

3 Growth and Capital Accumulation: The Solow Model

Chapter 3 Contents

A. Topics and Tools	2
B. Introduction to Growth Theory	2
<i>Efficient output for a single-person economy</i>	<i>3</i>
<i>Efficient output with multiple people and firms</i>	<i>4</i>
<i>Natural output vs. efficient output</i>	<i>5</i>
<i>Growth theory and behavior of natural output</i>	<i>7</i>
C. Growth in Continuous Time: Logarithmic and Exponential Functions.....	10
<i>Continuous-time vs. discrete-time models</i>	<i>10</i>
<i>Growth in discrete and continuous time</i>	<i>11</i>
<i>Exponentials, logs, and continuous growth</i>	<i>13</i>
D. Review of Some Basic Calculus Tools	16
<i>Derivatives of powers, sums, products, and quotients</i>	<i>19</i>
<i>Derivatives and maximization</i>	<i>20</i>
<i>Other rules of differentiation</i>	<i>21</i>
E. Calculus Applications in Macroeconomics	22
<i>An application: time derivatives</i>	<i>22</i>
<i>Growth rates of products, quotients, and powers</i>	<i>24</i>
<i>Multivariate functions and partial derivatives</i>	<i>25</i>
<i>Total differentials</i>	<i>25</i>
<i>Multivariate maximization and minimization</i>	<i>26</i>
F. Understanding Romer's Chapter 1	26
<i>Manipulating the production function</i>	<i>26</i>
<i>The Cobb-Douglas production function</i>	<i>28</i>
<i>The nature of growth equilibrium</i>	<i>30</i>
<i>Basic dynamic analysis of k</i>	<i>31</i>
<i>Using Taylor series to approximate the speed of convergence</i>	<i>32</i>
<i>Growth models and the environment</i>	<i>34</i>
G. Suggestions for Further Reading	35
<i>Expositions of the Solow model</i>	<i>35</i>

A. Topics and Tools

Romer's Chapter 1, covering the Solow growth model and related theories, presents several challenges that may be new to macroeconomics students. First and foremost, it may be the first time that you have used calculus and related mathematical methods to analyze economic models. Basic calculus concepts are reviewed in Section C of this chapter. If your calculus is shaky or rusty, this section may help, but you may also want to pursue remedial tutorial work through the Quantitative Skills Center.

The second novelty of this chapter is the concept of a dynamic equilibrium growth path rather than a static point of equilibrium. We construct the Solow model in continuous time, which enables us to describe rates of change in terms of "time derivatives" and to make extensive use of the logarithmic and exponential functions to model the movements of variables over time. These methods will be very familiar to you if you have taken a course covering differential equations, but otherwise might be quite new. Section B introduces you to some of the concepts of continuous-time modeling that we will use extensively.

The central element of growth theory is the feedback from current economic conditions to investment in new capital to increases in productive capacity that influence future economic conditions. This seems to suggest the possibility of self-sustaining growth through capital deepening. The Solow growth model examines a simple proposition: Can an economy that saves and invests a constant share of its income grow forever? The answer is no. With a constant saving rate, such an economy will converge to an equilibrium capital-labor ratio, after which any growth that occurs must originate in a growing labor force or improving technology.

B. Introduction to Growth Theory

Chapter 1 introduced the idea of growth and cyclical fluctuations in real GDP. In this chapter, we begin the analysis of the long-run growth path of the economy. This analysis largely ignores short-run fluctuations or deviations from the trend growth path in order to focus on the slope (growth rate) and level of the path.

Growth theory studies the evolution over time in the *natural level of output* in the economy—the amount of output that would be produced if the macroeconomy were in balance. To help clarify the nature of natural output, we shall start with the simplest possible economy and work our way toward greater complexity of economic interaction. We shall also consider the distinction between what we call the natural

level of output and *efficient output*. Efficient output measures the amount that an economy would produce if resources were fully and efficiently employed. Natural output measures the amount (usually less than efficient) that is produced when we account for microeconomic inefficiencies. In common microeconomic terms, efficient output corresponds to the production-possibilities frontier, whereas natural output recognizes that inefficiencies due to monopoly power, tax distortions, search costs, and other market imperfections will inevitably cause the economy to produce inside the PPF, even at “full employment.”

Efficient output for a single-person economy

Because it is so central to growth theory, it is worth considering the concept of natural output in more detail. We will start with the simplest possible economy: one person (Jane) living in isolation on an island and surviving by eating coconuts. Jane sleeps whenever the sun is down and eats for two hours each day. She has only two activities that occupy her remaining waking hours: picking coconuts, which she does not enjoy but she must pick them in order to eat, and playing Frisbee golf, which she likes. Coconuts are assumed to be perfectly perishable; any that are not eaten at the end of the day are stolen and consumed by monkeys.

Jane has a daily production function for coconuts with the single input factor being hours L she spends picking them. Because there are some already fallen coconuts and some lying low to the ground, her first hours of picking are very productive. As she spends more hours working, her marginal productivity (coconuts per additional hour of picking) declines. We can depict her production function $Y = F(L)$ by the upward-sloping curve in Figure 1, relating coconuts produced (Y) to hours spent picking (L). The production function represents her opportunity set: she can choose her number of coconut-picking hours (and the amount of coconuts picked) from any point on (or below) the production function in order to maximize her utility.

Given that she enjoys playing Frisbee golf more than picking coconuts, and that the number of hours spent playing golf (G) is the total number of daylight hours L_0 minus her time spent picking coconuts ($G = L_0 - L$), her utility is increasing as she moves up and to the left in Figure 1, as shown by the large arrow. Her ideal point would be few hours picking and lots of coconuts if that were feasible. Her indifference curves look like the three curves shown in Figure 1 and her utility maximizing point is where the highest attainable indifference curve is tangent to her production constraint, at (L^*, Y^*) . Her optimal choice is to work L^* hours, play Frisbee golf for $L_0 - L^*$ hours, and produce and eat Y^* coconuts.

Jane would be worse off if she produced and ate *more* than Y^* or *less* than Y^* coconuts. In macroeconomic terms, Y^* is her efficient level of output. More output is not better if it takes her above efficient output, because the marginal disutility of work would exceed the marginal utility of the resulting increase in output. Note that

Jane *could* produce more output than the efficient level; Y^* is not a physical limitation on the amount of output that can be produced. It is the level at which the benefits of increasing production no longer justify the costs.

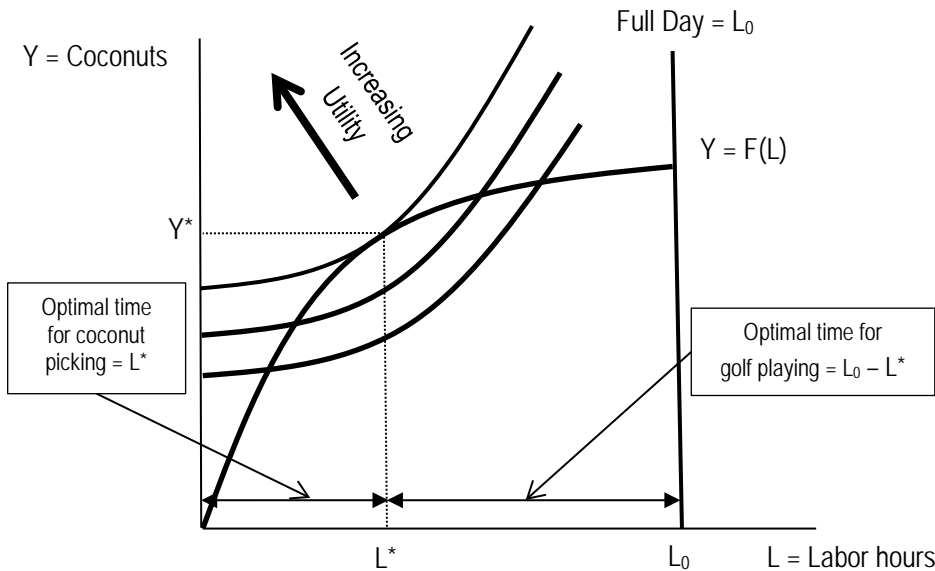


Figure 1. Jane's utility-maximizing equilibrium

In terms of familiar microeconomic concepts, the slope of Jane's production function is her *marginal product of labor*. The slope of her indifference curve is her *marginal rate of substitution* between coconuts and leisure (golf) time. Her utility is maximized where her marginal product of labor equals her marginal rate of substitution: $MPL = MRS$.

Efficient output with multiple people and firms

When we move from the one-person economy to a multi-person economy with exchange of goods and labor, the labor and production decisions seem more complicated, but the same underlying principles apply. Suppose that Jane arbitrarily decided to pay herself a wage for the hours worked picking coconuts then sell herself the coconuts at the price that would “clear the market.” Nothing would be different. She would balance the marginal utility of the goods that her wage would buy (her “real wage”) against the marginal utility of the Frisbee golf that she would have to forgo by working more. She would end up producing the same efficient level of output.

Introducing a labor market and a goods market does not fundamentally change her decision-making. As a “worker” she would work up to the point where the mar-

ginal rate of substitution between coconuts and leisure equals her real wage. As a “producer” or “firm” she would hire her work up to the point where her marginal product of labor equals the real wage. Thus, she would still produce where $MPL = MRS$. In terms of Figure 1, think of a line drawn through the equilibrium point that is tangent to both the indifference curve and the production function. As a worker, she chooses the point at which the line is tangent to her indifference curve; as a firm, she chooses the point at which the line is tangent to her production function.

The principles that Jane used to determine her efficient level of output in a single-person economy extend directly to economies involving many people and firms. As long as markets operate efficiently to equalize the marginal product of labor with workers’ marginal rate of substitution between goods and leisure, the economy will produce its optimal, or efficient, level of output. This is true even if there are many different goods and many different kinds of labor, as long as each good and each kind of labor is traded in a “perfect” or “efficient” market of the kind that we assume to exist under perfect competition.

