

test for that defect (consequences)—except that the test is not a perfect indicator of the defect. How certain can we be, given our observations, that the underlying condition does in fact exist?

We can do a simple numerical calculation to answer the question for our particular example. Consider a population of 10,000 persons in which 100 (1%) have the defect and 9,900 do not. Suppose they all take the test. Of the 100 persons with the defect, the test will be (correctly) positive for 99. Of the 9,900 without the defect, it will be (wrongly) positive for 99. So we have 198 positive test results, of which one-half are right and one-half are wrong. If a random person receives a positive test result, it is just as likely to be because the test is right as because the test is wrong, so the risk that the defect is truly present for a person with a positive result is only 50%. (That is why tests for rare conditions must be designed to have especially low rates of false positives.)

For general questions of this type, we use an algebraic formula called *Bayes' theorem* to help us set up the problem and do the calculations. To do so, we generalize our example, allowing for two alternative underlying conditions, A and B (genetic defect or not), and two observable consequences, X and Y (positive or negative test result). Suppose that, in the absence of any information (over the whole population), the probability that A exists is p , so the probability that B exists is $(1 - p)$. When A exists, the chance of observing X is a , so the chance of observing Y is $(1 - a)$. [In the language of probability theory, a is the probability of X conditional on A , and $(1 - a)$ is the probability of Y conditional on A .] Similarly, when B exists, the chance of observing X is b , so the chance of observing Y is $(1 - b)$.

This description shows us that four alternative combinations of events could arise: (1) A exists and X is observed, (2) A

exists and Y is observed, (3) B exists and X is observed, and (4) B exists and Y is observed. The likelihood of these combined events can be determined using what is known in probability theory as the *modified multiplication rule*: the probability of two events both happening is equal to the probability of the first event times the probability of the second event conditional on the first event having happened. Using this rule, we find the probabilities of the four combinations to be, respectively, (1) pa , (2) $p(1 - a)$, (3) $(1 - p)b$, and (4) $(1 - p)(1 - b)$.

Now suppose that X is observed: A person has the test for the genetic defect and gets a positive result. Then we restrict our attention to a subset of the four preceding possibilities—namely, the first and the third, both of which include the observation of X . These two possibilities have a total probability of $pa + (1 - p)b$; this is the probability that X is observed. Within this subset of outcomes in which X is observed, the probability that A *also* exists is just pa , as we have already seen. So we know how likely we are to observe X alone and how likely it is that both X and A exist.

But we are more interested in determining how likely it is that A exists, given that we have observed X —that is, the probability that a person has the genetic defect, given that the test is positive. This calculation is the trickiest one. Using the modified multiplication rule, we know that the probability of both A and X happening equals the product of the probability that X alone happens times the probability of A conditional on X ; it is this last probability that we are after. Using the formulas for “ A and X ” and for “ X alone,” which we just calculated, the general formula

$\text{Prob}(A \text{ and } X) = \text{Prob}(X \text{ alone}) \times \text{Prob}(A \text{ conditional on } X)$
becomes

$pa = [pa + (1 - p)b] \times \text{Prob}(A \text{ conditional on } X)$; so,

$$\frac{pa}{pa + (1-p)b}.$$

Prob(A conditional on X) =

This formula gives us an assessment of the probability that A has occurred, given that we have observed X (and have therefore conditioned everything on this fact); it is known as *Bayes' theorem* (or rule or formula).

In our example of testing for the genetic defect, we had Prob(A) = p = 0.01, Prob(X conditional on A) = a = 0.99, and Prob(X conditional on B) = b = 0.01. We can use Bayes' theorem to compute Prob(A conditional on X), the probability that a genetic defect exists when the test comes back positive:

Prob(A conditional on X)

$$= \frac{(0.01)(0.99)}{(0.01)(0.99) + (1-0.01)(0.01)}$$

$$= \frac{0.0099}{0.0099 + 0.0099}$$

$$= 0.5$$

OBSERVATION	Sum of
<div> <div>X</div> <div>Y</div> </div> You may need to scroll left and right to see the full figure.	row

OBSERVATION				Sum of row
		X	Y	
TRUE CONDITION	A	pa	$p(1-a)$	p
	B	$(1 - p)b$	$(1 - p)(1 - b)$	$1 - p$
Sum of column		$pa + (1 - p)b$	$p(1 - a) + (1 - p)(1 - b)$	
You may need to scroll left and right to see the full figure.				

The probability algebra employing Bayes' theorem confirms the arithmetical calculation that we used earlier, which was based on an enumeration of all the possible cases. The advantage of the formula is that once we have it, we can apply it mechanically; this saves us the lengthy and error-susceptible task of enumerating every possibility and determining each of the necessary probabilities.

We show Bayes' theorem in Figure 9A.1 in tabular form, which may be easier to remember and to use than the preceding formula. The rows of the table show the alternative true conditions that might exist—for example, “genetic defect” and “no genetic defect.” Here, we have just two, A and B , but the method generalizes immediately to any number of possibilities. The columns show the observed events—for example, “test positive” and “test negative.”

Each cell in the table shows the overall probability of that combination of the true condition and the observation; these are just the probabilities for the four alternative combinations listed above. The last column on the right shows the sum across the first two columns for each of the top two rows. This sum is the total probability of each true condition (so, for instance, A 's probability is p , as we have seen). The last row shows the sum of the first two rows

in each column. This sum gives the probability that each observation occurs. For example, the entry in the last row of the X column is the total probability that X is observed, either when A is the true condition (a true positive in our genetic test example) or when B is the true condition (a false positive).

To find the probability of a particular condition, given a particular observation, then, Bayes' theorem says that we should take the entry in the cell corresponding to the combination of that condition and that observation and divide it by the column sum in the last row for that observation. As an example, $\text{Prob}(B \text{ conditional on } X) = (1 - p) b / [pa + (1 - p) b]$.

10 ■ The Prisoners' Dilemma and Repeated Games

IN THIS CHAPTER, we turn our attention to an in-depth analysis of the prisoners' dilemma game. As a classic example of the theory of strategy and its implications for predicting the behavior of game players, the prisoners' dilemma is familiar to anyone who has learned even just a little bit of game theory, and many who have never formally studied the subject know the basic story behind the game. The prisoners' dilemma is a game in which each player has a dominant strategy, but the equilibrium that arises when all players use their dominant strategies provides a worse outcome for every player than would arise if they all used their dominated strategies instead. The paradoxical nature of this equilibrium leads to several more complex questions about the nature of these interactions that only a more thorough analysis, like the one we provide in this chapter, can begin to address.

We introduced you to the prisoners' dilemma in [Section 3 of Chapter 4](#). There, we took note of the curious nature of an equilibrium that is actually a bad outcome for the players. The “prisoners” can find another outcome that both prefer to the equilibrium outcome, but they find it difficult to bring about. The focus of this chapter is the potential for achieving that better outcome. That is, we consider whether and how the players in a prisoners' dilemma can attain and sustain a mutually beneficial cooperative outcome, overcoming their separate incentives to defect for individual gain. We first review the standard prisoners' dilemma game, then develop some broad categories of solutions. One group of solutions involves changes in the order of moves, or in the way moves are made, that lead to increased opportunities for cooperation. Another group of solutions involves mechanisms or forces outside the game that can change the payoffs for

one or both players and can therefore change the equilibrium outcome. The use of strategic moves—specifically, promises (with warnings)—and repetition of the game are solutions of the first type. Mechanisms such as penalty (or reward) schemes, often offered by entities that are not players in the game, are solutions of the second type. We consider both types, but begin with the most studied solution, repeated play.

Prisoners' dilemma payoff structures are commonly observed in games that repeat on a regular basis. The pricing game between Xavier's Tapas Bar and Yvonne's Bistro, introduced in [Chapter 5](#), is an example. Firms typically set prices each week or month, so their interaction is a repeated one. The general theory of such repeated games was the contribution for which Robert Aumann was awarded the 2005 Nobel Prize in economics (jointly with Thomas Schelling). As usual at this introductory level, we look at only a few simple examples of the general theory, but use them to motivate broader conclusions about behavior in repeated games.

This chapter concludes with a discussion of some experimental evidence on player behavior in prisoners' dilemma games as well as several examples of actual dilemmas in action. Experiments that put live players in a variety of prisoners' dilemma games show some perplexing as well as some more predictable behavior; experiments conducted with the use of computer simulations yield additional interesting outcomes. Our examples of real-world dilemmas that end the chapter are provided to give a sense of the diversity of situations in which prisoners' dilemmas arise and to show at least one application of a solution method in action.

1 THE BASIC GAME (REVIEW)

Before we consider methods for avoiding the bad outcome in the prisoners' dilemma, we briefly review the basics of the game. Recall our example from [Chapter 4](#) of the husband and wife suspected of murder. Each is interrogated separately and can choose to confess to the crime or to deny any involvement. The payoff matrix that they face was originally presented as Figure 4.4 and is reproduced here as Figure 10.1. The numbers shown indicate years in jail; therefore, lower numbers are better for both players.

Both players here have a dominant strategy: Each does better to confess, regardless of what the other player does. The equilibrium outcome entails both players deciding to confess and each getting 10 years in jail. If they had both chosen to deny any involvement, however, they would both have been better off, with only 3 years of jail time to serve.

In a general prisoners' dilemma game, the two available strategies are often labeled as *Cooperate* and *Defect*. In Figure 10.1, Deny is the cooperative strategy and could be labeled *Cooperate* (with each other, not with the police); both players using that strategy yields a better outcome for both players than when both players use the other strategy. Confess is the defecting strategy and could be labeled *Defect*; when the players do not cooperate with each other, they choose to Confess in the hope of attaining individual gain at the other player's expense. Thus, players in a prisoners' dilemma can always be labeled, according to their choice of strategy, as either *defectors* or *cooperators*. We will use this labeling system throughout the discussion of potential solutions to the dilemma.

WIFE

		Confess (Defect)	WIFE Deny (Cooperate)
HUSBAND	Confess (Defect)	10 yr, 10 yr	1 yr, 25 yr
	Deny (Cooperate)	25 yr, 1 yr	3 yr, 3 yr
You may need to scroll left and right to see the full figure.			

FIGURE 10.1 Payoffs for the Standard Prisoners' Dilemma

We want to emphasize that, although we label one of the strategies here as Cooperate, the prisoners' dilemma game is *noncooperative*, in the sense explained in [Chapter 2](#), because players make their decisions and implement their choices individually. If the two players could discuss, choose, and play their strategies jointly—as if, for example, the prisoners were in the same room and could give a joint answer to the question of whether they were both going to confess—there would be no difficulty about their achieving the outcome that both would prefer. The essence of the questions of whether, when, and how a prisoners' dilemma can be resolved is the difficulty of achieving a cooperative (jointly preferred) outcome through noncooperative (individual) actions.

2 CHANGING THE WAY MOVES ARE MADE: REPETITION

Of all the mechanisms that can sustain cooperation in the prisoners’ dilemma, the best known and the most natural is [repeated play](#) of the game. Repeated or ongoing relationships among players impart special characteristics to the games that they play against one another. In the prisoners’ dilemma, this result manifests itself in the fact that each player fears that one instance of defecting will lead to a collapse of cooperation in the future. If the value of future cooperation is large and exceeds what can be gained in the short term by defecting, then the long-term individual interests of the players can automatically and tacitly keep them from defecting, without the need for any additional punishments or enforcement by third parties.

Consider the meal-pricing dilemma faced by the two restaurants, Xavier’ s Tapas Bar and Yvonne’ s Bistro, introduced in [Chapter 5](#). For our purposes here, we have chosen to simplify that game by supposing that only two choices of price are available: the jointly best (collusive) price of \$26 or the Nash equilibrium price of \$20. The payoffs (profits measured in hundreds of dollars per month) for each restaurant can be calculated by using the profit expressions derived in [Section 1.A](#) of [Chapter 5](#); these payoffs are shown in Figure 10.2. As in any prisoners’ dilemma, each restaurant has a dominant strategy to defect and price its meals at \$20, although both restaurants would prefer the outcome in which each cooperates and charges the higher price of \$26 per meal.

