approx f(x^*, u^*) + \frac{\partial f}{\partial x}(x^*, u^*) \cdot \delta x + \frac{\partial f}{\partial u}(x^*, u^*) \cdot \delta u \quad (44)$$

$$+ \frac{1}{2} \frac{\partial^2 f}{\partial x^2}(x^*, u^*) \cdot (\delta x)^2 + \frac{1}{2} \frac{\partial^2 f}{\partial u^2}(x^*, u^*) \cdot (\delta u)^2 + \frac{\partial^2 f}{\partial x \partial u}(x^*, u^*) \cdot (\delta x)(\delta u). \quad (45)$$

This is very close to symmetric in all the partial derivatives, with one exception: the so-called “mixed partial” $\frac{\partial^2 f}{\partial x \partial u}$ has coefficient 1 while all other partial derivatives have coefficient $\frac{1}{2}$ (similarly to the single-variable case), and there is no partial $\frac{\partial^2 f}{\partial u \partial x}$.

We can change both of those at the same time, by using Clairaut's theorem to interchange the partial derivatives. Since $\frac{\partial^2 f}{\partial x \partial u} = \frac{\partial^2 f}{\partial u \partial x}$, we can write $\frac{\partial^2 f}{\partial x \partial u} = \frac{1}{2} \frac{\partial^2 f}{\partial x \partial u} + \frac{1}{2} \frac{\partial^2 f}{\partial u \partial x}$, and obtain the much nicer symmetric expression

$$f(x^* + \delta x, u^* + \delta u) \approx f(x^*, u^*) + \frac{\partial f}{\partial x}(x^*, u^*) \cdot \delta x + \frac{\partial f}{\partial u}(x^*, u^*) \cdot \delta u \quad (46)$$

$$+ \frac{1}{2} \frac{\partial^2 f}{\partial x^2}(x^*, u^*) \cdot (\delta x)^2 + \frac{1}{2} \frac{\partial^2 f}{\partial u^2}(x^*, u^*) \cdot (\delta u)^2 \quad (47)$$

$$+ \frac{1}{2} \frac{\partial^2 f}{\partial x \partial u}(x^*, u^*) \cdot (\delta x)(\delta u) + \frac{1}{2} \frac{\partial^2 f}{\partial u \partial x}(x^*, u^*) \cdot (\delta u)(\delta x). \quad (48)$$

In fact this formula holds in general if we have a general analytic scalar-valued function f in two variables.

Writing it in a similar fashion to the previous equations, we have

$$f(x, u) \approx f(x^*, u^*) + \frac{\partial f}{\partial x}(x^*, u^*) \cdot (x - x^*) + \frac{\partial f}{\partial u}(x^*, u^*) \cdot (u - u^*) \quad (49)$$

$$+ \frac{1}{2} \frac{\partial^2 f}{\partial x^2}(x^*, u^*) \cdot (x - x^*)^2 + \frac{1}{2} \frac{\partial^2 f}{\partial u^2}(x^*, u^*) \cdot (u - u^*)^2 \quad (50)$$

$$+ \frac{1}{2} \frac{\partial^2 f}{\partial x \partial u}(x^*, u^*) \cdot (x - x^*)(u - u^*) + \frac{1}{2} \frac{\partial^2 f}{\partial u \partial x}(x^*, u^*) \cdot (u - u^*)(x - x^*). \quad (51)$$

It even generalizes if we have a general analytic scalar-valued function f of a vector $\vec{x} \in \mathbb{R}^n$:

$$f(\vec{x}) = f(\vec{x}^*) + \sum_{i=1}^n \frac{\partial f}{\partial x_i}(\vec{x}^*) \cdot (x_i - x_i^*) + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 f}{\partial x_i \partial x_j}(\vec{x}^*) \cdot (x_i - x_i^*)(x_j - x_j^*) \quad (52)$$

Remember that when we derived linear approximations, we developed a nice vector-based expression for the terms involving the first derivatives:

$$\sum_{i=1}^n \frac{\partial f}{\partial x_i}(\vec{x}^*) \cdot (x_i - x_i^*) = \frac{\partial f}{\partial \vec{x}}(\vec{x}^*) \cdot (\vec{x} - \vec{x}^*). \quad (53)$$

Now we want to come up with a nice, idiomatic expression for the second derivatives:

$$\sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 f}{\partial x_i \partial x_j}(\vec{x}^*) \cdot (x_i - x_i^*)(x_j - x_j^*) =? \quad (54)$$

To do this in the linear case, we defined a new object – the row vector of first partial derivatives $\frac{\partial f}{\partial \vec{x}}$, such that $\frac{\partial f}{\partial \vec{x}}(\vec{x}^*) \cdot (\vec{x} - \vec{x}^*)$ gave us a scalar.