Like Jane, a complex macroeconomy is physically capable of producing more than the efficient level of output, but this would not be optimal. The benefits of the added goods and services would not justify the cost of the forgone leisure required to produce them. Once again, the efficient level of output is *not* the maximum possible flow of production; it is the one that an economy in which all decisions are made efficiently would choose.

Natural output vs. efficient output

The phrase above, “as long as each good and each kind of labor is traded in a perfect or efficient market,” should raise your eyebrows: It is simply never true. Every complex economy is made up of imperfect markets in which many kinds of market failures combine to prevent perfect equalization of MPL with MRS . In the presence of microeconomic market failures, the macroeconomy will generally settle at a ***natural level of output*** lower than the efficient level. The sources of market failure are well examined in most microeconomics courses and will only be summarized here.

Any sort of monopoly power in the market would cause the monopoly firm to produce less output than the efficient amount because it maximizes profit where price exceeds marginal cost. Firms in all modern economies enjoy some degree of market power, even if they are not strictly monopolies, so we expect that natural output will normally fall below the efficient level as a result.

Providing for public goods requires the collection of taxes by the government. It is not usually feasible to use lump-sum taxes, so the existence of taxes creates a market distortion driving a wedge between MPL and MRS .¹ Nearly all taxes in modern

¹ A lump-sum tax is one in which the amount that someone pays does not depend in any way on his or her actions. Such taxes are regressive and highly unpopular, so they are rarely im-

economies lead to a reduction in the amount that people choose to work, produce, or consume. Like monopolies, these kinds of taxes reduce the natural level of output below the efficient level.

Perfect competition also assumes perfect information and complete adjustment of prices to clear markets. Both of these assumptions are more plausible in the long run (when prices have time to respond to excess demand and supply and people have time to learn about the price changes) than in the short run. Relaxing these two assumptions has had an important role in the evolution of modern macroeconomic theory, as we discuss in Romer's Chapters 6 and 7.

In many economies, more subtle economic policies also reduce efficiency. If the political system directs resources to industries or firms that are favored by government policymakers rather than those that the market would select as most efficient, this misallocation of resources will lower economic efficiency and cause natural output to lie below efficient output.

Finally, and perhaps most importantly, labor markets rarely if ever approximate perfect competition. Workers (and jobs) differ more one from another than any other traded good or service. Differences in skills divide "the labor market" into thousands of sub-markets. Jobs and workers have specific preferred locations, which further sub-divides the market. And when you add in the personal preferences of workers (and employers), hiring of workers is often more like "match-making" than it is like a "market."

When we combine labor heterogeneity with imperfect information about the jobs available to a searching worker and about the workers available to fill a vacant job, it becomes costly to "make matches" in the labor market. Both job-seeking workers and employers with vacant jobs must invest time and effort into job search. The more effort they put in, the more efficiently they will allocate workers into jobs. But search, like almost everything else, is costly and subject to diminishing returns, so prospective workers and hiring firms will not search with "infinite intensity." The result is that there is, at every moment, a pool of unemployed workers and a pool of vacant jobs that have not yet been "matched up." We call this *frictional unemployment*; it reduces actual employment, and therefore actual output, below the efficient level.

Another source of unemployment in every economy is mismatching of skills and locations between the pools of unemployed workers and vacant jobs. The frictional unemployment discussed above assumes that there exists a suitable match between

posed. Most common forms of taxation apply taxes to actions: an income tax is a tax on earning income, a sales or excise tax is a tax on purchasing goods, etc. These kinds of taxes serve to reduce the activity being taxed below the efficient level. If Jane faced an income tax on her collection of coconuts (or on her consumption—it's the same here), it would shift her after-tax production function downward, causing her (most likely) to substitute more Frisbee golf for coconut production and consumption.

unemployed workers and vacant jobs, but that they simply have not yet found one another. What if the vacant jobs all require skills in Web design and the unemployed workers are all former assembly-line workers who lack these skills? What if the vacant jobs are all in Atlanta and the unemployed workers are in Detroit? These mismatches lead to *structural unemployment*. Specific situations of structural unemployment may tend to go away over time (as workers retrain for new jobs or relocate to places where jobs are available), but if the structure of labor demand changes rapidly across industries, locations, and occupations, it is likely that there will be long-lasting gaps of structural mismatch.

Frictional and structural unemployment cause the *natural rate of unemployment* to be positive. This leads the actual “equilibrium” or natural level of employment to fall below the efficient level that would prevail in perfect markets. Some people whose utility-maximizing choice is to work are left unemployed while searching for a job; some firms whose profit-maximizing choice is to hire another worker are left with vacant jobs until a suitable worker can be found.

The natural level of output is the amount that the economy produces when all of these market imperfections are taken into account. Because monopolies, taxes, distortions, and natural sources of unemployment reduce the amount that the economy produces, natural output is generally below efficient output. The magnitude of the shortfall depends on the pervasiveness of these microeconomic sources of inefficiency in the economy. Economies that have fewer monopolies, lower tax rates, and more homogeneous labor markets are likely to achieve more of their economic capacity than others, and have a higher natural level of output relative to efficient output.

Growth theory and behavior of natural output

Theories of economic growth examine the long-run evolution of an economy’s natural output along its *growth path*, which is, visually, a graph with natural output (or, for reasons to be discussed in the next section, the log of natural output) on the vertical axis and time on the horizontal. Growth theories, such as those Romer discusses in Chapters 1, 2, and 3, typically focus on factors that cause the efficient level of output to increase over time. A companion set of theories, some of which are discussed in Romer’s Chapter 4, look at the determinants of the gap between efficient and natural output.

For solitary Jane, there is no possibility of market failure, so the efficient and natural levels of output coincide. Her natural level of output depends on three things: (1) her endowment of labor (the number of daily daytime hours available for picking or golf), (2) her ability to transform work into coconuts (her production function), and (3) her preferences toward work and leisure (her indifference map). A change in

any of those factors could cause her efficient output to change. For example, we would expect her to pick more coconuts in the summer in response to a longer.

In modern macroeconomies, we think of the efficient and natural levels of output as being determined by three factors:

- **Labor input.** For Jane, this was the number of daytime hours coupled with her preferences toward work. For macroeconomies, it is the size of the adult population, multiplied times the share of the adult population that chooses to work, multiplied by the average number of hours worked. The efficient level of output incorporates the effects of demographic changes as well as changes in individual preferences about work. Natural output is also affected by distortions such as taxes that affect work decisions.
- **Capital input.** Imagine that Jane devotes some of her time to constructing a ladder to reduce the time needed to harvest coconuts. This would shift her production function (as a function of labor input) upward. For macroeconomies, we think of the capital stock as the available stock of productive structures and equipment. This stock increases when individuals decide to use some of their current resources to produce factories and machines rather than consumption goods—saving and investment.
- **Technology.** Suppose that Jane discovers a better method of climbing trees. Like the ladder, this would shift Jane’s production function upward. In macroeconomies, improvements in knowledge about production techniques (such as new inventions and innovations) likewise shift the production function upward.

Growth theory models the natural level of output using an aggregate production function that depends on these three factors: $Y = F(K, L, A)$, where K is the aggregate capital stock, L is the aggregate labor force, and A is an index of technological progress or “knowledge.” The growth path of natural output depends on the changes over time in K , L , and A .

Capital increases over time through the saving and investment choices of households and firms in the economy. One key element of every growth model is the decision of how much of current output is consumed immediately and how much is saved or invested in the form of capital goods to be used in future production. Because saving and investment are key economic decisions, the evolution of the capital stock is endogenous in most growth models.

Labor increases over time as the population grows, which is usually an exogenous variable. Secondary factors could be changes in the rate of labor-force participa-

tion and in average hours worked per worker, but growth theory usually either ignores these effects or takes them as exogenous.

Technology improves through advances in applied knowledge about production methods. The neoclassical growth models that were developed before the 1980s took the rate of technological improvement to be exogenous, determined by the progress of scientific and engineering knowledge. Modern, “endogenous” growth models have recognized that technological progress often results from economic decisions to invest resources in research and development. These models treat the “stock of knowledge” A as another kind of capital, in which economies can invest by choosing to pay for R&D activities rather than buying consumer goods or traditional, physical capital goods.

Every growth model has two key components—careful attention to these components is essential to understanding how differences in the models’ assumption affect their conclusions:

- A production function that describes how output is related to the inputs.
- Laws of motion that describe how the inputs to production (L , K , and A) change over time, either exogenously or endogenously in response to economic variables.

Given these two components, solving a growth model consists of determining the equilibrium growth path along which natural real output Y moves over time. Although this is not necessarily the case, the growth models that we shall study all end up having a “steady-state growth path” along which the economy eventually settles onto a path with a constant growth rate. Our task in understanding the implications of growth theory is to explore how changes in key parameters—such as the amount or current resources that economies choose to dedicate to investing in capital and/or research and development—lead to changes in the steady-state growth path.

In particular, we shall see examples of three kinds of effects that a change in the economy can have on its growth path:

- No effect at all.
- A “level effect,” shifting the level of the growth path upward but not changing the steady-state growth rate (the slope of the path).
- A “growth effect,” increasing the steady-state growth rate of the economy.

C. Growth in Continuous Time: Logarithmic and Exponential Functions

Continuous-time vs. discrete-time models

When we construct any kind of dynamic model such as a growth model, we must decide whether time should pass in discrete intervals or as a continuous flow. *Discrete-time models* assume that there is an interval of time—one period—for the duration of which the values of all variables remain unchanged. When a period ends, all variables may jump to different values for the next period, but they then remain unchanged during that period. The step function shown in Figure 2 shows a typical graphical representation of a time path of a variable in a discrete-time model.

In continuous-time models, time flows continuously and variables can change to new values at any moment. A typical variable in a continuous-time model might have a time path like the smooth line in Figure 2.