YVONNE’ S BISTRO	
\$20 (Defect)	\$26 (Cooperate)
You may need to scroll left and right to see the full figure.	

		YVONNE' S BISTRO	
		\$20 (Defect)	\$26 (Cooperate)
XAVIER' S TAPAS	\$20 (Defect)	288, 288	360, 216
	\$26 (Cooperate)	216, 360	324, 324
You may need to scroll left and right to see the full figure.			

FIGURE 10.2 Payoff Table for Restaurant Prisoners' Dilemma (in Hundreds of Dollars per Month)

Let us start our analysis by supposing that the two restaurants are initially in the cooperative mode, each charging the higher price of \$26. If one restaurant owner—say, Xavier—deviates from this pricing strategy, he can increase his profit from 324 to 360 (from \$32,400 to \$36,000) for one month. But then cooperation will have dissolved, and Xavier's rival, Yvonne, will see no reason to cooperate with him from then on. Once cooperation has broken down, presumably permanently, Xavier's profit is 288 (\$28,800) each month instead of the 324 (\$32,400) it would have been if he had not defected. By gaining 36 (\$3,600) in one month of defecting, he gives up 36 (\$3,600) each month thereafter by destroying cooperation. Even if the two restaurants' relationship lasts as little as three months, it seems that defecting is not in Xavier's best interest. A similar argument can be made for Yvonne's. Thus, if the two restaurants compete on a regular basis for at least three months, it seems that we might see cooperative behavior and high prices rather than the defecting behavior and low prices predicted by theory for the one-shot game.

A. Finite Repetition

But the solution of the dilemma is not actually that simple. What if the relationship did last exactly three months? Then strategic restaurants would want to analyze the full three-month game and choose their optimal pricing strategies. Each would use rollback to determine what price to charge each month. Starting their analyses with the third month, they would realize that, at that point, there was no future relationship to consider. Each restaurant would find that it had a dominant strategy: to defect. Given that, there would be effectively no future to consider in the second month either. Each player would know that there would be mutual defecting in the third month, and therefore both would defect in the second month; defecting would become the dominant strategy in that month, too. Then the same argument would apply to the first month as well. Knowing that both would defect in the second and third months anyway, neither player would see any future value in cooperation in the first month. Both players defect right from the start, and the dilemma is alive and well.

This result is a very general one. As long as the relationship between the two players in a prisoners' dilemma game lasts a fixed and known length of time, the Nash equilibrium with defecting should prevail in the last period of play. When the players arrive at the end of the game, there is never any value to continued cooperation, and so they defect. Then rollback predicts mutual defecting all the way back to the very first period of play. In practice, however, players in finitely repeated prisoners' dilemma games show a lot of cooperation; more on this to come.

B. Infinite Repetition

Analysis of the finitely repeated prisoners' dilemma shows that even repetition of the game cannot guarantee the players a solution to their dilemma. But what happens if the relationship does not have a predetermined length? What if the two restaurants in our example expect to continue competing with each other indefinitely? Then our analysis must change to incorporate this new aspect of their interaction, and we will see that the incentives of the players change too.

In repeated games of any kind, the sequential nature of the relationship means that players can adopt strategies that depend on behavior in preceding periods of play. Such strategies are known as [contingent strategies](#), and several specific examples are used frequently in the theory of repeated games. Of special note are the contingent strategies known as [trigger strategies](#). A player using a trigger strategy plays cooperatively as long as her rival(s) do so, but any defection on their part triggers a period of [punishment](#), of specified length, in which she defects in response. Two of the best-known trigger strategies are the grim strategy and tit-for-tat. The [grim strategy](#) entails cooperating with your rival until such time as she defects from cooperation; once she has done so, you punish her (by choosing the Defect strategy) on every play for the rest of the game. ¹[Tit-for-tat \(TFT\)](#) is not so harshly unforgiving as the grim strategy and is well known for its ability to solve the prisoners' dilemma without requiring permanent punishment. Playing TFT involves cooperating in the first period of play and then choosing, in each future period, the action chosen by your rival in the preceding period of play. Thus, when playing TFT, you cooperate with your rival if she cooperated during the most recent play of the game and you defect (as punishment) if she defected. The punishment phase lasts only as long as your rival continues to defect; you will return to cooperation one period after she chooses to do so.

Let us consider how play might proceed in the infinitely repeated restaurant pricing game if one of the players uses the contingent strategy tit-for-tat. We have already seen that if Xavier defects in one month, he can add 36 to his profits (which will be 360 instead of 324). But if Yvonne is playing TFT, then Xavier's defection induces Yvonne to punish him the next month in retaliation. At that point, Xavier has two choices. One option is to continue to defect by pricing his meals at \$20 and to endure Yvonne's continued punishment according to TFT; in this case, Xavier loses 36 (288 rather than 324) in every month thereafter in the foreseeable future. This option appears quite costly. But Xavier *can* get back to cooperation, too, if he so desires. By reverting to the cooperative price of \$26 after one month's defection, Xavier will incur only one month's punishment from Yvonne. During that month, Xavier will suffer a loss in profit of 108 (earning 216 rather than the 324 that he would have earned without any defection). In the second month after Xavier's defection, both restaurants could be back at the cooperative price, earning 324 each month. This one-time defection yields Xavier an extra 36 in profit, but costs him an additional 108 during the punishment.

It is important to realize here, however, that Xavier's extra 36 from defecting is gained in the first month, and his losses are incurred in the future. Therefore, the relative importance of the two depends on the relative importance of the present and the future. Here, because payoffs are calculated in dollar terms, an objective comparison can be made. Generally, money (or profit) that is earned today is better than money that is earned later, because even if you do not need (or want) the money until later, you can invest it now and earn a return on it until you need it. So Xavier should be able to calculate whether it is worthwhile to defect on the basis of the total rate of return on his investment (including interest and/or capital gains and/or dividends, depending on the type of investment). We use the symbol r to denote this rate of return. Thus, one dollar invested generates a return of r dollars, and 100 dollars generates $100r$; therefore, the rate of return is sometimes also said to be $(100r)\%$.

Note that we can calculate whether it is in Xavier's interest to defect because the restaurants' payoffs are given in dollar terms, rather than as ordinal rankings of outcomes. This means that payoff values in different cells are directly comparable: A payoff of 4 (dollars) is twice as good as a payoff of 2 (dollars) here, whereas in any two-by-two game in which the four possible outcomes are ranked from 1 (worst) to 4 (best), a payoff of 4 is not necessarily exactly twice as good as a payoff of 2. As long as the payoffs to the players are given in measurable units, we can calculate whether defecting in a prisoners' dilemma game is worthwhile.

I. IS IT WORTHWHILE TO DEFECT ONLY ONCE AGAINST A RIVAL PLAYING TFT? One of Xavier's options when playing repeatedly against a rival using TFT is to defect just once and then return to cooperating in subsequent periods. This particular strategy gains him 36 in the first month (the month during which he defects) but loses him 108 in the second month. By the third month, cooperation is restored. Is defecting for only one month worth it?

We cannot directly compare the 36 gained in the first month with the 108 lost in the second month because the monetary value of time must be incorporated into the calculation. That is, we need a way to determine how much the 108 lost in the second month is worth during the first month. Then we can compare that number with 36 to see whether defecting once is worthwhile. What we are looking for is the present value (PV) of 108, or how much profit earned this month (in the present) is equivalent to (has the same value as) 108 earned next month. We need to determine the number of dollars earned this month that, with interest, would give us 108 next month; we call that number PV, the present value of 108.

Given that the (monthly) total rate of return is r , getting PV this month and investing it until next month yields a total next month of $PV + rPV$, where the first term is the principal being paid back and the second term is the return. When the total is exactly 108, then PV equals the present value of 108. Setting $PV + rPV = 108$ yields a solution for PV:

$$PV = \frac{108}{1 + r}.$$

For any value of r , we can now determine the exact number of dollars that, earned this month, would be worth 108 next month.

From Xavier's perspective, the question remains whether the gain of 36 this month is offset by the loss of 108 next month. The answer depends on the value of PV. Xavier must compare the gain of 36 with the PV of the loss of 108. To defect once (and then return to cooperation) is worthwhile only if $36 > 108/(1 + r)$. This is the same as saying that defecting once is beneficial only if $36(1 + r) > 108$, which reduces to $r > 2$. Thus, Xavier should choose to defect once against a rival playing TFT only if the monthly total rate of return exceeds 200%. This outcome is very unlikely; for example, prime lending rates rarely exceed 12% per year, which translates into a monthly interest rate of no more than 1% (compounded annually, not monthly)—well below the 200% just calculated. So, if Yvonne is playing TFT, then it is better for Xavier to continue cooperating than to try a single instance of defecting.

II. IS IT WORTHWHILE TO DEFECT FOREVER AGAINST A RIVAL PLAYING TFT? What about the possibility of defecting once and then continuing to defect forever? This second option of Xavier's gains his restaurant 36 in the first month, but loses it 36 in every month thereafter into the future if the rival restaurant plays TFT. Whether such a strategy is in Xavier's best interest again depends on the present value of the losses incurred. But this time the losses are incurred over an [infinite horizon](#) of future months of competition.

Xavier's option of defecting forever against a rival playing TFT yields a payoff (profit) stream equivalent to what Xavier would get if he were to defect against a rival using the *grim strategy*. Recall that the grim strategy requires players to punish any defection with retaliatory defection in all future periods. In that case, it is not worthwhile for Xavier to attempt any return to cooperation after his initial defection because the rival firm

will be choosing to defect, as punishment, forever. Any defection on Xavier's part against a rival playing the grim strategy would then lead to a gain of 36 in the first month and a loss of 36 in all future months, exactly the same outcome as if he defected forever against a rival playing TFT. The analysis that follows is therefore the same analysis one would complete to assess whether it is worthwhile to defect at all against a rival playing the grim strategy.

To determine whether a defection of this type is worthwhile, we need to figure out the present value of all of the 36s that are lost in future months, add them all up, and compare them with the 36 gained during the month of defecting. The PV of the 36 lost during the first month of punishment and continued defecting on Xavier's part is just $36/(1+r)$; the calculation is identical to that used in [Section 2.B.1](#) to find that the PV of 108 was $108/(1+r)$. For the next month, the PV must be the dollar amount needed this month that, with two months of [compound interest](#), would yield 36 in two months. If the PV is invested now, then in one month the investor would have that principal amount plus a return of rPV , for a total of $PV + rPV$, as before; leaving this total amount invested for the second month means that at the end of two months, the investor has the amount invested at the beginning of the second month ($PV + rPV$) plus the return on that amount, which would be $r(PV + rPV)$. The PV of the 36 lost two months from now must then solve the equation $PV + rPV + r(PV + rPV) = 36$. Working out the value of PV here yields $PV(1+r)^2 = 36$, or $PV = 36/(1+r)^2$. You should see a pattern developing. The PV of the 36 lost in the third month of continued defecting is $36/(1+r)^3$, and the PV of the 36 lost in the fourth month is $36/(1+r)^4$. In fact, the PV of the 36 lost in the n th month of continued defecting is just $36/(1+r)^n$. Xavier loses an infinite sum of 36s, and the PV of each of them gets smaller each month.