Now we would like to get a scalar from some product of $\vec{x} - \vec{x}^*$ and some function of the second derivatives. One guess we can make is to go entirely based off object types, in the following manner.

One sees that a matrix times a column vector is a column vector, and that a row vector times a column vector is a scalar. So a row vector times a matrix times a column vector should be a scalar. Let us naively guess that the matrix should contain all our second partial derivatives, and that our column and row vectors should be $\vec{x} - \vec{x}^*$, respectively, since they seem like the only candidate quantities we can work with.

This actually turns out to work. More specifically, define the *Hessian* of f to be the matrix containing all its second partial derivatives:

$$\frac{\partial^2 f}{\partial \vec{x}^2} = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} \quad (55)$$

We can also use some other notations like $H_{\vec{x}}f$ that you might see in discussion. One thing to note here is that the *Jacobian*, i.e., $\frac{\partial \vec{f}}{\partial \vec{x}}$, the first derivatives of a vector-valued function of a vector, is different from the *Hessian*, i.e., $\frac{\partial^2 f}{\partial \vec{x}^2}$, the second derivative of a scalar-valued function of a vector. But they are connected in the following way: if f is a scalar-valued function of a vector, and $\vec{g}(\vec{x}) = \left(\frac{\partial f}{\partial \vec{x}}(\vec{x})\right)^T$ (so that g is a column vector), then $\frac{\partial \vec{g}}{\partial \vec{x}} = \frac{\partial^2 f}{\partial \vec{x}^2}$. This is analogous to the way one takes second derivatives by taking the derivative

of the first derivative, and so on.

Back to our guess. We plug in and try to compute:

$$(\vec{x} - \vec{x}^*)^\top \left(\frac{\partial^2 f}{\partial \vec{x}^2}(\vec{x}^*) \right) (\vec{x} - \vec{x}^*) = \sum_{i=1}^n \sum_{j=1}^n \left(\frac{\partial^2 f}{\partial \vec{x}^2}(\vec{x}^*) \right)_{ij} (\vec{x} - \vec{x}^*)_i (\vec{x} - \vec{x}^*)_j \quad (56)$$

$$= \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 f}{\partial x_i \partial x_j}(\vec{x}^*) (x_i - x_i^*) (x_j - x_j^*). \quad (57)$$

This turns out to be exactly the quantity we need! So we have found out how to simplify our original expression for the second derivatives:

$$\sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 f}{\partial x_i \partial x_j}(\vec{x}^*) \cdot (x_i - x_i^*) (x_j - x_j^*) = (\vec{x} - \vec{x}^*)^\top \left(\frac{\partial^2 f}{\partial \vec{x}^2}(\vec{x}^*) \right) (\vec{x} - \vec{x}^*). \quad (58)$$

Through this, we can simplify our expression for quadratic approximation:

$$f(\vec{x}) \approx f(\vec{x}^*) + \frac{\partial f}{\partial \vec{x}}(\vec{x}^*) (\vec{x} - \vec{x}^*) + \frac{1}{2} (\vec{x} - \vec{x}^*)^\top \left(\frac{\partial^2 f}{\partial \vec{x}^2}(\vec{x}^*) \right) (\vec{x} - \vec{x}^*). \quad (59)$$

Quadratic approximation will be useful to us in the next few sections on classification and machine learning, where sometimes linearization is not helpful to us because our functions have some kind of “essential non-linearity” that must be taken into account during analysis. However, linearization works just fine for many applications, such as robotics, control, circuits, and so on.

Contributors:

- Rahul Arya.
- Anant Sahai.
- Ashwin Vangipuram.
- Druv Pai.