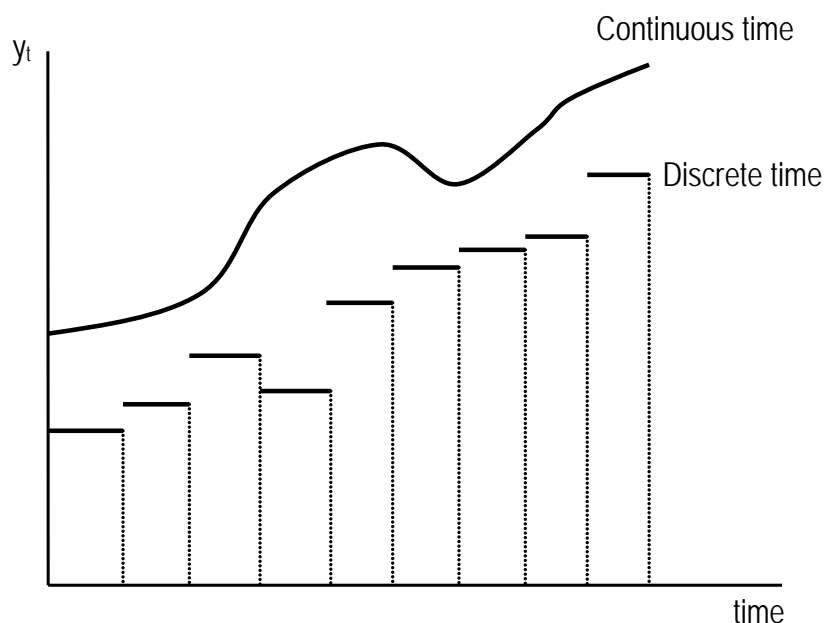


Figure 2. Discrete- and continuous-time variables

Although we usually think of time as flowing continuously, there are actually many examples of discrete time in real economies. The price of gold is fixed twice daily, for example, and banks reckon one's deposit balances once a day at the close of business. Moreover, all macroeconomic data are published only at discrete inter-

vals such as a day, month, quarter, or year, even when the underlying variables move continuously. In these cases, the single monthly value assigned to the variable might be an average of its values on various days of the month (as with some time-aggregated measures of interest rates and exchange rates) or its value on a particular day in the month (as with estimates of the unemployment rate and consumer price index).

The world we are modeling has elements of both continuous and discrete time so neither type of model is obviously preferable. We usually choose the modeling strategy that is most convenient for the particular analysis we are performing. Empirical models are nearly always discrete because of the discrete availability of data, while many theoretical models are easier to analyze in continuous time. We shall examine models of both kinds during this course. The first growth models we encounter are in continuous time, so we shall preface that analysis with some discussion of the mathematical concepts used to model continuous growth.

Growth in discrete and continuous time

You are probably more used to thinking of growth rates, inflation rates, and other rates of change over time in terms of discrete, period-to-period changes. Empirically, this is a natural way of thinking about growth and inflation because macroeconomic data are published for discrete periods. We typically calculate the discrete-time growth rate of real output y from year t to year $t + 1$ as $g_y = (y_{t+1} - y_t) / y_t = \Delta y / y$, where Δy is defined to be the change in y from one year to the next. As we discussed above, such discrete growth calculations correspond to a world where the flow of output is constant throughout a period (year), then moves to a possibly different level for the next period.

In the discrete case, a variable growing at a constant rate g increases its value by $100g$ percent each year. If $g = 0.04$, then each year's value is 4% higher than the previous year's, or $y_{t+1} = (1 + g) y_t = 1.04 y_t$. Applying this formula year after year (with the growth rate assumed to be constant) yields $y_{t+2} = (1 + g)y_{t+1} = (1 + g)^2 y_t$ and, in general, $y_{t+n} = (1 + g)^n y_t$.

However, one ambiguity with discrete growth rates (and discrete-time analysis in general) is that the length of the period is, in principle, arbitrary. To see how this affects the calculation of growth rates, suppose that we have quarterly data so that there are four observations for each year. The value of the variable in the first quarter of the first year is y_1 , y_2 is the value in the second quarter of the first year, and so on through the years, with y_5 through y_8 being the observations for the four quarters of the second year, etc. Can we use the formula $g_y = (y_{t+1} - y_t) / y_t$ for this case? Yes and no. Although this formula gives us a growth rate, that growth rate is now expressed as a rate of *growth per quarter* rather than the conventional *growth per year*—a value of 0.04 now means that the variable increases by 4% each quarter, not 4% per year. For ease of comparison, we prefer to express growth rates, inflation rates, and interest

rates in “annual” rates (percent per year), so the quarterly growth rate calculated by this formula would not give a number comparable to our usual growth-rate metric.

To convert the quarterly (percent per quarter) growth rate to an annual rate (percent per year), we must think about how much a variable would grow over four quarters if its quarterly rate of growth was, say, g_q . In other words, we want to know how much bigger y_{t+4} is in percentage terms than y_t if y grows by g_q per quarter. By the reasoning above, $y_{t+4} = (1 + g_q)^4 y_t$, so if g is the annual growth rate,

$$1 + g = (1 + g_q)^4. \tag{1}$$

Using basic laws of exponents, $1 + g_q = (1 + g)^{1/4}$, so we can express the value of y for n quarters after date t as $y_{t+n} = (1 + g_q)^n y_t = (1 + g)^{n/4} y_t$.

One obvious question is whether formula (1) is the same as $g = 4g_q$. The answer is no. For example, if $g_q = 0.01 = 1\%$, then $1 + g = (1.01)^4 = 1.04060401$, so $g = 4.060401\% > 4\%$. This is because of the **compounding** of growth—the effect of the expansion over time in the base to which the growth rate is applied. The formula $g = 4g_q$ reflects no compounding: a fraction g_q of the *initial* quarter’s value of y is added in each quarter. But by the second quarter, the value of y has grown, so the amount of increase in y in the second quarter will be larger than in the first quarter. Similarly, the third and fourth quarters will have even larger amounts of absolute increase in y . The cumulative effect of this compounding causes the annual growth rate of the variable to be more than four times the quarterly growth rate, though when the growth rates are small this difference may not be very substantial over short periods of time.

So now we have a formula that allows us to translate between quarterly and annual growth rates. However, there is nothing particularly special about quarterly growth. If we considered one month to be the time period, then by similar reasoning the annual growth rate g would be related to the monthly growth rate g_m by $1 + g = (1 + g_m)^{12}$. Using a weekly time period, $1 + g = (1 + g_w)^{52}$, and if we have a daily period, $1 + g = (1 + g_d)^{365}$ (except in leap years). Using logic parallel to that used above, the level of the daily-growth variable n days after date t would be related to the date t value by $y_{t+n} = (1 + g_d)^n y_t = (1 + g)^{n/365} y_t$.

As you can see, the algebra varies depending on the choice of time units: years, quarters, months, weeks, or days. In empirical applications, we are usually restricted to these discrete time units by the constraints of the available data. National-account statistics are published only as quarterly or annual averages; the consumer price index is published monthly; exchange rates and prices of financial assets are usually available daily or even hourly.

In a purely theoretical model, we are not constrained by data availability and it is often more convenient and intuitive to think of variables as moving continuously through time rather than jumping from one level to another as one finite period ends and the next begins. Analytically, continuous-time modeling allows us to think of

our variables as continuous functions of the time variable t , which means that the methods of calculus and differential equations can be applied.

In continuous-time models, t can take on any value, not just integer values. If $t = 0$ is defined to be midnight at the beginning of January 1, 2001 and periods are normalized at one year, then $t = 0.5$ would be exactly one-half year later, $t = 1.0$ would be one year later, etc. To reflect this continuous variation, we typically use the notation $y(t)$ rather than y_t to denote the value of variable y at moment t . The change in y per unit time at moment t is a “time derivative” $dy(t)/dt$, which is commonly denoted by $\dot{y}(t)$.²

The time derivative measures the amount of change per period (year) in a variable as time passes, so it is analogous to the discrete-time “first difference” $\Delta y = y_{t+1} - y_t$. The time derivative or first difference tells the *amount* of growth in y , but not the *rate* of growth. In order to convert the time derivative or first difference into a growth rate (percentage change per year), we divide it by the level of the variable. In discrete time this gives us $g = (y_{t+1} - y_t) / y_t = \Delta y / y$. In continuous time, the corresponding growth rate is $\gamma = \dot{y}(t) / y(t)$ or just \dot{y} / y . The continuous-time growth rate incorporates “continuous compounding,” which is the limiting case as the period of compounding shrinks from a year to a month to a day and down to zero.

So if a variable grows continuously (with continuous compounding) for n years, how much bigger will it get? In discrete time (with an annual time unit and annual compounding), we used the formula $y_{t+n} = (1 + g)^n y_t$ to calculate this. In the continuous case, the corresponding formula is

$$y(t+n) = e^{gn} y(t), \tag{2}$$

where e is the constant (approximately 2.71) that is the base of the natural logarithms.

Exponentials, logs, and continuous growth

Equation (2) shows that the value of a variable growing at a constant rate is an exponential function of time. In Romer’s analysis of the Solow growth model, we assume that the labor force L and the stock of knowledge A both grow at given constant rates. Applying our equation (2) from above gives Romer’s equations (1.13) and (1.14) on page 14.

In graphical terms, a variable following a constant-growth path looks like the one shown in Figure 3, which begins in 1900 with a value of 100 and increases by 5 percent per year until 2000. The formula for the value of this variable is

² The next section discusses time derivatives and presents some useful rules for calculating the growth rates of products, quotients, and exponential functions of variables.

$$y(t) = 100 e^{0.05t}, \quad (3)$$

where t is defined as a “trend” variable with value zero in 1900, one in 1901, two in 1902, etc.