More precisely, Xavier loses the sum, from $n = 1$ to $n = \infty$ (where n is the number of months of continued defecting after the initial month, which is month 0), of $36/(1+r)^n$. Mathematically, it is written as the sum of an infinite number of terms:²

$$\frac{36}{1+r} + \frac{36}{(1+r)^2} + \frac{36}{(1+r)^3} + \frac{36}{(1+r)^4} + \dots$$

Because r is a rate of return and presumably a positive number, the ratio $1/(1+r)$ will be less than 1; this ratio is generally called the [discount factor](#) and is denoted by the Greek letter δ . With $\delta = 1/(1+r) < 1$, the mathematical rule for infinite sums tells us that this sum converges to a specific value, in this case $36/r$. (The appendix to this chapter contains a detailed discussion of the solution of infinite sums.)

It is now possible to determine whether Xavier will choose to defect forever. He compares his restaurant's gain of 36 with the PV of all the lost 36s, or $36/r$. Then he defects forever only if $36 > 36/r$, or $r > 1$; defecting forever is beneficial in this particular game only if the monthly rate of return exceeds 100%, another unlikely event. Thus, we would not expect Xavier to defect against a cooperative rival when both are playing tit-for-tat. (Nor would we expect defection against a cooperative rival when both are playing the grim strategy.) When both Yvonne and Xavier play TFT, the cooperative outcome in which both set a high price is a Nash equilibrium outcome. Both playing TFT is a Nash equilibrium, and use of this contingent strategy solves the prisoners' dilemma for the two restaurant owners.

Remember that tit-for-tat is only one of many trigger strategies that players could use in repeated prisoners' dilemmas. And it is one of the "nicer" ones. Thus, if TFT can be used to solve the dilemma for the two restaurant owners, other, harsher trigger strategies should be able to do the same. As noted, the grim strategy can also be used to sustain cooperation in this infinitely repeated game, and in others as well.

C. Games of Unknown Length

In addition to considering games of finite or infinite length, we can incorporate a more sophisticated tool to deal with games of unknown length. It is possible that, in some repeated games, players will not know for certain exactly how long their interaction will continue. They may, however, have some idea of the *probability* that the game will continue for another period. For example, our restaurant owners might believe that their repeated competition will continue only as long as their customers find prix fixe menus to be the dining-out experience of choice; if there were some probability each month that à la carte dinners would take over that role, then the nature of the game would be altered.

Recall that the present value of a loss next month is already worth only $\delta = 1/(1+r)$ times the amount earned. If, in addition, there is only a probability p (less than 1) that the relationship will actually continue to the next month, then next month's loss is worth only p times δ times the amount lost. For Xavier's Tapas Bar, this means that the PV of the 36 lost with continued defecting is $36 \times \delta$ [the same as $36/(1+r)$] when the game is assumed to be continuing with certainty, but is only $36 \times p \times \delta$ when the game is assumed to be continuing with probability p . Incorporating the probability that the game may end next period means that the present value of the lost 36 is smaller, because $p < 1$, than it is when the game is definitely expected to continue (when p is *assumed* to equal 1).

The effect of incorporating p is that we now effectively discount future payoffs by the factor $p \times \delta$ instead of simply by δ . We can then generate an effective rate of return, R , where $1/(1+R) = p \times \delta$, and R depends on p and δ as shown:³

$$\frac{1}{1 + R} = p\delta$$

$$1 = p\delta(1 + R)$$

$$R = \frac{1 - p\delta}{p\delta}.$$

With a 5% actual rate of return on investments ($r = 0.05$, and so $\delta = 1/1.05 = {}^{20}/_{21}$) and a 50% chance that the game continues for an additional month ($p = 0.5$), then $R = [1 - (0.5)({}^{20}/_{21})]/(0.5)({}^{20}/_{21}) = 1.1$, or 110%.

The high rates of return required to destroy cooperation (encourage defection) in these examples now seem more realistic, if we interpret them as effective rather than actual rates of return. It becomes conceivable that defecting forever, or even once, might actually be to one's benefit if there is a large enough probability that the game will end in the near future. Consider Xavier's decision whether to defect forever against a TFT-playing rival. Our earlier calculations showed that permanent defecting is beneficial only when r exceeds 1, or 100%. If Xavier's faces the 5% actual rate of return and the 50% chance that the game will continue for an additional month, as we assumed in the preceding paragraph, then the effective rate of return of 110% will exceed the critical value needed for him to continue defecting. Thus, the cooperative behavior sustained by the TFT strategy can break down if there is a sufficiently large chance that the repeated game might be over by the end of the next period of play—that is, when the value of p is sufficiently small.

D. General Theory

We can easily generalize these ideas about when it is worthwhile to defect against TFT-playing rivals so that you can apply them to any prisoners’ dilemma game that you encounter. To do so, we can use a table with general payoffs (delineated in appropriately measurable units) that satisfy the standard structure of payoffs in a prisoners’ dilemma, as in Figure 10.3. The payoffs in the table must satisfy the relation $H > C > D > L$ for the game to be a prisoners’ dilemma, where C is the payoff in the *cooperative* outcome, D is the payoff when both players *defect* from cooperation, H is the *high* payoff that goes to the defector when one player defects while the other cooperates, and L is the *low* payoff that goes to the cooperator in the same situation.

		COLUMN	
		Defect	Cooperate
ROW	Defect	$D, \textcolor{teal}{D}$	$H, \textcolor{teal}{L}$
	Cooperate	$L, \textcolor{teal}{H}$	$C, \textcolor{teal}{C}$

FIGURE 10.3 General Version of the Prisoners’ Dilemma

In this general version of the prisoners’ dilemma, a player’ s one-time gain from defecting is $(H - C)$. The single-period loss for being punished while you return to cooperation is $(C - L)$, and the per-period loss for perpetual defecting is $(C - D)$. To be as general as possible, we will allow for situations in which there is a probability $p < 1$ that the game continues beyond the next period, and so we will discount payoffs using an effective rate of return of R per period. If $p = 1$, as would be the case when the game is guaranteed to continue, then $R = r$, the rate of return used in our preceding calculations. Replacing r with R , we find that the results attained earlier generalize almost immediately.

We found earlier that a player defects exactly once against a rival playing TFT only if the one-time gain from defecting $(H -$

C) exceeds the present value of the single-period loss from being punished (the PV of $C - L$). In this general game, that means that a player defects once against a TFT-playing opponent only if $(H - C) > (C - L)/(1 + R)$, or $(1 + R)(H - C) > C - L$, or

$$R > \frac{C - L}{H - C} - 1.$$

Similarly, we found that a player defects forever against a rival playing TFT only if the one-time gain from defecting exceeds the present value of the infinite sum of the per-period losses from perpetual defecting (where the per-period loss is $C - D$). In the general game, then, a player defects forever against a TFT-playing opponent, or defects at all against a grim strategy-playing opponent, only if $(H - C) > (C - D)/R$, or

$$R > \frac{C - D}{H - C}.$$

The three critical elements in a player's decision to defect, as seen in these two expressions, are the immediate gain from defection ($H - C$), the future losses from punishment ($C - L$ or $C - D$ per period of punishment), and the effective rate of return (R , which measures the importance of the present relative to the future). Defecting from cooperation becomes more attractive when the immediate gain from defection is higher, the future losses from punishment are smaller, and the effective rate of return is higher. But how high or low do these three quantities need to be, exactly, in order for players to prefer defecting?

First, assume that the values of the gains and losses from defecting are fixed. Then changes in R determine whether a player defects, and defection is more likely when R is large. Large values of R are associated with small values of p and small values of δ (and large values of r), so defection is more likely when the probability of continuation is low or the discount factor is low (or the interest rate is high). Another way to

think about it is that defection is more likely when the future is less important than the present, or when there is little future to consider; that is, defection is more likely when players are impatient or when they expect the game to end quickly.

Second, consider the case in which the effective rate of return is fixed, as is the one-period gain from defecting. Then changes in the per-period loss associated with punishment determine whether defecting is worthwhile. Here, it is smaller values of $C - L$ or $C - D$ that encourage defection. In this case, defection is more likely when punishment is not very severe.⁴

Finally, assume that the effective rate of return and the per-period loss associated with punishment are held constant. Now players are more likely to defect when the gains, $H - C$, are high. This situation is more likely when defecting garners a player large and immediate benefits.

This discussion also highlights the importance of the detection of defecting. Decisions about whether to continue along a cooperative path depend on how long defecting might be able to go on before it is detected, on how accurately it is detected, and on how long any punishment can be made to last before an attempt is made to revert to cooperation. Although our model does not incorporate these considerations explicitly, if defecting can be detected accurately and quickly, its benefit will not last long, and the subsequent cost will have to be paid more surely.

Therefore, the success of any trigger strategy in resolving a repeated prisoners' dilemma depends on how well (in terms of both speed and accuracy) players can detect defecting. This is one reason why the TFT strategy is often considered dangerous: Slight errors in the execution of actions or in the perception of those actions can send players into continuous rounds of punishment from which they may not be able to escape for a long time, until a slight error of the opposite kind occurs.

You can use all of these ideas to guide you in when to expect more cooperative behavior between rivals and when to expect more defecting and cutthroat actions. If times are bad and an entire

industry is on the verge of collapse, for example, so that businesses feel that there is no future, competition may become fiercer (less cooperative behavior may be observed) than in normal times. Even if times are temporarily good, but good conditions are not expected to last, firms may want to make a quick profit while they can, so cooperative behavior may again break down. Similarly, in an industry that emerges temporarily because of a quirk of fashion and is expected to collapse when fashion changes, we should expect less cooperation. Thus, a particular beach resort might become the place to go, but all the hotels there will know that such a situation cannot last, and that they cannot afford to collude on pricing. If, in contrast, the shifts in fashion are among products made by an unchanging group of companies in long-term relationships with one another, cooperation might persist. For example, even if all the children want cuddly bears one year and Transformers Rescue Bots the next, collusion in pricing may occur if the same small group of manufacturers makes both items.

In [Chapter 11](#), we will look in more detail at prisoners' dilemmas that arise in games with many players. We will examine when and how players can overcome such dilemmas and achieve outcomes better for them all.

Endnotes

- Defecting as retaliation under the requirements of a trigger strategy is often termed *punishing* to distinguish it from the original decision to defect. [Return to reference 1](#)
- The key formula provided in the appendix to this chapter states that, for any x , the infinite sum

$$\frac{x}{1+r} + \frac{x}{(1+r)^2} + \frac{x}{(1+r)^3} + \dots = \frac{x}{r}.$$

[Return to](#)

[reference 2](#)

- We could also express R in terms of r and p , in which case $R = (1 + r)/p - 1$. [Return to reference 3](#)
- The costs associated with defection may also be smaller if information transmission is not perfect, as might be the case if there are many players, so that difficulties might arise in identifying the defector and in coordinating a punishment scheme. Similarly, gains from defection may be larger if rivals cannot identify a defection immediately. [Return to reference 4](#)

Glossary

repeated play

A situation where a one-time game is played repeatedly in successive periods. Thus, the complete game is mixed, with a sequence of simultaneous-move games.

contingent strategy

In repeated play, a plan of action that depends on other players' actions in previous plays. (This is implicit in the definition of a strategy; the adjective "contingent" merely reminds and emphasizes.)

trigger strategy

In a repeated game, this strategy cooperates until and unless a rival chooses to defect, and then switches to noncooperation for a specified period.

punishment

We reserve this term for costs that can be inflicted on a player in the context of a repeated relationship (often involving termination of the relationship) to induce him to take actions that are in the joint interests of all players.

grim strategy

A strategy of noncooperation forever in the future, if the opponent is found to have cheated even once. Used as a threat of punishment in an attempt to sustain cooperation.

tit-for-tat (TFT)

In a repeated prisoners' dilemma, this is the strategy of [1] cooperating on the first play and [2] thereafter doing each period what the other player did the previous period.

present value (PV)

The total payoff over time, calculated by summing the payoffs at different periods each multiplied by the appropriate discount factor to make them all comparable with the initial period's payoffs.

infinite horizon

A repeated decision or game situation that has no definite end at a fixed finite time.

compound interest

When an investment goes on for more than one period, compound interest entails calculating interest in any one period on the whole accumulation up to that point, including not only the principal initially invested but also the interest earned in all previous periods, which itself involves compounding over the period previous to that.

discount factor

In a repeated game, the fraction by which the next period's payoffs are multiplied to make them comparable with this period's payoffs.

effective rate of return

Rate of return corrected for the probability of noncontinuation of an investment to the next period.