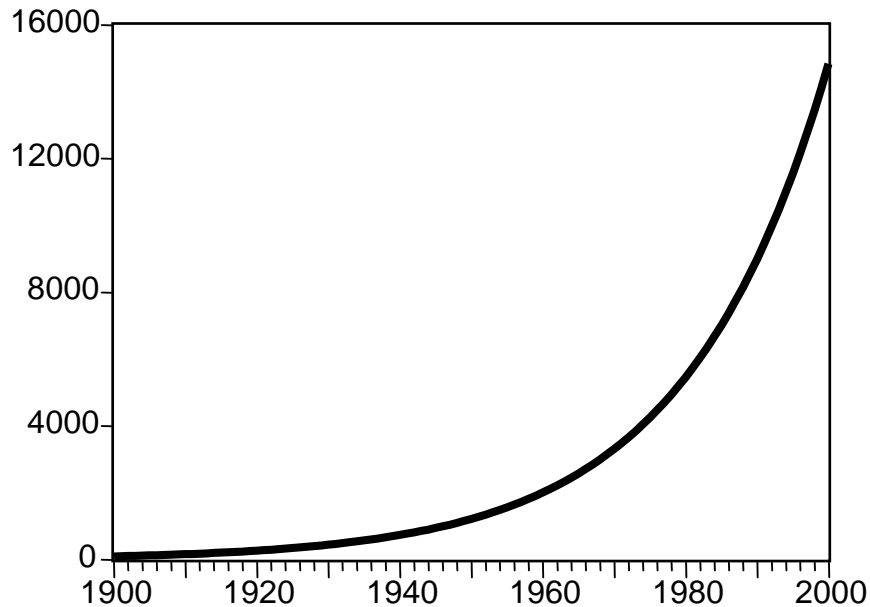


Figure 3. Continuous growth at constant rate

Two things are apparent from Figure 3. First, exponential growth, even at a fairly modest rate such as 5 percent, leads to huge increases in a variable over a long period of time. This is the “miracle of compound growth,” that allows modest sums invested early in life to provide large retirement incomes through compound interest.

The second notable, if obvious, feature of the time path in Figure 3 is that it is not a straight line. This can make life difficult, not only for economics professors who are used to drawing (or trying to draw, ☺) straight lines on the blackboard but also because it makes it hard to tell constant-growth paths from other paths where the growth rate varies over time.

Because straight lines are very convenient, it would be nice to find a way to represent a constant-growth-rate path as a straight line. The natural logarithm function provides a way to transform the exponential growth path into a line. The natural log, which we sometimes write as \ln , is the inverse function to the exponential function: by definition, $\ln e^x = x$. Logarithms also have the well-known property that the log of

a product (quotient) is the sum (difference) of the logs of the two things being multiplied (divided).

Applying these rules to the formula in equation (3) allows us to write the natural logarithm of that variable y as $\ln y = \ln(100) + gt = 4.605 + gt$, which is a linear function of time. Figure 4 shows a plot of the time path of $\ln y$; you can see that it is a straight line. However, the fact that the numbers on the vertical axis are values of $\ln y$

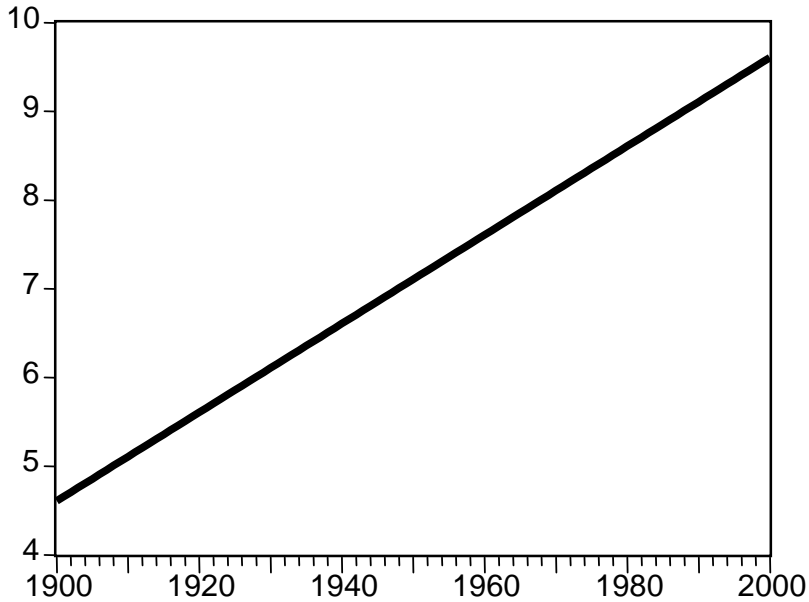


Figure 4. Constant growth rate in logarithmic space

rather than y is a disadvantage when we try to interpret Figure 4. To make the numbers easier to interpret, we sometimes use the values of y rather than $\ln y$ on the vertical axis as in Figure 5. (Note that Figure 4 and Figure 5 are identical except for the numbers and tick marks on the vertical axis.)

The disadvantage of using a “log scale” as in Figure 5 is that a given vertical distance in the graph represents a particular amount of percentage change in y rather than a particular absolute change. In Figure 5, the vertical tick marks for 4000 and 8000 are farther apart than those for 8000 and 12000. Depending on the circumstance, we may find it easiest to use a “normal” graph like Figure 3, a graph of the log like Figure 4, or a log-scale graph like Figure 5. However, the main point here is that if x grows at a constant rate in continuous time, then the plot of $\ln x$ against time will be a straight line whose slope equals the growth rate of x .

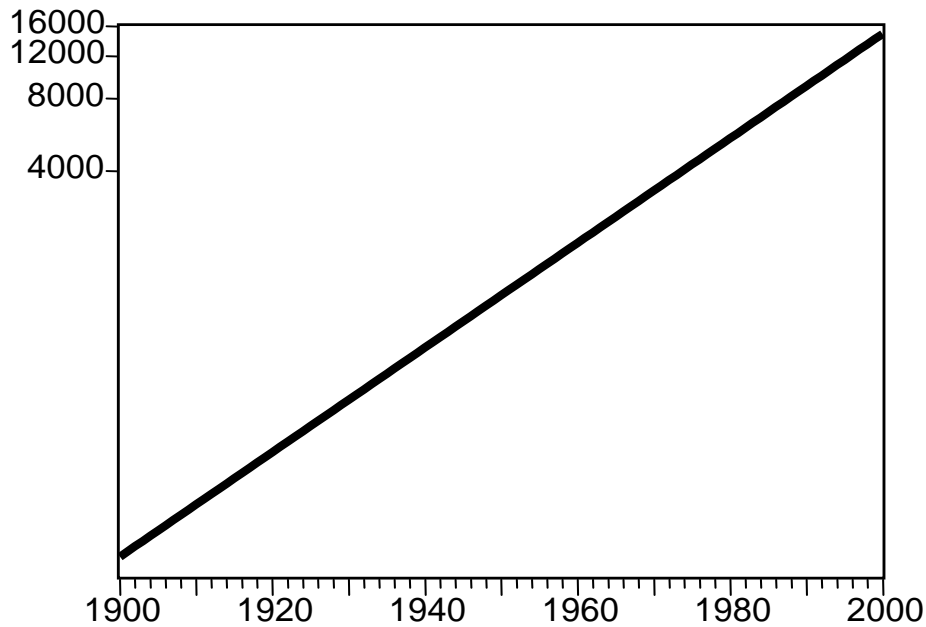


Figure 5. Using a "log scale" on the vertical axis

D. Review of Some Basic Calculus Tools

Although calculus is a fundamental tool of economics, most undergraduate courses sidestep using it by relying on graphs and algebraic analysis of linear models. However, the concepts of calculus are so intimately related to the task of economic modeling that it is often intuitively clearer, as well as analytically more elegant, to talk about economics using the language of calculus. This section and the section on constrained optimization in the next chapter develop some basic tools and notation, so that you will be more comfortable reading and understanding the texts. They do not attempt to teach you any but the most elementary properties of derivatives and integrals. A deeper knowledge of calculus such as that presented in Math 111 (and higher-level math courses) at Reed is an important part of the economics major's tool kit.

Calculus is concerned with relationships between two or more variables. The particular kind of relationship for which we can employ calculus tools is called a *function*. A function relates one variable (the dependent variable) to one or more others in a particular way: if f is a function relating a dependent variable y to a set of inde-

pendent variables x_1, x_2, \dots, x_n , then any admissible set of values for the x variables must correspond to a unique value of y . We write the functional notation as $y = f(x_1, x_2, \dots, x_n)$. The x variables are called the “arguments” of the function. The simplest functions are *univariate*; they have only one variable as an argument, so $y = f(x)$. We begin by developing some calculus concepts for univariate functions, then we extend the analysis to *multivariate* functions.

In economics and other sciences, we frequently want to know how a change in a function’s independent variable affects the dependent variable. In particular, we are interested in the magnitude $\Delta y/\Delta x$, where we use the capital Greek letter delta (Δ) to mean “a small change in.” The ratio $\Delta y/\Delta x$ tells the amount of change that is induced in y for each unit of change in x . In macroeconomics, we sometimes call such a ratio a “multiplier.” If we graph a function with the dependent variable on the vertical axis and the independent (argument) variable on the horizontal axis, then $\Delta y/\Delta x$ is the slope of the function.

Unless the function is linear, a slope measured between two points on the curve will depend on which two points are chosen. For example, in Figure 6 we could measure the slope between points a and b , which gives $50 - 40 = 10$ for Δy and $14 - 10 = 4$ for Δx , with a slope of $10/4 = 2.5$. We could also measure slope between points a and c , which gives a slope of $10/10 = 1$, or between points a and d , which gives a slope of $-15/20 = -0.75$.

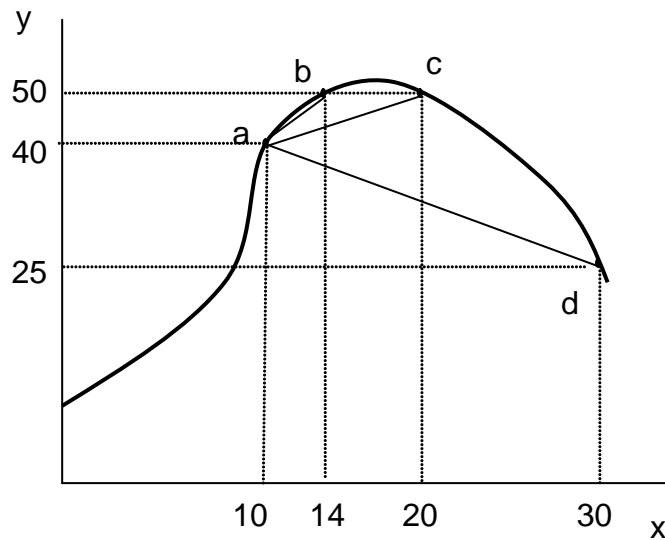


Figure 6. Slopes and derivatives

We calculated all of these slopes by a general formula for the slope between two points. Let's call the value of x at the initial point x_0 and the value after the change $x_0 + \Delta x$. Then the slope of the function between x_0 and $x_0 + \Delta x$ is

$$f^*(x_0, x_0 + \Delta x) = \frac{f(x_0 + \Delta x) - f(x_0)}{(x_0 + \Delta x) - x_0} = \frac{f(x_0 + \Delta x) - f(x_0)}{\Delta x}. \quad (4)$$

The function f^* is “derived” from the function f —for any function f there is a unique function f^* that gives the slope of f between any pair of x values.

The f^* function defined in (4) gives the slope of the chord connecting two points on the curve: $(x_0, f(x_0))$ and $(x_0 + \Delta x, f(x_0 + \Delta x))$. However, we are often interested in the behavior of the function only in a small neighborhood around a point such as a . For this, we use another “derived” function—called the *derivative function*—that gives the slope of the line that is *tangent* to the curve at a particular point. The tangent line is the line that touches the curve at exactly one point with the tangent line (usually) lying entirely on one side of the curve as in Figure 7. In contrast to the f^* function above, which depended on both x_0 and $x_0 + \Delta x$, the derivative function takes only one argument: the value of x at the point at which the tangent line touches the curve.

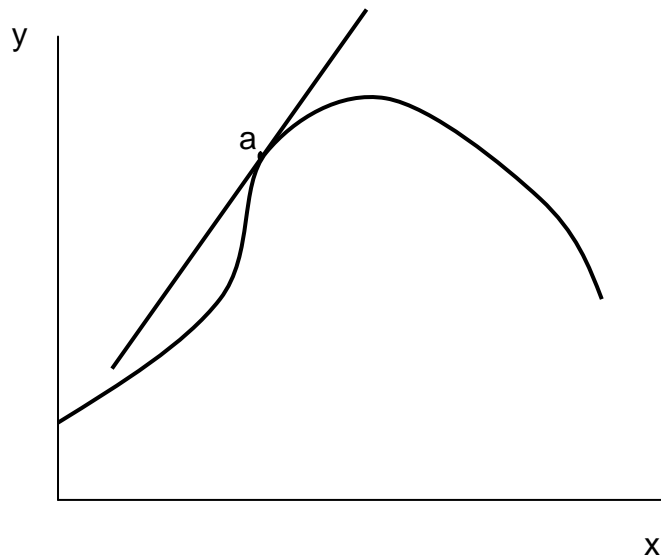


Figure 7. Tangent line to a continuous function

We can think about the tangent line at point a as the limit of a sequence of chords connecting a with other points on the curve, such as the line segments drawn

in Figure 7. In terms of algebra, the slope of this limiting line is obtained by taking the limit of equation (4) as the two points get very close together, *i.e.*, as Δx gets close to zero. The derivative function, often denoted by $f'(x)$ or by dy/dx , is given by

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}. \quad (5)$$

Equation (5) is the formal definition of the derivative of the function f . However, it would be awkward to have to take a formal limit every time we want to find the slope of a function. There are some simple formulas for the derivatives of common functions so that you will not need to take limits. We will examine some of these formulas below.

An alternative notation that is commonly used for the derivative is dy/dx , which is a direct analogy to the $\Delta y/\Delta x$ notation used for slope. You must be careful with this definition, though, because the dy and dx terms are not really numbers, they are infinitesimal changes that are sometimes called *differentials*. Thus, while we sometimes multiply both sides of $dy/dx = f'(x)$ by dx to get $dy = f'(x) dx$, we must remember that this formula holds exactly only for infinitesimal changes in x and y .

Derivatives of powers, sums, products, and quotients

Finding the derivative of a function is called *differentiation*. The following basic rules of differentiation apply to all functions that have derivatives. In these rules, f , g , and h are all functions of a single variable.

1. The derivative of a constant times a function is the constant times the derivative of the function: If $g(x) = c f(x)$, then $g'(x) = c f'(x)$.
2. If the function is the variable raised to a power, then the derivative is the number of the power multiplied by the variable raised to one less power: If $f(x) = x^n$, then $f'(x) = n x^{n-1}$. This formula works for *all* values of n , positive or negative, integer or not. For example, the derivative of the function $f(x) = x^2$ is $f'(x) = 2x$; the derivative of $f(x) = x^3$ is $f'(x) = 3x^2$; the derivative of $f(x) = x^1 = x$ is $f'(x) = 1$; the derivative of the constant function $f(x) = ax^0 = a$ is $f'(x) = 0$; the derivative of $f(x) = x^{-1}$ is $f'(x) = -x^{-2}$; and the derivative of $f(x) = x^{1/2}$ is $f'(x) = \frac{1}{2} x^{-1/2}$.
3. The derivative of a sum of two functions is the sum of the derivatives of the functions. If $h(x) = f(x) + g(x)$, then $h'(x) = f'(x) + g'(x)$.
4. The derivative of a product of two functions is given by the following formula: If $h(x) = f(x) g(x)$, then $h'(x) = f(x) g'(x) + g(x) f'(x)$.

5. The derivative of a quotient of two functions is given by the following formula: If $h(x) = f(x) / g(x)$, then $h'(x) = [g(x)f'(x) - f(x)g'(x)] / [g(x)]^2$.

Using these formulas, we can calculate the derivatives of a wide variety of functions. For example, if $h(x) = (4x^2 - 3x + 7) / (x^3 + 7x + 4)$, then we can apply the quotient rule letting the function in the numerator be $f(x)$ and the denominator be $g(x)$. Using the rules for powers, sums, and multiplication by a constant, $f'(x) = 8x - 3$ and $g'(x) = 3x^2 + 7$. Thus,

$$h'(x) = [(x^3 + 7x + 4)(8x - 3) - (4x^2 - 3x + 7)(3x^2 + 7)] / (x^3 + 7x + 4)^2,$$

a complicated expression, but one that was obtained a lot more easily by the formulas than by taking limits of everything.

Derivatives and maximization

Since the derivative gives the slope of a function at each point, we can use the derivative to tell whether the value of the function is increasing, decreasing, or flat at that point. If the derivative is positive at a particular value x_0 , *i.e.*, $f'(x_0) > 0$, then the function is upward-sloping or increasing at x_0 . Similarly, a negative derivative indicates a downward-sloping or decreasing function at that particular point. At a point where the derivative is zero, the tangent line to the function is horizontal.

In economics we often want to find the maximum or minimum value of a function. For example, we often model households as maximizing utility and firms as maximizing profit or minimizing cost. At a point where a function reaches a maximum or minimum relative to the points around it, its slope is zero. To the left of a maximum (minimum) point it has positive (negative) slope and to the right it has negative (positive) slope. Thus, finding the values for which a function's derivative is zero identifies all the values that might be (local) maxima or minima.

Suppose that a firm's profit is related to its level of output by the function

$$P(q) = 1000 + 500q - 2q^2.$$

We can identify the possible maximum or minimum points of this function by taking its derivative and setting it equal to zero: $P'(q) = 500 - 4q = 0$. Solving this equation for q gives $4q = 500$ or $q = 125$. Thus, profit may be at a maximum or minimum when 125 units are produced.

Since the derivative of a function is zero at both maximum and minimum points, how are we to know whether $q = 125$ is a point where profit is maximized or mini-

mized? There are two ways we could do this. One would be to examine the derivative just above and below 125. When $q = 124$, the derivative is

$$P'(q) = 500 - 4q = 500 - 4(124) = 4,$$

so the curve is upward sloping to the left of $q = 125$. When $q = 126$, $P'(126) = 500 - 4(126) = -4$, so the curve slopes downward to the right of 125. Thus, we are assured that producing 125 units maximizes the firm's profit.

A more precise way (because we can never know how "close" to 125 we need to be) of distinguishing maxima from minima is to use the *second derivative*. Just as the derivative function tells the rate at which the value of the function changes as x changes, we can take the derivative of the derivative to find out how the derivative, or slope, function is changing as x changes. If the slope is increasing at a point where it is zero, then it is going from negative to positive and the function is at a minimum. If the slope is decreasing, then it is going from positive to negative and the function is at a maximum.

The second derivative, denoted $f''(x)$, is found by applying the rules of differentiation to the first derivative function $f'(x)$. In the case of the profit function, $P'(q) = 500 - 4q$, so $P''(q) = -4 < 0$. The second derivative of the profit function is negative, so the function is surely at a maximum.

The second derivative tells us about the curvature of the function. A negative second derivative means that the function opens downward, or is concave. A positive second derivative indicates a function that opens upward, which is called a convex function.³

Other rules of differentiation

There are several other rules of differentiation that we will need later in the course. Since we will be working with (natural) logarithms frequently, the derivative of the log function will often be important. If $f(x) = \ln x$, then $f'(x) = 1/x = x^{-1}$. Since we saw above that power functions typically differentiate into other power functions, it may seem surprising that the log function also differentiates into a power function. However, recall that differentiating a power function gives a power function with the exponent reduced by one. Thus, the power function that could possibly give $1/x = x^{-1}$ would be x^0 . However, the derivative of x^0 is $0 \cdot x^{-1} = 0$. Thus, there is no power function that gives a derivative involving $1/x$; the log function does so instead.

The inverse of the natural log function is the exponential function $f(x) = e^x$, where e is the natural constant equal to approximately 2.71. This function has the

³Note that some texts use the opposite definitions for convex and concave.

unique property that it is its own derivative: $f'(x) = f(x) = e^x$. It is a function whose slope is equal to the value of the function at every point.