3 CHANGING THE ORDER OF MOVES: PROMISES

A key aspect of the classic prisoners' dilemma game is that players make their moves simultaneously, deciding whether to cooperate (by not confessing, by setting a high price, etc.) without being able to observe the other player's choice. However, many games with a prisoners' dilemma structure unfold sequentially, with one player irreversibly and observably making a move before the other. In this section, we focus on such *sequential-move prisoners' dilemma games*, showing how players can achieve the mutually preferred outcome in which both cooperate as long as either of them can be trusted.

Stanford economist Avner Greif studied the sequential-move prisoners' dilemma (also referred to as a "one-sided prisoners' dilemma") in the historical context of *exchange* in the medieval world at marketplaces known as bazaars.⁵ Once a buyer and seller found each other and agreed on the terms of trade, who would be the first to complete their side of the bargain? If the seller first hands over the product, the buyer might run away without paying. On the other hand, if the buyer first hands over the money, the seller might refuse to hand over the product, denying that any money was paid or, even more insidiously, might hand over a defective product whose defects the buyer will discover only after it is too late to return for a refund.

Greif refers to this strategic challenge as "the fundamental problem of exchange" and argues that societies' abilities to resolve this problem (or not) contributed significantly to their growth and prosperity (or lack thereof) during the medieval period.⁶ To understand the challenge, consider the sequential-move game shown in Figure 10.4, in which the buyer (Bob) first decides whether to pay, and then the seller (Ann) decides whether to hand over the product. Each benefits individually by choosing

to *break the trust*—not following through on their own end of the bargain—but both are better off when both *keep the trust*. The players' payoffs in this game are exactly the same as in a standard prisoners' dilemma, but now one player moves first. (You can compare the payoff structure here with the ordinal payoffs from Figure 4.10 to confirm that this game is a prisoners' dilemma.) As the second mover, Ann's dominant strategy is breaking the trust; anticipating this, Bob will choose to break the trust as well. Thus, in the rollback equilibrium, both players break the trust, and no trade gets made; this is exactly what players do in the Nash equilibrium in the standard simultaneous-move version of the prisoner's dilemma.

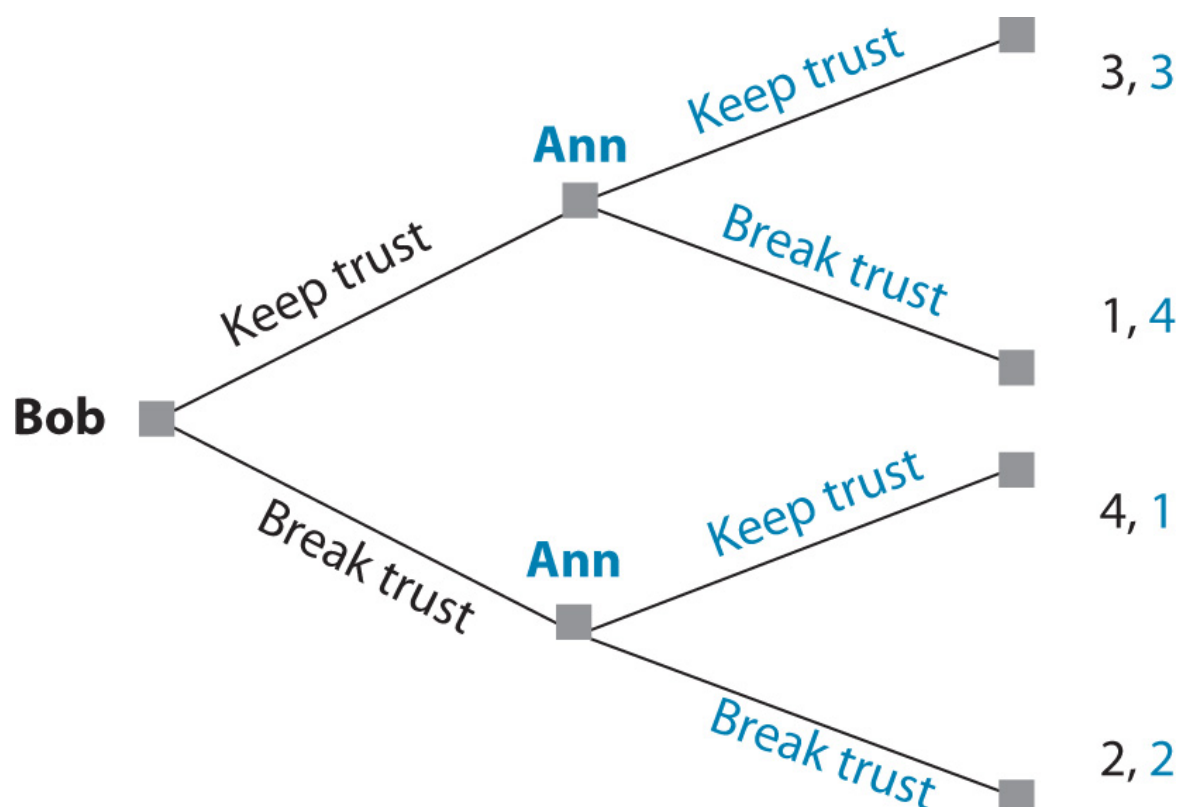


Figure 10.4 Game Tree for a Sequential-Move Prisoners' Dilemma

The problem here is less severe than in the simultaneous-move game, however, because only the second mover needs to establish her trustworthiness in order to solve the dilemma. In particular,

suppose that Ann can *make a promise*, declaring that she will keep the trust (as the second mover) if Bob first keeps the trust himself. (Ann also implicitly warns Bob that she will break the trust if he does.) So long as Ann's promise is credibly made, Bob will prune the branch of the tree in which Ann breaks the trust after he keeps the trust. Bob's choice is then between the outcomes in which both keep the trust or both break it, of which he prefers for both to keep the trust.

To summarize, the fundamental problem of exchange can be solved, so long as *either* side in an exchange can be trusted, by allowing the trusted party to move last in the upcoming game and to make a credible promise in the pregame. Establishing a reputation for trustworthiness is therefore enormously valuable for sellers, as it allows them to make profitable trades where other less trusted sellers cannot. Viewed this way, trust is an asset—and one that can be loaned as well. In medieval bazaars, merchants would typically appoint one of their own to serve as judge to resolve disputes and impose penalties for misbehavior. So long as the judge had a reputation for fairness (to punish wrongdoers) and discernment (to learn the truth in any situation), no *other* merchant would dare to cheat a buyer, knowing that he would face punishment far more costly than whatever gain he could have gotten by cheating. In this way, the trustworthiness of one merchant (appointed as judge) could be extended to all of them.

Such solutions worked in medieval times, allowing trade to flourish even on the frontiers of civilization, but what about today, in the new frontier of electronic commerce? On eBay, for example, buyer *feedback* provides a means by which sellers can establish their trustworthiness. Unfortunately, unscrupulous sellers can easily create false positive feedback for themselves and, even after getting kicked off the site, can purchase on the dark Web a new “verified, ready to sell” eBay identity with a sterling *false* history of positive transactions. (For more on the market for false eBay identities, see Exercise S9.)

To give consumers confidence that they will not be cheated, eBay offers a limited money-back guarantee for buyers who don't receive an item or receive one that is not as described. Such a

guarantee ensures that consumers don't suffer from the most obvious sorts of fraud, such as when a man in Texas fleeced 46 customers of \$191,000 in a "hot-tub scam" in which he never delivered the promised tubs.⁷ However, in subtler cases, when the seller has delivered an item that is not *quite* as described, or is defective in some way that is itself difficult to prove, the buyer may be unable to prove her case—especially if the seller claims that the buyer is herself attempting to commit fraud (or extortion) by submitting a false claim. Clearly, trust on eBay and other e-commerce marketplaces is a work in progress; the fundamental problem of exchange remains.

Endnotes

- Avner Greif, *Institutions and the Path to the Modern Economy* (New York: Cambridge University Press, 2006). See also Avinash Dixit, *Lawlessness and Economics: Alternative Modes of Governance* (Princeton, N.J.: Princeton University Press, 2007). [Return to reference 5](#)
- Thomas Hobbes considered the same problem in 1651, writing in his classic treatise *Leviathan*, “If a Covenant be made, wherein neither of the parties performe presently, but trust one another . . . upon any reasonable suspition, it is Voyd. . . . For he that performeth first, has no assurance that the other will performe after; because the bonds of words are too weak to bridle mens . . . avarice.” See Thomas Hobbes, *Leviathan* (Urbana, Ill.: Project Gutenberg, 2009), Chapter 14, Section 17; retrieved April 30, 2019, from www.gutenberg.org/ebooks/3207. [Return to reference 6](#)
- Mary Flood, “Prison for Houston Man Who Ran eBay Hot Tub Scam,” *Houston Chronicle*, May 19, 2010. [Return to reference 7](#)

4 CHANGING PAYOFFS: PENALTIES AND REWARDS

Changes in the way moves are made and changes in the order of moves are the major vehicles for the solution of the prisoners’ dilemma. But there are several other mechanisms that can achieve the same purpose by changing payoffs for the players. One of the simplest such mechanisms, which averts the prisoners’ dilemma in the one-shot version of the game, arises when an agent or entity outside the game has the power to inflict some direct [penalty](#) on the players when they defect. Once the payoffs have been altered to incorporate the cost of the penalty, players may find that the dilemma has been resolved.

		WIFE	
		Confess	Deny
HUSBAND	Confess	10 yr, 10 yr	21 yr, 25 yr
	Deny	25 yr, 21 yr	3 yr, 3 yr

FIGURE 10.5 Prisoners’ Dilemma with Penalty for the Lone Defector

Consider the husband-wife dilemma from [Section 1](#) of this chapter. If only one player defects, the game’ s outcome entails 1 year in jail for the defector and 25 years for the cooperator (see Figure 10.1). The defector, though, upon getting out of jail early, might find the cooperator’ s friends waiting outside the jail. The physical harm caused by those friends might be equivalent to an additional 20 years in jail. If so, and if the players account for the possibility of this harm, then the payoff structure of the original game has changed.

The new game, with the physical penalty included in the payoffs, is illustrated in Figure 10.5. With the 20 years in jail added to each player's sentence when that player confesses while the other denies, the game is completely different. Best-response analysis shows that there are now two pure-strategy Nash equilibria: One of them is (Confess, Confess); the other is (Deny, Deny). Now each player finds that it is in his or her best interest to cooperate if the other is going to do so. The game has changed from a prisoners' dilemma to an assurance game, which we also studied in [Chapter 4](#). Solving the new game requires selecting an equilibrium from the two that exist. One of them—the cooperative outcome—is clearly better than the other from the perspective of both players. Therefore, it may be easy to sustain it as a focal point if some convergence of expectations can be achieved.

Notice that the penalty in this scenario is inflicted on a defector only when his or her rival does *not* defect. However, stricter penalty mechanisms can be incorporated into the prisoners' dilemma, including ones in which there are penalties for *any* defection. Such discipline typically must be imposed by a third party with some power over *both* players, rather than by the other player's friends, because the friends would have little authority to penalize the first player for defecting when their associate also defects. If both prisoners are members of a special organization (such as a gang or a crime mafia), and the organization has a standing rule of never confessing to the police under penalty of extreme physical harm, then the game changes again to the one illustrated in Figure 10.6.

		WIFE	
		Confess	Deny
HUSBAND	Confess	30 yr, 30 yr	21 yr, 25 yr
	Deny	25 yr, 21 yr	3 yr, 3 yr

FIGURE 10.6 Prisoners' Dilemma with Penalty for Any Defecting

Now, the equivalent of an additional 20 years in jail is added to *all* payoffs associated with the Confess strategy (compare Figure 10.6 with Figure 10.1). In the new game, each player has a dominant strategy, as in the original game. The difference is that the change in the payoffs makes Deny the dominant strategy for each player. And (Deny, Deny) becomes the unique pure-strategy Nash equilibrium. The stricter penalty scheme achieved with the third-party enforcement mechanism makes defecting so unattractive to players that the cooperative outcome becomes the new equilibrium of the game.

Larger prisoners' dilemma games are more difficult to solve with mechanisms that change payoffs by exacting penalties for defection. In particular, if there are many players and some uncertainty exists, penalty schemes may be more difficult to maintain. It becomes harder to decide whether actual defecting is taking place or if it's just bad luck or a mistaken move. In addition, if there really is defecting, it is often difficult to determine the identity of the defector from among the larger group. And if the game is a one-shot one, there is no opportunity in the future to correct a penalty that is too severe or to inflict a penalty once a defector has been identified. Thus, penalties may be less successful in large one-shot games than in the two-person games we consider here. We study prisoners' dilemmas with a large number of players in greater detail in [Chapter 11](#).

A further interesting possibility arises when a prisoners' dilemma that has been solved with a penalty mechanism is considered in the context of the larger society in which the game is played. It might be the case that, although the equilibrium outcome is bad for the players, it is actually good for the rest of society, or for some subset of persons within the rest of society. If so, social or political

pressures might arise to try to minimize the ability of players to break out of the dilemma. When third-party penalties are the solution to a prisoners' dilemma, as is the case with crime mafias that enforce a no-confession rule, society can come up with its own strategy to reduce the effectiveness of the penalty mechanism. The U.S. Federal Witness Protection Program, in which the government removes the threat of penalty in return for confessions and testimony in court, is an example of a system that has been set up for just this purpose.

Similar opportunities for changing payoffs can be seen in other prisoners' dilemmas, such as the pricing game between our two restaurants. The equilibrium in that case entails both restaurants charging the low price of \$20, even though they would enjoy higher profits when charging the higher price of \$26. Although the restaurants want to break out of this bad equilibrium—and we have already seen how the use of trigger strategies could help them do so—their customers are happier with the low price offered in the Nash equilibrium of the one-shot game. The customers, then, have an incentive to try to destroy the efficacy of any enforcement mechanism or solution process the restaurants might use. For example, because some firms facing pricing games attempt to solve the dilemma through the use of price-matching campaigns, customers might want to press for legislation banning such practices. We analyze the effects of price-matching strategies in [Section 7.B](#).

Just as a prisoners' dilemma can be resolved by way of mechanisms that reduce the payoffs for defectors, it can also be resolved by mechanisms that increase the payoffs for, or reward, cooperators. Because such solutions are more difficult to implement in practice, we mention them only briefly.

The most important question is who is to pay the rewards. If it is a third party, that person or group must have

sufficient interest of its own in the cooperation achieved by the players to make it worth its while to pay out the rewards. A rare example of this solution approach was used when the United States brokered the Camp David Accords between Israel and Egypt by offering large promises of aid to both. If the rewards are to be paid by the players themselves to each other, the trick is to make the rewards contingent (paid out only if the other player cooperates) and credible (guaranteed to be paid if the other player cooperates). Meeting these criteria may require some unusual arrangements; Exercise U6 shows one such example.

Glossary

penalty

We reserve this term for one-time costs (such as fines) introduced into a game to induce the players to take actions that are in their joint interests.

5 CHANGING PAYOFFS: LEADERSHIP

The final type of situation we consider involves a prisoners’ dilemma in which the standard payoff structure is changed due to differences in size between the two players. These situations are frequently ones in which one player takes on the role of leader in the interaction. In most examples of the prisoners’ dilemma, the game is assumed to be symmetric; that is, all the players stand to lose (or gain) the same amount from defecting (or cooperating). However, in actual strategic situations, one player may be relatively “large” (a leader) and the other “small.” If the sizes of the payoffs are unequal enough, so much of the harm from defecting may fall on the larger player that she acts cooperatively, even while knowing that the other will defect. Saudi Arabia, for example, played such a role as the “swing producer” in OPEC (Organization of Petroleum Exporting Countries) for many years; to keep oil prices high, it cut back on its output when one of the smaller producers, such as Libya, expanded its production.

SOPORIA			
		Participate	Not
DORMINICA	Participate	-1, -1	-2, 0
	Not	0, -2	-1.6, -1.6
You may need to scroll left and right to see the full figure.			

FIGURE 10.7 Payoffs for Equal-Population SANE Research Game (in Billions of Dollars)

As with the OPEC example, [leadership](#) tends to be observed more often in games between nations than in games between firms or individual persons. Thus, our example of a game in which leadership may be used to solve the prisoners' dilemma is one played between countries. Imagine that the populations of two countries, Dorminica and Soporina, are threatened by a disease, Sudden Acute Narcoleptic Episodes (SANE). This disease strikes 1 person in every 2,000, or 0.05% of the population, and causes the victim to fall into a deep-sleep state for a year.⁸ There are no aftereffects from the disease, but the cost of a worker being removed from the economy for a year is \$32,000. Each country has a population of 100 million workers, so the expected number of cases in each is 50,000 ($0.0005 \times 100,000,000$), and the expected cost of the disease is \$1.6 billion to each ($50,000 \times 32,000$). The total expected cost of the disease worldwide—that is, in both Dorminica and Soporina—is then \$3.2 billion.

Scientists are confident that a crash research program costing \$2 billion will lead to a vaccine that is 100% effective in preventing the disease. Comparing the cost of the research program with the worldwide cost of the disease shows that, from the perspective of the entire population, the research program is clearly worth pursuing (\$2 billion cost < \$3.2 billion savings). However, the government in each country must consider whether to participate in this research program. If both participate, they will share the cost (\$1 billion each), and both will get the benefit (avoiding the \$1.6 billion cost). If only one government chooses to participate, it must fund the entire research program (at a cost of \$2 billion), while the population of the other country can make and use the vaccine for its own population without incurring the research cost.

The payoff matrix for the noncooperative game between Dorminica and Soporina is shown in Figure 10.7. Each country chooses from two strategies, Participate or Not; the payoff

matrix shows the costs to the countries, in billions of dollars, of the various strategy combinations. It is straightforward to verify that each country has a dominant strategy not to participate ($0 > -1$ and $-1.6 > -2$). In addition, the equilibrium payoffs are -1.6 to each, but the two countries could each get -1 if they both participate. So the game is a prisoners' dilemma.

		SOPORIA	
		Participate	Not
DORMINICA	Participate	-1.5, -0.5	-2, 0
	Not	0, -2	-2.4, -0.8
You may need to scroll left and right to see the full figure.			

FIGURE 10.8 Payoffs for Unequal-Population SANE Research Game (in Billions of Dollars)

But now suppose that the two countries have unequal populations of workers, with 150 million in Dorminica and 50 million in Soporina. Then, if no research is funded by either government, the cost to Dorminica of SANE will be \$2.4 billion ($0.0005 \times 150,000,000 \times 32,000$), and the cost to Soporina will be \$0.8 billion ($0.0005 \times 50,000,000 \times 32,000$). If both choose Participate, they will share the cost of research in proportion to their populations, \$1.5 billion for Dorminica and \$0.5 billion for Soporina. The payoff matrix changes to the one illustrated in Figure 10.8.

In this version of the game, Soporina still has a dominant strategy not to participate. But Dorminica's best response to that strategy is now Participate. What has happened to change Dorminica's choice of strategy? Clearly, the answer lies in the unequal distribution of the population in this revised version of the game. Dorminica now stands to suffer such a large portion of the total cost of the disease that it

finds it worthwhile to do the research on its own. This is true even though Dorminica knows full well that Soporina is going to be a free rider and get a share of the full benefit of the research.

The research game in Figure 10.8 is no longer a prisoners' dilemma. Here we see that the dilemma has, in a sense, been solved by the size asymmetry between the countries. The larger country chooses to take on a leadership role and provide the benefit for the whole world.

Examples of leadership in what would otherwise be prisoners' dilemma games are common in international diplomacy. The role of leader often falls naturally to the biggest or best established of the players, a phenomenon labeled "the exploitation of the great by the small." ⁹ For instance, the United States has carried a disproportionate share of the expenditures of our defense alliances, such as NATO, and has maintained a policy of relatively free international trade even while our trading partners, such as Japan and China, have been much more protectionist. In such situations, our theory suggests that the large player would actually act more cooperatively, carrying the burden and tolerating free riding by smaller players. NATO expects each member country to spend 2% of its gross domestic product (GDP) on defense. But our example shows that even when the cost is shared proportionately, small players may still find it optimal not to carry their fair share of the burden and to free ride instead. Whether President Trump's threats in such situations will prove more credible than previous presidents' exhortations to allies to carry their share of the burden is uncertain. Theory suggests not, but only time will tell.

Endnotes

- Think of Rip Van Winkle or of Woody Allen in the movie *Sleeper*, but the duration is much shorter. [Return to reference 8](#)
- Mancur Olson, *The Logic of Collective Action* (Cambridge, Mass.: Harvard University Press, 1965), p. 29. [Return to reference 9](#)

Glossary

leadership

In a prisoners' dilemma with asymmetric players, this is a situation where a large player chooses to cooperate even though he knows that the smaller players will cheat.

6 EXPERIMENTAL EVIDENCE

Numerous people have conducted experiments in which subjects compete in prisoners' dilemma games against each other.^{[10](#)} Such experiments show that cooperation can and does occur in such games, even in repeated versions of known and finite length. Many players start off by cooperating and continue to cooperate for quite a while, as long as the rival player reciprocates. Only in the last few plays of a finite game does defecting seem to creep in. Although this behavior goes against the reasoning of rollback, it can be "profitable" if sustained for a reasonable length of time. The cooperating pairs get higher payoffs than would rational, calculating strategists who defect from the very beginning.

The idea that some level of cooperation may constitute rational—that is, equilibrium—behavior has theoretical backing. Consider the fact that when asked about their reasons for cooperating in the early rounds, players will usually say something such as, "I was willing to try and see if the other player was nice, and when this proved to be the case, I continued to cooperate until the time came to take advantage of the other's niceness." Of course, the other player may not have been genuinely nice, but thinking along similar lines. A rigorous analysis of a finitely repeated prisoners' dilemma shows that this type of asymmetric information can actually be another solution to the dilemma. As long as there is some chance that players are nice rather than selfish, even a selfish player can gain by pretending to be nice. She can reap the higher payoffs from cooperation for a while and then also hope to exploit the gains from double-crossing her opponent near the end of the sequence of plays. For a thorough explication of the case in which just one of the players has the choice between being selfish and being nice, see the online appendix to this chapter.^{[11](#)} Such

cooperative behavior in lab experiments can also be rationalized without relying on this type of asymmetric information. Perhaps the players are not sure that the relationship will actually end at the stated time. Perhaps they believe that their reputations for cooperation will carry over to other similar games against the same opponent or other opponents. Perhaps they think it possible that their opponents are naive cooperators, and they are willing to risk a little loss in testing this hypothesis for a couple of plays. If successful, this experiment will lead to higher payoffs for a sufficiently long time.