The final rule of differentiation that we study here is a rule for taking the derivative of a function of a function. Suppose that $h(x) = g[f(x)]$. The rule for differentiating such chains of functions is called the **chain rule** and is $h'(x) = g'[f(x)] \cdot f'(x)$. For example, if $h(x) = \ln x^2$, then we can think of the log function as g and the square function as f . Applying the chain rule yields $h'(x) = (1/x^2) 2x = 2/x$.

E. Calculus Applications in Macroeconomics

An application: time derivatives

In the study of economic growth, our primary interest is on how variables change over time. Using our rules of differentiation, we can think of the amount by which a variable y changes per period as a **time derivative**, dy/dt , where t is time. When we work with variables in continuous time, the time derivative plays a role analogous to the role played in discrete time by the **(first) difference** of the variable, $\Delta y = y(t+1) - y(t)$.⁴ If we plot the path of the variable with time on the horizontal axis as in Figure 8, then the first difference of the variable is the slope of the line segment connecting the point $(t, y(t))$ with the point $(t+1, y(t+1))$; the time derivative at time t is the slope of the line that is tangent to the time path at the point $(t, y(t))$. To economize on notation, we often represent the time derivative of y as $\dot{y} \equiv dy/dt$.

⁴It may seem like Δy is analogous to only the numerator of dy/dt . To see why the denominator disappears, note that the difference can be thought of as the ratio of the change from one year to the next in y (i.e., Δy) to the change in t (which is $t - (t-1) = 1$). Because the change in t is exactly one, the denominator vanishes.

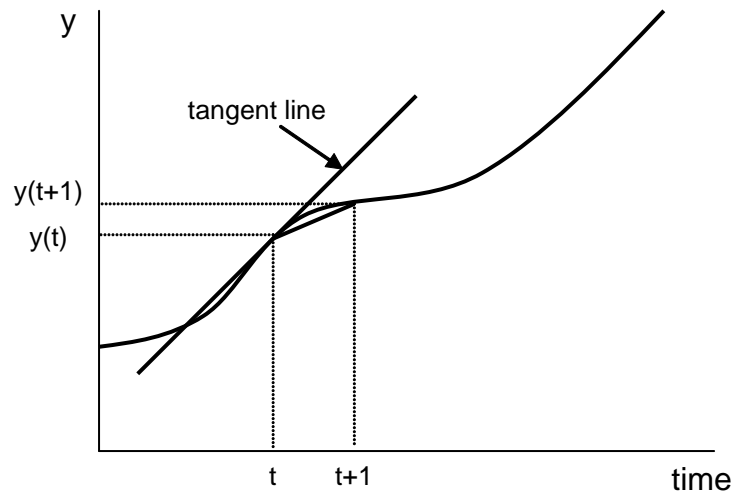


Figure 8. Time derivatives

The relationship between a variable and its time derivative (or its difference) is analogous to the relationship between a stock and a flow. For example, if $K(t)$ represents the size of an economy's capital stock at time t , then $\dot{K}(t)$ is the rate of change of K at time t . Suppose first that capital never wears out. Then the capital stock increases at the rate that new investment is put in place. If we denote the flow of investment in new capital by $I(t)$, then $\dot{K}(t) = I(t)$. Taking account of capital wearing out is only a little more complicated. Depreciation is the flow of reductions in the value of the capital stock due to wearing out and obsolescence. If the flow of depreciation of capital is $D(t)$, then the change in the capital stock is the flow of *net* investment, $\dot{K}(t) = I(t) - D(t)$.

While time derivatives are very useful for many applications, it is often more helpful to measure the change in a variable over time as a *percentage* change or **growth rate** rather than as an *absolute* change. We convert differences into growth rates by dividing the change in the variable by the level of the variable. For time derivatives, this means that the (annualized) growth rate of a variable y at time t is $\dot{y}(t) / y(t)$. This ratio is a “pure number” such as 0.03 (or 3%) that can be compared directly with the growth rates of other variables.

Recalling the laws of derivatives, consider the time derivative of $\ln(y(t))$. Using the chain rule and the rule about derivatives of logarithms, $d \ln(y(t)) / dt = (1/y(t)) \cdot dy/dt = \dot{y}(t) / y(t)$. Thus, the slope of the time path of the logarithm of a variable is equal to the variable's growth rate. If a variable has a constant growth rate over time, the path of its *log* will be a straight line, which we discussed in section B of this chap-

ter. As we noted above, it is often more convenient to plot the path of a variable's log rather than plotting the level of a variable, since periods of faster and slower growth will show up readily to our eyes as regions with steeper and flatter slopes.

Growth rates of products, quotients, and powers

It often happens in growth theory that we know the growth rates of two variables, x and y , and we want to know the growth rate of another variable z that is a function of x and y . It turns out that there are very easy rules for the relationship between the growth rates if z is a product, quotient, or power function of x and y .

Rule 1. If $z = xy$, then $\dot{z}/z = \dot{x}/x + \dot{y}/y$. In words, the growth rate of a product is the sum of the growth rates of the variables being multiplied.

A common application is the relationship between nominal and real GDP. Nominal GDP is the product of real GDP and the GDP price index, so the growth rate of nominal GDP is the sum of the growth rate of real GDP and the GDP inflation rate.

To prove Rule 1, take the derivative of z with respect to time, using the product rule for derivatives described earlier. That gives us $\dot{z} \equiv dz/dt = \dot{x}y + \dot{y}x$. Dividing the left side of this equation by z and the right side by the equivalent expression xy yields $\dot{z}/z = \dot{x}y/xy + \dot{y}x/xy = \dot{x}/x + \dot{y}/y$.

Rule 2. If $z = x/y$, then $\dot{z}/z = \dot{x}/x - \dot{y}/y$. Again, in words, the growth rate of a quotient is growth rate of the numerator minus the growth rate of the denominator.

The same application can be used here. The GDP price index is nominal GDP divided by real GDP. Thus, the inflation rate is the difference between the growth rate of nominal GDP and the growth rate of real GDP. Rule 2 follows directly from Rule 1.

Rule 3. If $z = x^n$, where n is a constant, then $\dot{z}/z = n(\dot{x}/x)$. If one variable is equal to another variable raised to a power, then the growth rate of the first is the growth rate of the second times the power.

This rule has applications involving elasticities. The function $z = x^n$ is a constant-elasticity function with the elasticity of z with respect to x being equal to n . Thus, if x is growing at rate g_x and the elasticity of z with respect to x is n , then z will grow at rate ng_x .

Rule 3 can be proved using the chain rule and the rule for taking the derivative of a power. Differentiating with respect to time yields $\dot{z} = nx^{n-1}\dot{x}$. Dividing the left side by z and the right side by the equivalent expression x^n yields the result that $\dot{z}/z = n\dot{x}x^{n-1}/x^n = n\dot{x}/x$.

Multivariate functions and partial derivatives

All of the applications we have discussed above have related to situations in which the dependent variable under consideration (y) could be related to a single other variable (an independent variable x or time t). In most economic models, each variable is affected by many other variables, not just one. To analyze such models, we must extend the idea of a derivative to accommodate multiple variables.

The concept of a derivative is essentially bivariate: it involves a “dependent” variable whose change is in the numerator and an “independent” variable whose change is in the denominator. To use bivariate derivatives in a multivariate context, we examine the relationship between the dependent variable and each of several independent variables one at a time. In doing this, we explicitly assume that all of the independent variables other than the one we are currently examining do not change.

For example, suppose that production Y is assumed to depend on the levels of two inputs, labor L and capital K , according to a production function $Y = F(K, L)$. We can use the tools of calculus to examine the effect of an increase in capital on production holding labor constant (the *marginal product of capital*) or the effect of an increase in labor on production holding capital constant (the *marginal product of labor*). The *partial derivative* of the production function with respect to capital (labor) is defined to be the derivative of the production function taking capital (labor) as the independent variable and holding labor (capital) constant. We denote this partial derivative using the curly ∂ rather than d , e.g., $\partial Y/\partial K$ or $\partial Y/\partial L$, to signal that other variables are being held constant.

Since most economic relationships are multivariate, the partial derivative is used extensively in economic analysis. All the usual rules of differentiation that we studied above apply to partial derivatives as well. You must be careful, however, to remember which variables are allowed to change and which are being held constant.

Total differentials

When we are considering multivariate relationships among variables, the concept of the *total differential* is often useful. In the production function example, the level of output is related to the levels of the inputs by the production function $Y = F(K, L)$. The partial derivatives $\partial Y/\partial K$ and $\partial Y/\partial L$ measure how Y changes if *either* K or L changes, but what happens if *both* K and L change?

The total differential of Y , which we write as dY , relates the change in Y to changes in both K and L . The formula for the total differential is

$$dY = \frac{\partial Y}{\partial K} dK + \frac{\partial Y}{\partial L} dL,$$

where dK and dL represent changes in K and L . The total differential applies exactly only for infinitesimally small changes in K and L . We sometimes use the total differential to evaluate the relationships among the changes in variables in non-linear multiple-equation models where explicit solutions are impossible.

Multivariate maximization and minimization

Now that we have generalized the concept of the derivative to allow multiple independent variables, we can consider how to find the maximum and minimum of multivariate functions. A function such as $y = f(x)$ may reach a maximum or minimum only at a value of x where tiny changes in x have no effect on y . This occurs where the first derivative is zero: $f'(x) = 0$.

Similarly, a multivariate function $y = F(x_1, x_2)$ can have a maximum or minimum only where *both* partial derivatives are zero: $\partial y / \partial x_1 = 0$ and $\partial y / \partial x_2 = 0$. Geometrically, this means that the three-dimensional surface described by the function is flat looking both in the x direction and in the y direction.⁵

F. Understanding Romer's Chapter 1

Romer's Chapter 1 introduces you to the neoclassical growth model developed by Robert Solow and elaborated by him and many others in the 1950s and 1960s. The math in Chapter 1 is not very high-powered, but it contains some subtle applications that you may find tricky if you have not seen them before. This section will help you understand those points.

Manipulating the production function

On pages 10 through 12, Romer starts with the assumption that aggregate output depends on inputs of labor and capital and on an index of technology. He then moves rather quickly through some mathematical assumptions and manipulations that lead him to express the level of output per effective unit of labor input as a function of the capital/effective-labor ratio.