In some laboratory experiments, players engage in multiple-round games, each round consisting of a given finite number of repetitions. All the repetitions in any one round are played against the same rival, but each new round is played against a new opponent. Thus, there is an opportunity to develop cooperation with an opponent in each round and to learn from preceding rounds when devising one's strategy against new opponents as the rounds continue. These situations have shown that cooperation lasts longer in early rounds than in later rounds. This result suggests that the theoretical argument on the unraveling of cooperation, based on the use of rollback, is being learned from experience of the play itself over time as players begin to understand the benefits and costs of their actions more fully. Another possibility is that players learn simply that they want to be the first to defect, and so the timing of the initial defection occurs earlier as the number of rounds played increases.

Suppose you were playing a game with a prisoners' dilemma structure and found yourself in a cooperative mode with the known end of the relationship approaching. When should you decide to defect? You do not want to do so too early, while a lot of potential future gains remain. But you also do not want to leave it until too late in the game, because then

your opponent might preempt you and leave you with a low payoff for the periods in which she defects. Similar calculations are relevant when you are in a finitely repeated relationship with an uncertain end date. Your decision about when to defect cannot be deterministic. If it were, your opponent would figure it out and defect in the period before you planned to do so. If no deterministic choice is feasible, then the unraveling of cooperation must include some element of uncertainty, such as mixed strategies, for both players. Many thrillers whose plots hinge on tenuous cooperation among criminals, or between informants and police, acquire their suspense precisely because of this uncertainty.

Examples of the collapse of cooperation as players near the end of a repeated game are observed in numerous situations in the real world, as well as in the laboratory. The story of a long-distance bicycle (or foot) race is one such example. There may be a lot of cooperation for most of the race, as players take turns leading and letting others ride in their slipstreams; nevertheless, as the finish line looms, each participant will want to make a dash for the tape. Similarly, signs saying "No checks accepted" often appear in stores in college towns each spring near the end of the semester.

Computer-simulation experiments have matched a range of very simple to very complex contingent strategies against each other in two-player prisoners' dilemmas. The most famous of them were conducted by Robert Axelrod at the University of Michigan. He invited people to submit computer programs that specified a strategy for playing a prisoners' dilemma repeated a finite but large number (200) of times. There were 14 entrants. Axelrod held a "league tournament" that pitted pairs of these programs against each other, in each case for a run of the 200 repetitions. The point scores for each pairing and its 200 repetitions were recorded, and each program's scores over all its runs against different opponents were added up to see which program did best in the

aggregate against all other programs. Axelrod was initially surprised when “nice” programs did well; none of the top eight programs were ever the first to defect. The winning strategy turned out to be the simplest program: tit-for-tat, submitted by the Canadian game theorist Anatole Rapoport. Programs that were eager to defect in any particular run got the defecting payoff early, but then suffered repetitions of mutual defections and poor payoffs. In contrast, programs that were always nice and cooperative were badly exploited by their opponents. Axelrod explains the success of TFT in terms of four properties: It is at once forgiving, nice, provokable, and clear.

In Axelrod’s words, one does well in a repeated prisoners’ dilemma to abide by these four simple rules: “Don’t be envious. Don’t be the first to defect. Reciprocate both cooperation and defection. Don’t be too clever.” [12](#) The TFT strategy embodies each of these four ideals for a good, repeated prisoners’ dilemma strategy. It is not envious; it does not continually strive to do better than the opponent, only to do well for itself. In addition, TFT clearly fulfills the admonitions not to be the first to defect and to reciprocate, defecting only in retaliation to the opponent’s preceding defection and always reciprocating in kind. Finally, TFT does not suffer from being overly clever; it is simple and understandable to the opponent. In fact, it won the tournament not because it helped players achieve high payoffs in any individual game—the contest was not about “winner takes all”—but because it was always close; it simultaneously encourages cooperation and avoids exploitation, whereas other strategies cannot.

Axelrod then announced the results of his tournament and invited submissions for a second round. Here, people had a clear opportunity to design programs that would beat TFT. The result: TFT won again! The programs that were cleverly designed to beat it could not beat it by very much, and they

did poorly against one another. Axelrod also arranged a tournament of a different kind. Instead of having each program face each other program exactly once, he ran a tournament intended to simulate the evolutionary dynamics underlying the “survival of the fittest.” Each program submitted to the tournament had a certain number of copies included in a larger population, and met an opponent randomly chosen from this population. Those programs that did well were then given a larger proportion of the population in the next round; those that did poorly had their proportion reduced. This was a game of evolution and natural selection, which we will study in greater detail in [Chapter 12](#). But the idea is simple in this context, and the results are fascinating. At first, nasty programs did well at the expense of nice ones. But as the population became nastier and nastier, each nasty program met other nasty programs more and more often, and they began to do poorly and fall in numbers. Then TFT started to do well and eventually triumphed.

However, TFT has some flaws. Most importantly, it assumes no errors in execution of the strategy. If there is some risk that a player intends to play the cooperative action but plays the defecting action in error, then this action can initiate a sequence of retaliatory defecting actions that locks two TFT programs playing one another into a bad outcome; another error is required to rescue them from this sequence. When Axelrod ran a third variant of his tournament, which provided for such random mistakes, TFT could be beaten by even “nicer” programs that tolerated an occasional episode of defecting to see whether it was a mistake or a consistent attempt to exploit them, and retaliated only when convinced that it was not a mistake. [13](#)

Interestingly, a twentieth-anniversary competition modeled after Axelrod’s original contest and run in 2004 and 2005 generated a new winning strategy. [14](#) Actually, the winner was a set of strategies designed to recognize one another during

play so that one would become docile in the face of the other's continued defections. (The authors likened their approach to a situation in which prisoners manage to communicate with each other by tapping on their cell walls.) This collusion meant that some of the strategies submitted by the winning team did very poorly, whereas others did spectacularly well, a testament to the value of working together. Of course, Axelrod's original contest did not permit multiple submissions, so such strategy sets were ineligible, but the winners of the recent competition argue that with no way to preclude coordination, strategies such as those they submitted should have been able to win the original competition as well.

Endnotes

- The literature on experiments involving the prisoners' dilemma game is vast. A brief overview is given by Alvin Roth in *The Handbook of Experimental Economics* (Princeton, N.J.: Princeton University Press, 1995), pp. 26 – 28. Journals in both psychology and economics can be consulted for additional references. For some examples of the outcomes that we describe, see Kenneth Terhune, “Motives, Situation, and Interpersonal Conflict Within Prisoners' Dilemmas,” *Journal of Personality and Social Psychology Monograph Supplement*, vol. 8, no. 30 (1968), pp. 1 – 24; R. Selten and R. Stoecker, “End Behavior in Sequences of Finite Prisoners' Dilemma Supergames,” *Journal of Economic Behavior and Organization*, vol. 7 (1986), pp. 47 – 70; and Lisa V. Bruttel, Werner Güth, and Ulrich Kamecke, “Finitely Repeated Prisoners' Dilemma Experiments Without a Commonly Known End,” *International Journal of Game Theory*, vol. 41 (2012), pp. 23 – 47. Robert Axelrod's *Evolution of Cooperation* (New York: Basic Books, 1984) presents the results of his computer-simulation tournament for the best strategy in an infinitely repeated dilemma. [Return to reference 10](#)
- The corresponding version in which both players can choose between selfish and nice is solved in full in the original article: David Kreps, Paul Milgrom, John Roberts, and Robert Wilson, “Rational Cooperation in a Finitely Repeated Prisoner's Dilemma,” *Journal of Economic Theory*, vol. 27 (1982), pp. 245 – 52. [Return to reference 11](#)
- Axelrod, *Evolution of Cooperation*, p. 110. [Return to reference 12](#)
- For a description and analysis of Axelrod's computer simulations from the biological perspective, see Matt Ridley, *The Origins of Virtue* (New York: Penguin Books,

1997), pp. 61, 75. For a discussion of the difference between computer simulations and experiments using human players, see John K. Kagel and Alvin E. Roth, *Handbook of Experimental Economics* (Princeton, N.J.: Princeton University Press, 1995), p. 29. [Return to reference 13](#)

- See Wendy M. Grossman, “New Tack Wins Prisoner’s Dilemma,” *Wired*, October 13, 2004, available at <http://archived.wired.com/culture/lifestyle/news/2004/10/65317> (accessed August 1, 2014). [Return to reference 14](#)

7 REAL-WORLD DILEMMAS

Games with the prisoners' dilemma structure arise in a surprisingly varied number of contexts in the real world. Although we would be foolish to try to show you every possible instance in which the dilemma can arise, we take the opportunity in this section to consider in detail three specific examples from a variety of fields of study. One example comes from evolutionary biology, a field that we will study in greater detail in [Chapter 12](#). A second example describes the policy of “price matching” as a solution to a prisoners' dilemma pricing game. And a final example concerns international environmental policy and the potential for repeated interactions to mitigate the prisoners' dilemma in this situation.

A. Evolutionary Biology

In our first example, we consider a game known as the bowerbirds' dilemma, from the field of evolutionary biology.¹⁵ Male bowerbirds attract females by building intricate nesting spots called bowers, and female bowerbirds are known to be particularly choosy about the bowers built by their prospective mates. For this reason, male bowerbirds often go out on search-and-destroy missions aimed at ruining other males' bowers. While they are out, however, they run the risk of losing their own bowers to the beak of another male. The ensuing competition between male bowerbirds, in which they can choose whether to maraud others' bowers or guard their own, has the structure of a prisoners' dilemma game.

Ornithologists have constructed a table that shows the payoffs in a two-bird game with two possible strategies, Maraud and Guard (Figure 10.9). The GG payoff represents the benefits associated with Guarding when the rival bird also Guards; GM represents the payoff from Guarding when the rival bird is a Marauder. Similarly, MM represents the benefits associated with Marauding when the rival bird also is a Marauder; MG represents the payoff from Marauding when the rival bird Guards. Careful scientific study of bowerbird matings led to the discovery that $MG > GG > MM > GM$. In other words, the payoffs in the bowerbird game have exactly the same structure as the prisoners' dilemma. The birds' dominant strategy is to Maraud, but when both choose that strategy, they end up in equilibrium with each worse off than if they had both chosen to Guard.

In reality, the strategy used by any particular bowerbird is not actually the result of a process of rational choice on the part of the bird. Rather, in evolutionary games,

strategies are assumed to be genetically “hardwired” into individual organisms, and payoffs represent reproductive success for the different types. Then equilibria in such games define the proportions of types that naturalists can expect to observe in the population—all Marauders, for instance, if Maraud is a dominant strategy, as in Figure 10.9. This equilibrium outcome is not the best one, however, given the existence of the dilemma. In constructing a solution to the bowerbirds’ dilemma, we can appeal to the repetitive nature of the interaction in the game. In the case of the bowerbirds, repeated play against the same or different opponents in the course of several breeding seasons can allow a bird to choose a flexible strategy based on his opponent’s last move. Contingent strategies such as tit-for-tat can be, and often are, adopted in evolutionary games to solve exactly this type of dilemma. We will return to the idea of evolutionary games and provide detailed discussions of their structure and equilibrium outcomes in [Chapter 12](#).

		BIRD 2	
		Maraud	Guard
BIRD 1	Maraud	MM, MM	MG, GM
	Guard	GM, MG	GG, GG

FIGURE 10.9 Bowerbirds’ Dilemma

B. Price Matching

Now we return to a pricing game, in which we consider two specific stores engaged in price competition with each other, using identical price-matching policies. The stores in question, Staples and Office Depot, are both national chains that regularly advertise their prices for name-brand office supplies. In addition, each store maintains a published policy that guarantees customers that it will match the advertised price of any competitor on a specific item (model and item numbers must be identical) as long as the customer provides the competitor's printed advertisement.¹⁶

For the purposes of this example, we assume that the firms have only two possible prices that they can charge for a particular package of printer ink (Low or High). In addition, we use hypothetical profit numbers, and we further simplify the analysis by assuming that Staples and Office Depot are the only two competitors in the office supplies market in a particular city—Billings, Montana, for example.

OFFICE DEPOT			
		Low	High
STAPLES	Low	2, 500, 2, 500	5, 000, 0
	High	0, 5, 000	3, 400, 3, 400

FIGURE 10.10 Staples and Office Depot Ink Pricing Game

OFFICE DEPOT				
		Low	High	Match
STAPLES	Low	2, 500, 2, 500	5000, 0	2, 500, 2, 500
You may figure.	need to scroll left and right to see the full			

	OFFICE DEPOT			
		Low	High	Match
	High	0, 5,000	3,400, 3,400	3,400, 3,400
	Match	2,500, 2,500	3,400, 3,400	3,400, 3,400
You may need to scroll left and right to see the full figure.				

FIGURE 10.11 Ink Pricing Game with Price Matching

Suppose, then, that the basic structure of the game between the two firms can be illustrated as in Figure 10.10. If both firms advertise low prices, they split the available customer demand, and each earns \$2,500. If both advertise high prices, they split a market with lower sales, but their markups end up being large enough to let them each earn \$3,400. Finally, if they advertise different prices, then the one advertising a high price gets no customers and earns nothing, whereas the one advertising a low price earns \$5,000.

The game illustrated in Figure 10.10 is clearly a prisoners' dilemma. Advertising and selling at a low price is the dominant strategy for each firm, although both would be better off if each advertised and sold at the high price. But, as mentioned earlier, each firm actually makes use of a third pricing strategy: a price-matching guarantee to its customers. How does the inclusion of such a policy alter the prisoners' dilemma that would otherwise exist between these two firms?

Consider the effects of allowing firms to choose among pricing low, pricing high, and price matching. The Match strategy entails advertising a high price, but promising to match any lower price advertised by a competitor. A firm using Match then benefits from advertising high if the rival firm does so also, but it does not suffer any harm from advertising a high price if the rival advertises a low price.

We can see this in the payoff structure for the new game, shown in Figure 10.11. In that table, we see that a combination of one firm playing Low while the other plays Match is equivalent to both playing Low, while a combination of one firm playing High while the other plays Match (or both playing Match) is equivalent to both playing High.

Using our standard tools for analyzing simultaneous-move games shows that High is weakly dominated by Match for both players, and that once High is eliminated, Low is weakly dominated by Match also. The resulting Nash equilibrium entails both firms using the Match strategy. In equilibrium, both firms earn \$3,400—the profit level associated with both firms pricing high in the original game. The addition of the Match strategy has allowed the firms to emerge from the prisoners’ dilemma that they faced when they had only the choice between two simple pricing strategies, Low or High.

How did this happen? The Match strategy acts as a penalty mechanism. By guaranteeing to match Office Depot’s low price, Staples substantially reduces the benefit that Office Depot achieves by advertising a low price while Staples is advertising a high price. Promising to match Office Depot’s low price hurts Staples, too, because Staples has to accept the lower profit associated with the low price. Thus, the price-matching guarantee is a method of penalizing both players whenever either one defects. This penalty scheme is just like that in the crime mafia example discussed in [Section 4](#), except that this scheme—and the higher equilibrium prices that it supports—is observed in markets in virtually every city.

Actual empirical evidence of the effects of price-matching policies on prices is available but limited, and some research has found evidence of lower prices in markets with such policies.¹⁷ However, more recent experimental evidence does support the collusive effect of price-matching policies.

This result should put all customers on alert.¹⁸ Even though stores that match prices promote their policies in the name of competition, the ultimate outcome when all firms use such policies can be better for the firms than if there were no price matching at all, and so customers can be the ones who are hurt.

C. International Environmental Policy

Our final example of a real-world prisoners’ dilemma pertains to international climate control agreements. In 2003, the American Geophysical Union stated that “humanity is the major influence on the global climate change observed over the past 50 years,” but that “rapid societal responses can significantly lessen negative outcomes.” Since then, scientists have only grown more certain that emissions of carbon dioxide and other greenhouse gases (GHGs) are causing atmospheric temperatures to rise to levels that threaten the well-being, and perhaps the very existence, of the human species on earth.¹⁹ Except for a nuclear Armageddon, it is hard to think of a worse “global bad.”

		THEM	
		Cut emissions	Don’ t cut
US	Cut emissions	−1, −1	−20, 0
	Don’ t cut	0, −20	−12, −12

FIGURE 10.12 Greenhouse Gas Emissions Game

The difficulty in achieving a global reduction in GHG emissions comes in part from the nature of the interaction among the nations of the world: It is a prisoners’ dilemma. No individual country has any incentive to reduce its own emissions, knowing that if it does so alone, it will bear significant costs with little benefit to overall climate change reduction. And if other countries do reduce their emissions, the first country cannot be stopped from enjoying the benefits of the others’ actions.

Consider the emissions reduction problem as a game played between two countries, Us and Them. Estimates generated by the British government's Office on Climate Change suggest that coordinated action may come at a cost of about 1% of GDP per nation, whereas coordinated inaction could cost each nation between 5% and 20% of its GDP, perhaps 12% on average.²⁰ By extension, the cost to one country of cutting emissions on its own may be at the high end of the inaction estimate (20%), but holding back and letting the other country cut emissions could entail virtually no cost at all. We can then summarize the situation between Us and Them using the game table in Figure 10.12, where payoffs represent changes in GDP for each country.

The game in Figure 10.12 is indeed a prisoners' dilemma. Both countries have a dominant strategy to refuse to cut their emissions. The single Nash equilibrium occurs when neither country cuts emissions, but they suffer as a group as a result of the ensuing climate change.

Some attempts have been made to resolve the dilemma. The two that have progressed farthest toward agreement and implementation are the Kyoto Protocol and the Paris Agreement. However, both have major shortcomings that are already apparent. We outline these agreements briefly here and also discuss some ideas and solutions proposed by economists.

The Kyoto Protocol was negotiated by the United Nations Framework Convention on Climate Change (UNFCCC) in 1997, and went into effect in 2005 to last until 2012. More than 170 countries signed on and more joined later, although the United States was noticeably absent from the list. The protocol was extended in December 2012 to last through 2020. Of 192 member countries, only 124 accepted the extension. Canada withdrew from the protocol in 2012, and others, including Belarus, Ukraine, and Kazakhstan, have stated their

intention to withdraw. Thus the future of this agreement is unclear at best.

The Kyoto Protocol had some commendable features. First, it was based on the principle of “common but differentiated responsibilities.” This meant that developed countries, which had historically contributed more to GHG accumulation, and at the same time had greater economic capability to combat climate change, should bear a greater share of the burden. Second, the protocol imposed binding targets for emissions reduction. Unfortunately, some countries did not take on the new targets in the post-2012 phase, and the protocol lacked any effective provisions to enforce participation.

Discussions in the UNFCCC on measures to be taken after 2020 led to the separate Paris Agreement of 2015, signed by 195 countries and ratified by 185. Its goal is to keep the global average temperature rise in the twenty-first century below 2° C, and preferably below 1.5° C. For this purpose, each member country is to regularly report on its GHG emission reduction target and on the actions it will take to achieve the target. Unfortunately, these targets and actions are purely voluntary, and despite many ambitious statements, especially from European countries, actual performance has been disappointing. Recent studies, summarized by the American Geophysical Union,^{[21](#)} find that the commitments of the United States, European Union, and China would not achieve the 1.5° C goal, even if the rest of the world reduced its emissions to zero. President Trump has announced his intention to withdraw the United States from the Paris Agreement, and this action may start a cascade of other withdrawals or relaxation of targets. We saw in [Section 5](#) that large players can serve an important leadership role to resolve prisoners’ dilemmas and should be willing to bear a disproportionately large share of the burden for this purpose. In this instance, the opposite seems to be

happening: The largest players are lagging with their targets and actions. That does not bode well for the outcome.

Of course, this prisoners' dilemma is not a one-period game, but a repeated one. That observation implies that there are ways to sustain a good outcome, as Michael Liebrieck has suggested.²² He argues that the repeated nature of this game makes it amenable to solution by way of contingent strategies, and that countries should use strategies that embody the four critical properties of TFT, as outlined by Axelrod and described in [Section 6](#). Specifically, countries are encouraged to employ strategies that are “nice” (signing on to the protocol and beginning emissions reductions), “retaliatory” (employing mechanisms to punish those that do not do their part), “forgiving” (welcoming those newly accepting the protocol), and “clear” (specifying actions and reactions).

Liebrieck assesses the actions of current players, including the European Union, the United States, and developing countries (as a group), and provides some suggestions for improvements. He explains that the European Union does well with nice, forgiving, and clear, but not with retaliatory, so other countries will do best to defect when interacting with the European Union. One solution would be for the European Union to institute carbon-related import taxes or another retaliatory policy for dealing with recalcitrant trade partners. The United States, in contrast, ranks high on retaliatory and forgiving, given its history of such behavior following the end of the Cold War. But it has not been nice or clear, at least on the national level (individual states may behave differently), giving other countries an incentive to retaliate against it quickly and painfully, if possible. The solution is for the United States to make a meaningful commitment to GHG emission reductions, a standard conclusion in most policy circles. Developing countries are described as not nice (negotiating no carbon limits for themselves),

retaliatory, unclear, and quite unforgiving. A more beneficial strategy, argues Liebreich, would be for these countries—particularly China, India, and Brazil—to make clear their commitment to sharing in international efforts to reduce climate change; this approach would leave them less subject to retaliation and more likely to benefit from a global improvement in climatic outlook.

The general conclusion is that the process of international GHG emissions reduction fits the profile of a prisoners' dilemma game. But the future of the global climate should not be considered a lost cause simply because of this aspect of the nations' one-time interaction. Repeated play among the nations involved in the Kyoto Protocol negotiations makes the game amenable to solutions by way of contingent (nice, clear, and forgiving, but also retaliatory) strategies.

Endnotes

- Larry Conik, “Science Classics: The Bowerbird’ s Dilemma,” *Discover*, October 1994. [Return to reference 15](#)
- The Staples policy applies to in-store purchases only, while Office Depot includes online purchases in its guarantee. Both stores give customers 14 days after purchase to find a lower price, and both include details of their policies on their Web sites. Similar policies exist in many industries, including that for credit cards, where “interest rate matching” has been observed. See Aaron S. Edlin, “Do Guaranteed-Low-Price Policies Guarantee High Prices, and Can Antitrust Rise to the Challenge?” *Harvard Law Review*, vol. 111, no. 2 (December 1997), pp. 529 – 75. [Return to reference 16](#)
- J. D. Hess and Eitan Gerstner present evidence of increased prices as a result of price-matching policies in “Price-Matching Policies: An Empirical Case,” *Managerial and Decision Economics*, vol. 12 (1991), pp. 305 – 15. Contrary evidence is provided by Arbatskaya, Hviid, and Shaffer, who find that the effect of matching policies is to lower prices; see Maria Arbatskaya, Morten Hviid, and Greg Shaffer, “Promises to Match or Beat the Competition: Evidence from Retail Tire Prices,” *Advances in Applied Microeconomics*, vol. 8, *Oligopoly* (New York: JAI Press, 1999), pp. 123 – 38. [Return to reference 17](#)
- See Subhasish Dugar, “Price-Matching Guarantees and Equilibrium Selection in a Homogeneous Product Market: An Experimental Study,” *Review of Industrial Organization*, vol. 30 (2007), pp. 107 – 19. [Return to reference 18](#)
- As NASA atmospheric scientist Kate Marvel explained during an expert-panel discussion in September 2018, “We are more sure that greenhouse gas is causing climate change than we are that smoking causes cancer.” See Abel Gustafson and Matthew Goldberg, “Even Americans Highly

Concerned about Climate Change Dramatically Underestimate the Scientific Consensus,” Yale Program on Climate Change Communication, October 18, 2018. [Return to reference 19](#)

- See Nicholas Stern, *The Economics of Climate Change: The Stern Review* (Cambridge: Cambridge University Press, 2007). [Return to reference 20](#)
- See <https://eos.org/scientific-press/new-studies-highlight-challenge-of-meeting-paris-agreement-climate-goals>. [Return to reference 21](#)
- Michael Liebreich presents his analysis of the Kyoto Protocol as a repeated prisoners’ dilemma in his paper “How to Save the Planet: Be Nice, Retaliatory, Forgiving and Clear,” New Energy Finance White Paper, September 11, 2007, available for download from www.bnef.com/InsightDownload/7080/pdf/ (accessed August 1, 2014). [Return to reference 22](#)

SUMMARY

The prisoners' dilemma is probably the most famous game of strategy. Each player has a dominant strategy (Defect), but the equilibrium outcome is worse for all players than when each uses her dominated strategy (Cooperate). The dilemma can be solved in some cases when the way the game is played changes, or when mechanisms exist that can change the payoff structure.