⁵ The first partial derivatives equaling zero is the "first-order condition" for a maximum or minimum. To be sure that a point at which the partial derivatives are zero is an extremum and to determine whether it is a maximum or a minimum requires second-order conditions. We will not be concerned with second-order conditions in this course—they are satisfied in all the models with which we shall work.

The initial assumption is that the production function has constant returns to scale. In mathematical terms, this condition is written as Romer's equation (1.2). Since c in equation (1.2) can be *any* positive number, we can choose a particular one. It turns out to be convenient to choose $c = 1/AL$, the reciprocal of the amount of "effective labor" in the economy.⁶ The reason that this is a convenient choice is that it implies that $F(K/AL, AL/AL) = F(K/AL, 1) = F(K, AL)/AL = Y/AL$. In words, this equation says that output per unit of effective labor depends *only* on the amount of capital per effective unit of labor. (If we ignore the presence of A for the moment, this says that output per worker depends only on how much capital each worker has to work with.) We simplify the notation by writing $y - f(k)$ rather than $Y/AL = F(K/AL, 1)$ with the $k \equiv K/AL$, $y \equiv Y/AL$, and $f(\bullet) \equiv F(\bullet, 1)$.

The partial derivatives of the production function have important economic interpretations. The marginal product of capital is defined to be the amount of additional output that can be obtained if the amount of capital input rises by one unit holding the amounts of the other inputs (labor, in this case) constant. Mathematically, this corresponds to the partial derivative: the amount by which the dependent variable changes when one of the independent variables changes by one unit with the others unchanged. Thus,

$$\text{MPK} = \partial Y / \partial K = \partial F(K, AL) / \partial K.$$

Similarly, the marginal product of labor is

$$\text{MPL} = \partial Y / \partial L = \partial F(K, AL) / \partial L.$$

On page 12, Romer shows that the marginal product of capital is equal to the first derivative of the intensive form of the production function, that is, $\text{MPK} = f'(k)$.⁷ Thus, the assumption that $f'(k) > 0$ is the natural economic assumption that capital's marginal product is positive—that more capital allows more output to be produced. The assumption that $f''(k) < 0$ asserts that as an economy gets more capital relative to (effective) labor, the marginal product of capital declines. This is nothing more or

⁶You can think of L as measuring the number of workers in the economy and A as measuring how effectively each worker works. The product AL is the amount of effective labor input. As we shall see below, AL grows for two reasons: the labor force usually expands over time with the population and each worker becomes more effective (or productive) as technology improves.

⁷This may seem obvious, but remember that MPK is $\partial Y / \partial K$, while $f'(k)$ is $\partial y / \partial k = \partial(Y/AL) / \partial(K/AL)$.

less than the standard microeconomic assumption of diminishing marginal returns, dressed up in fancy calculus clothes.

The “polar” Inada conditions discussed on page 12 also have easy intuitive interpretations. The condition that $\lim_{k \rightarrow 0} f'(k) = \infty$ says that as the capital/effective-labor ratio gets close to zero, the marginal product of capital gets extremely large. In other words, if workers have practically no tools at all, then an extremely large increase in production occurs if they acquire a small amount of capital. Similarly, the condition that $\lim_{k \rightarrow \infty} f'(k) = 0$ refers to the other extreme, when workers have huge amounts of capital. If the marginal product of capital goes to zero in this situation, it means that once a very large amount of capital is in place for each worker, additional units of capital eventually have only vanishingly small effects on production.

Both of the Inada conditions are natural extensions to extreme cases of the idea of diminishing marginal returns. The effect that they have on the production function shown in Romer’s Figure 1.1 is to assure that the slope of the curve at the origin is vertical and that if you follow the curve far enough to the right, it will become arbitrarily close to horizontal. These conditions (together with the assumption that the MPK is everywhere diminishing) assure that for any positive value r , there is some level of k at which the MPK is equal to r . The Inada conditions are important in assuring the existence and uniqueness of a steady-state equilibrium.

The Cobb-Douglas production function

Economists usually prefer to work at the greatest possible level of generality in order to assure that specific assumptions do not lead to conclusions that would not be valid in more general cases. For this reason, most of the analysis of the Solow model does not specify a particular functional form for the production function. However, sometimes we specify a particular form either because analysis in the general case is too difficult or in order to provide a specific example for expositional purposes.

In the short section on pages 12 and 13, Romer examines the properties of the *Cobb-Douglas production function*. This functional form is a workhorse of economics because it is one of the simplest functional forms having the basic properties that we require: constant returns to scale and positive but diminishing marginal products for the factors of production.

The constant-returns-to-scale Cobb-Douglas function with labor-augmenting or “Harrod-neutral” technological progress is written as in Romer’s equation (1.5):

$$Y = F(K, AL) = K^\alpha (AL)^{1-\alpha}, \quad (6)$$

where α is a parameter between zero and one. Romer's equation (1.6) shows that this function has constant returns to scale. In equation (1.7) he shows that the intensive form of the Cobb-Douglas is $f(k) = k^\alpha$.

Let's consider some other properties of the Cobb-Douglas that will be useful on the many occasions that we use it in this course. First of all, the marginal product of capital is the partial derivative of the production function with respect to capital. Thus,

$$MPK = \frac{\partial F}{\partial K} = \alpha K^{\alpha-1} (AL)^{1-\alpha} = \alpha \left(\frac{K}{AL} \right)^{\alpha-1} = \alpha k^{\alpha-1}. \quad (7)$$

To get the marginal product of labor, we differentiate with respect to L (not AL) to get

$$MPL = \frac{\partial F}{\partial L} = A(1-\alpha)K^\alpha (AL)^{-\alpha} = A(1-\alpha) \left(\frac{K}{AL} \right)^\alpha = A(1-\alpha)k^\alpha. \quad (8)$$

An interesting property of the Cobb-Douglas function emerges when we assume that each unit of labor and capital employed is paid an amount equal to its marginal product, as occurs under perfect competition and profit maximization. If this is the case, then the total amount paid to owners of capital is $MPK \times K$ and the share of total GDP paid to capital is $\alpha_K = (MPK \times K) / Y$. Using the first part of equation (7),

$$\alpha_K = \frac{\alpha K^{\alpha-1} (AL)^{1-\alpha} K}{K^\alpha (AL)^{1-\alpha}} = \alpha.$$

Similarly, labor's share α_L is

$$\alpha_L = \frac{A(1-\alpha)K^\alpha (AL)^{-\alpha} L}{K^\alpha (AL)^{1-\alpha}} = 1 - \alpha.$$

Thus, the exponents of capital and labor in the Cobb-Douglas function are the shares of GDP that they receive in competitive equilibrium. Since α and $1 - \alpha$ sum to one, the competitive payments to capital and labor exactly exhaust total GDP.

Another interesting property of the Cobb-Douglas coefficients α and $1 - \alpha$ is that they are the elasticities of output with respect to the two factors. The elasticity of output with respect to capital is defined as

$$\varepsilon_K \equiv \frac{\partial Y}{\partial K} \cdot \frac{K}{Y} = MPK \cdot \frac{K}{Y}.$$

Using our marginal product formula from (7) gives

$$\varepsilon_K = \alpha K^{\alpha-1} (AL)^{1-\alpha} \frac{K}{K^\alpha (AL)^{1-\alpha}} = \alpha.$$

Similar analysis shows that $\varepsilon_L = 1 - \alpha$.

Finally, it is sometimes convenient to represent the Cobb-Douglas function in log form rather than in levels. Taking the natural logs of both sides of (6) gives

$$\ln Y = \alpha \ln K + (1 - \alpha) \ln A + (1 - \alpha) \ln L.$$

Thus, the Cobb-Douglas is equivalent to a log-linear production function—the log of output is a linear function of the logs of the inputs.

The nature of growth equilibrium

The aim in these opening chapters is to characterize economic growth. Since sustained growth implies an ongoing process of change, we need to think about the kind of equilibrium that would be appropriate for a growth model. The equilibrium we seek will be a stable “growth path” for the main variables of the model rather than a fixed level. By stable, we mean that an economy will tend to converge to this equilibrium path over time and, once on the path, will proceed along it.

There are many different kinds of growth paths that could be stable equilibrium paths. We could have equilibrium paths with constant growth rates or ones on which growth rates increase, decrease, or oscillate over time. We could have equilibrium paths on which some or all of the major variables grow at the same rate or paths on which growth rates of variables differ.

Most of the simple growth models that we study in this course have equilibrium ***balanced-growth paths*** on which at least some of the major variables grow at the same, constant rate. This suggests two possible strategies for analyzing the equilibrium growth path.

For some models, we can find a balanced-growth path by looking for conditions under which the ratio of two variables is constant. For example, if K/AL is constant (*i.e.*, has a zero growth rate), then K and AL must be growing at the same rate because the growth rate of a quotient is the difference between the growth rates of the numerator and denominator. Thus a situation in which $\dot{k}/k = 0$, where $k \equiv K/AL$, is a candidate as a possible balanced-growth path.

Another possible approach to finding an equilibrium growth path is to examine situations in which the growth rate of one of our “level variables” is constant. So we might look for a situation in which $\dot{g}_K = 0$, where $g_K \equiv \dot{K}/K$. Each of these strategies will be useful to us in our growth analysis. In the basic Solow model of Romer’s

Chapter 1, the equilibrium growth path is most easily characterized by the $\dot{k}/k = 0$ condition, which is equivalent to the simpler condition $\dot{k} = 0$. The models of Romer's Chapter 3 will often be easier to characterize using the second strategy.

Basic dynamic analysis of k

On page 14, Romer presents three basic "equations of motion" for the two factors of production and the index of productivity. Equations (1.8) and (1.9) define the exogenous and constant growth rates for labor and productivity, n and g . These equations are equivalently expressed as (1.11) and (1.12) or as (1.13) and (1.14). Equation (1.15) defines the change in capital stock to be the difference between the flows of new gross investment and depreciation. Gross investment is assumed to equal saving, which is proportional to income: $sY(t)$. Depreciation is assumed to be proportional to the existing stock, $\delta K(t)$.

Following the strategy suggested above, we try to represent the model in terms of a variable that might be expected to approach a constant value on the equilibrium balanced-growth path: k . Thus, we are looking for an expression for the growth or change over time in k . One way of obtaining the solution is direct differentiation with respect to time. Romer shows how this is done on page 15.

An alternative derivation makes use of the growth-rate rules discussed above. Since $k = K/AL$, we can use Rules 1 and 2 to calculate its growth rate as $\dot{k}/k = \dot{K}/K - \dot{A}/A - \dot{L}/L$. The growth rates of technology and labor are assumed to be the constants g and n respectively, while the change in the capital stock is given by Romer's equation (1.15). Thus, $\dot{k}/k = (sY - \delta K)/K - g - n$. Multiplying both sides of this equation by k yields $\dot{k} = (sY/K)k - \delta k - gk - nk$. But $Y/K = y/k = f(k)/k$, so

$$\dot{k} = sf(k) - (n + g + \delta)k,$$

which is equivalent to Romer's equation (1.18).

If we knew the specific form of the production function f , we might be able to use methods of differential equations to solve this expression for a time path for k given some starting value $k(0)$. This would tell us the value of k at any time as a function of the initial value and time t . We shall need to do something like this in order to analyze how the model converges to the equilibrium balanced-growth path.

However, we can characterize the properties of the equilibrium path itself without choosing a specific production function and without resorting to such sophisticated mathematics. Romer's Figure 1.3 is a ***phase diagram*** that depicts the relationship between the level of k and its change. We define a ***steady-state*** equilibrium (or a balanced-growth path) to be a situation where the value of k (the ratio of capital to

quality-adjusted labor) is stable over time. Mathematically, we seek a solution in which $\dot{k} = 0$. From Figure 1.3, you can see that there is a unique level of k at which this steady-state equilibrium occurs. For lower values of k , there are economic forces that will cause it to increase; for higher values, these forces will decrease k . If the economy's capital/effective-labor ratio is the value k^* shown in Figure 1.3, then it is on a stable, steady-state balanced-growth path.

Using Taylor series to approximate the speed of convergence

Beginning on page 25, Romer examines the speed at which a Solow-model economy would converge to the steady-state path. We can only really solve for the path of convergence if we know the functional form of the production function. For example, if the production function is Cobb-Douglas, then one could use differential-equation methods to calculate a path of convergence for k and y given any starting values.

However, we would rather not tie ourselves down to one, specific functional form unless it is truly necessary. An alternative procedure is to *approximate* the behavior of the unspecified production function using the method of ***Taylor series***. A *first-order* Taylor-series approximation of a function around a specific value approximates the behavior of the function as a *linear* function of its variable. Since Taylor-series methods are often covered toward the end of a calculus sequence, we shall digress briefly to introduce the mathematical ideas behind them.

Suppose that two variables are related by a function $z = g(x)$, such as the one shown in Figure 9. We assume that the first, second, and higher derivatives of g are continuous functions at some chosen point x^* . Further suppose that z is equal to the value z^* when x is x^* . If g were a linear function, having a constant slope, then we could calculate the value of z corresponding to any value of x as

$$z = z^* + g'(x^*) (x - x^*). \quad (9)$$

This equation expresses the value of z as z^* (its value when x is x^*) plus the slope of the function at x^* times the difference between x and x^* . If g were a linear function, then the slope would be constant and equation (9) would give the exact value of z for any value of x . If g is not linear, then the slope changes and the actual function curves away from the straight-line approximation given by (9) as x moves away from x^* . Figure 9 shows how the linear approximation z_2 to the true value $z_1 = g(x_1)$ is calculated as z^* plus the vertical distance $g'(x^*) (x_1 - x^*)$, which is the height of the “slope triangle” to the right of (x^*, z^*) .

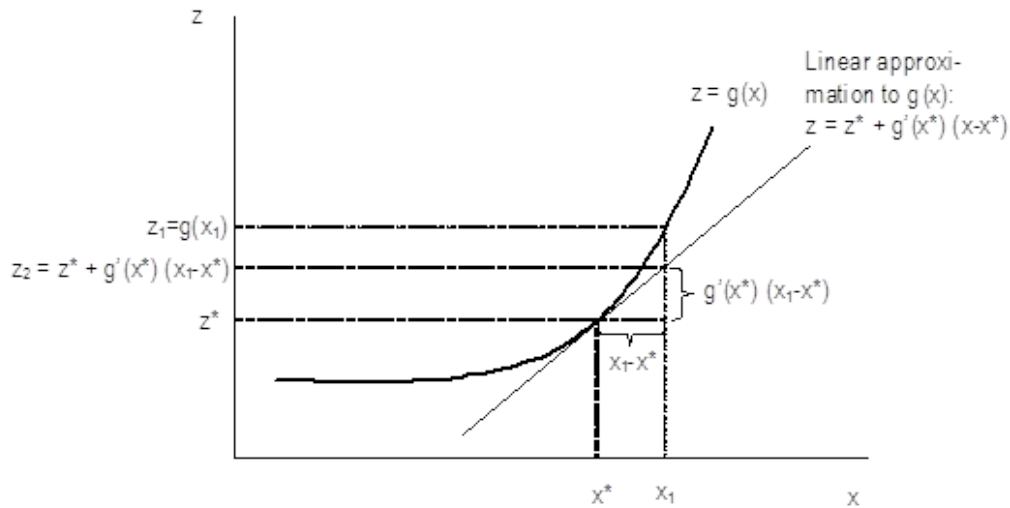


Figure 9. Taylor approximation

The linear approximation given by equation (9) and shown in Figure 9 is only a “first-order” approximation. A famous theorem of calculus called *Taylor’s Theorem* asserts that we can approximate any well-behaved function arbitrarily closely in the neighborhood around (x^*, z^*) by including higher and higher-order terms. For example, the second-order Taylor approximation of $g(x)$ in a neighborhood around x^* would be

$$z \cong z^* + g'(x^*)(x - x^*) + \frac{1}{2} g''(x^*)(x - x^*)^2. \quad (10)$$

Equation (10) approximates the g function as a parabola with both slope and curvature equal to those of g at the point (x^*, z^*) . The mathematical series that grows as the order of the approximation is increased is called a *Taylor series*. Equation (9) is called a first-order Taylor series and equation (10) is a second-order Taylor series.

For growth analysis, a first-order Taylor approximation is usually sufficient. Romer’s equation (1.28) on page 25 applies equation (9) to the phase diagram function of Romer’s Figure 1.3. The z variable is \dot{k} , the change in k , and the x variable is k . We choose k^* , the steady state value of k , as the specific value around which we approximate. We know that when $k = k^*$, its change is zero because we are in the steady state, so $\dot{k}(k^*) = 0$. Thus the point corresponding to (x^*, z^*) in Figure 9 is $(k^*, 0)$. The “ z^* ” term on the right-hand side of (9) does not appear in Romer’s equation (1.28) because it is zero.

To evaluate the derivative in (1.28), we differentiate the equation of motion (1.18) with respect to k and evaluate the resulting expression at the steady-state value k^* . The result of this differentiation is equation (1.31). Romer then denotes the negative of that derivative by λ ; it is just a constant number since it is evaluated at the steady-state point. Equation (1.29) shows that the gap between the current level of k and the steady-state level will decrease by a fraction approximately equal to λ each year. Using the formula for continuous growth (in this case, at a negative rate), we have $k(t) - k^* = e^{-\lambda t} [k(0) - k^*]$, where $k(0)$ is the value of the capital stock at which we begin the convergence process.

By appealing to some benchmark empirical estimates of the parameters of the model, Romer estimates λ to be approximately 4%. This means that 4% of the gap between actual and steady-state per-capita output (and capital per worker) will be eliminated in one year. Since this gap gets smaller over time, the absolute amount of change in k will diminish year-by-year during the convergence process (the flip-side of compound growth since the gap is diminishing rather than growing). This means that it would take about 18 years for one-half of the initial gap to be eliminated. He gets this number by noting that $e^{-0.04(18)} = 0.487 \cong 1/2$, so $k(18) - k^* \cong 1/2 [k(0) - k^*]$.

Growth models and the environment

In section 1.8, Romer presents a very basic introduction to how depletable resources and pollution can be introduced into the Solow model. The models Romer includes are a tantalizing introduction to a complex issue and should not be taken as the last word on the subject. (Then again, this warning could apply to almost everything in this course.)

The analysis of the natural resource model on pages 39–40 should be fairly straightforward if you have understood the basic mechanics of the Solow model. Note that the production function in Romer’s equation (1.41) has constant returns to scale in the four factors of production, and thus decreasing returns in labor and capital.

As Romer notes in the section titled “A Complication” starting on page 42, the assumption of a Cobb-Douglas production function is *not* an innocuous one in this case. The Cobb-Douglas (or any other production function) makes very specific assumptions about how production behaves as particular inputs become very scarce. We have little experience with entropic depletion of resources, so the reasonableness of the Cobb-Douglas for recent data should not endow us with great confidence that it is appropriate as resources run out. The literature on environmental effects in growth models is still young; no doubt many important theoretical results and empirical assessments will emerge in the coming decades.

G. Suggestions for Further Reading

Expositions of the Solow model

Solow, Robert M., "A Contribution to the Theory of Economic Growth," *Quarterly Journal of Economics* 70(2), February 1956, 65–94. (Solow's original exposition.)

Swan, Trevor W., "Economic Growth and Capital Accumulation," *Economic Record* 32, November 1956, 334–361. (An independent development of the same model by an Australian economist.)

Solow, Robert M., "Technical Change and the Aggregate Production Function," *Review of Economics and Statistics* 38(3), August 1957, 312–320. (Solow adds technical progress to the model and introduces the Solow residual.)

Barro, Robert, and Xavier Sala-i-Martin, *Economic Growth*, 2nd ed. (Cambridge, Mass.: MIT Press, 2004), Chapter 1. (A slightly more advanced exposition of the Solow model.)