The best-known solution to the dilemma is *repeated play*. In a game with a finite number of periods, the value of future cooperation is eventually zero, and rollback yields an equilibrium with no cooperative behavior. With infinite play (or an uncertain end date), cooperation can be achieved with the use of an appropriate *trigger strategy* such as *tit-for-tat (TFT)* or the *grim strategy*; in either case, cooperation is possible only if the *present value* of cooperation exceeds the present value of defecting. More generally, the prospect of no future relationship or of a short-term relationship leads to decreased cooperation among players.

The dilemma can also be solved with a change in the order of moves in which the second mover is trustworthy and makes a promise, or with *punishment* mechanisms that alter the payoffs for players who defect from cooperation when their rivals are cooperating or when others are also defecting. Another solution mechanism, *leadership*, exists when a large or strong player's loss from defecting is greater than the available gain from cooperative behavior on that player's part.

Experimental evidence suggests that players often cooperate longer than theory might predict. Such behavior can be explained by incomplete knowledge of the game on the part of the players or by their views regarding the benefits of

cooperation. Tit-for-tat has been observed to be a simple, nice, provokable, and forgiving strategy that performs very well on the average in repeated prisoners' dilemmas.

Prisoners' dilemmas arise in a variety of contexts. Specific examples from evolutionary biology, product pricing, and international environmental policy show how to explain and predict actual behavior by using the framework of the prisoners' dilemma.

KEY TERMS

[compound interest](#) ([382](#))

[contingent strategy](#) ([379](#))

[discount factor](#) ([383](#))

[effective rate of return](#) ([384](#))

[grim strategy](#) ([379](#))

[infinite horizon](#) ([381](#))

[leadership](#) ([393](#))

[penalty](#) ([390](#))

[present value \(PV\)](#) ([381](#))

[punishment](#) ([379](#))

[repeated play](#) ([377](#))

[tit-for-tat \(TFT\)](#) ([379](#))

[trigger strategy](#) ([379](#))

Glossary

[repeated play](#)

A situation where a one-time game is played repeatedly in successive periods. Thus, the complete game is mixed, with a sequence of simultaneous-move games.

[contingent strategy](#)

In repeated play, a plan of action that depends on other players' actions in previous plays. (This is implicit in the definition of a strategy; the adjective "contingent" merely reminds and emphasizes.)

[trigger strategy](#)

In a repeated game, this strategy cooperates until and unless a rival chooses to defect, and then switches to noncooperation for a specified period.

[punishment](#)

We reserve this term for costs that can be inflicted on a player in the context of a repeated relationship (often involving termination of the relationship) to induce him to take actions that are in the joint interests of all players.

[grim strategy](#)

A strategy of noncooperation forever in the future, if the opponent is found to have cheated even once. Used as a threat of punishment in an attempt to sustain cooperation.

[tit-for-tat \(TFT\)](#)

In a repeated prisoners' dilemma, this is the strategy of [1] cooperating on the first play and [2] thereafter doing each period what the other player did the previous period.

[present value \(PV\)](#)

The total payoff over time, calculated by summing the payoffs at different periods each multiplied by the

appropriate discount factor to make them all comparable with the initial period' s payoffs.

infinite horizon

A repeated decision or game situation that has no definite end at a fixed finite time.

compound interest

When an investment goes on for more than one period, compound interest entails calculating interest in any one period on the whole accumulation up to that point, including not only the principal initially invested but also the interest earned in all previous periods, which itself involves compounding over the period previous to that.

discount factor

In a repeated game, the fraction by which the next period' s payoffs are multiplied to make them comparable with this period' s payoffs.

effective rate of return

Rate of return corrected for the probability of noncontinuation of an investment to the next period.

penalty

We reserve this term for one-time costs (such as fines) introduced into a game to induce the players to take actions that are in their joint interests.

leadership

In a prisoners' dilemma with asymmetric players, this is a situation where a large player chooses to cooperate even though he knows that the smaller players will cheat.

SOLVED EXERCISES

1. “If a prisoners’ dilemma is repeated 100 times, and both players know how many repetitions to expect, they are sure to achieve their cooperative outcome.” True or false? Explain and give an example of a game that illustrates your answer.
2. Consider a two-player game between Child’s Play and Kid’s Korner, each of which produces and sells wooden swing sets for children. Each firm can set either a high or a low price for a standard two-swing, one-slide set. If they both set a high price, each receives profits of \$64,000 per year. If one sets a low price and the other sets a high price, the low-price firm earns profits of \$72,000 per year, while the high-price firm earns \$20,000. If they both set a low price, each receives profits of \$57,000.
 1. Verify that this game has a prisoners’ dilemma structure by looking at the ranking of the payoffs associated with the different strategy combinations (both cooperate, both defect, one defects, and so on). What are the Nash equilibrium strategies and payoffs in the simultaneous-move game if the players make price decisions only once?
 2. If the two firms decide to play this game for a fixed number of one-year periods—say, for four years—what will each firm’s total profits be at the end of the game? (Don’t discount.) Explain how you arrived at your answer.
 3. Suppose that the two firms play this game repeatedly forever. Let each of them use a grim strategy in which they both price high unless one of them defects, in which case they price low for the rest of the game. What is the one-time gain from defecting against an opponent playing such a strategy? How much

does each firm lose, in each future period, after it defects once? If $r = 0.25$ ($\delta = 0.8$), will it be worthwhile for them to cooperate? Find the range of values of r (or δ) for which this strategy is able to sustain cooperation between the two firms.

4. Suppose the firms play this game repeatedly year after year, with neither expecting any change in their interaction. If the world were to end after four years, without either firm having anticipated this event, what would each firm's total profits (not discounted) be at the end of the game? Compare your answer here with your answer in part (b). Explain why the two answers are different, if they are different, or why they are the same, if they are the same.
5. Suppose now that the firms know that there is a 10% probability that one of them will go bankrupt in any given year. If bankruptcy occurs, the repeated game between the two firms ends. Will this knowledge change the firms' actions when $r = 0.25$? What if the probability of a bankruptcy increases to 35% in any year?
3. A firm has two divisions, each of which has its own manager. These managers are paid according to their effort in promoting productivity in their divisions. The payment scheme is based on a comparison of their two outcomes. If both managers have expended "high effort," each earns \$150,000 a year. If both have expended "low effort," each earns "only" \$100,000 a year. But if one of the two managers shows high effort whereas the other shows low effort, the high-effort manager is paid \$150,000 plus a \$50,000 bonus, but the second (low-effort) manager gets a reduced salary (for subpar performance in comparison with her competition) of \$80,000. The managers make their effort decisions independently and without knowledge of the other manager's choice.

1. Assume that expending effort is costless to the managers and draw the payoff table for this game. Find the Nash equilibrium of the game and explain whether the game is a prisoners' dilemma.
2. Now suppose that expending high effort is costly to the managers (i.e., that it is a costly signal of quality). In particular, suppose that high effort costs an equivalent of \$60,000 a year to a manager who chooses this effort level. Draw the game table for this new version of the game and find the Nash equilibrium. Explain whether the game is a prisoners' dilemma and how it has changed from the game in part (a).
3. If the cost of high effort is equivalent to \$80,000 a year, how does the game change from that described in part (b)? What is the new equilibrium? Explain whether the game is a prisoners' dilemma and how it has changed from the games in parts (a) and (b).
4. You have to decide whether to invest \$100 in a friend's enterprise, where in a year's time the money will increase to \$130. You have agreed that your friend will then repay you \$120, keeping \$10 for himself. But instead he may choose to run away with the whole \$130. Any of your money that you don't invest in your friend's venture you can invest elsewhere safely at the prevailing rate of interest r and get $\$100(1+r)$ next year.
 1. Draw the game tree for this situation and show the rollback equilibrium. Next, suppose this game is played repeatedly for an infinite number of years. That is, each year you have the opportunity to invest another \$100 in your friend's enterprise, and you agree to split the resulting \$130 in the manner already described. From the second year onward, you get to make your decision of whether to invest with your friend in the light of whether he made the agreed-upon repayment the preceding year. The rate of interest between any two successive periods is r , the

same as the outside rate of interest, and the same for you and your friend.

2. For what values of r can there be an equilibrium outcome of the infinitely repeated game in which each period you invest with your friend and he repays you as agreed?
3. If the rate of interest is 10% per year, can there be an alternative profit-splitting agreement that is an equilibrium outcome of the infinitely repeated game, where each period you invest with your friend and he repays as agreed?
5. Recall the example from Exercise S3 in which two division managers' choices of high or low effort levels determine their salary payments. In part (b) of that exercise, the cost of exerting high effort is assumed to be \$60,000 a year. Suppose now that the two managers play the game in part (b) of Exercise S3 repeatedly for many years. Such repetition allows scope for an unusual type of cooperation in which one is designated to choose high effort while the other chooses low effort. This cooperative agreement requires that the high-effort manager make a side payment to the low-effort manager so that their payoffs are identical.
 1. What size of side payment guarantees that the final payoffs for the two managers are identical? How much does each manager earn in a year in which the cooperative agreement is in place?
 2. Cooperation in this repeated game entails each manager's choosing her assigned effort level and the high-effort manager making the designated side payment. Defection entails refusing to make the side payment. Under what values of the rate of return can this agreement sustain cooperation in the managers' repeated game?
6. Consider the game of chicken in [Chapter 4](#), with slightly more general payoffs (Figure 4.16 shows an example in which $k = 1$):

		DEAN	
		Swerve	Straight
JAMES	Swerve	0, 0	-1, k
	Straight	k , -1	-2, -2

Suppose this game is played every Saturday evening. If $k < 1$, the two players stand to benefit by cooperating to play (Swerve, Swerve) all the time, whereas if $k > 1$, they stand to benefit by cooperating so that one plays Swerve and the other plays Straight, taking turns to go Straight in alternate weeks. Can either type of cooperation be sustained?

7. Recall the example from Exercise S8 in [Chapter 5](#), where South Korea and Japan compete in the market for production of VLCCs. As in parts (a) and (b) of that exercise, the cost of building the ships is \$30 (million) in each country, and the demand for ships is $P = 180 - Q$, where $Q = q_{\text{Korea}} + q_{\text{Japan}}$.
 1. Previously, we found the Nash equilibrium for the game. Now find the collusive outcome. What total quantity should be set by the two countries in order to maximize their joint profit?
 2. Suppose the two countries produce equal quantities of VLCCs, so that they earn equal shares of this collusive profit. How much profit would each country earn? Compare this profit with the amount they would earn in the Nash equilibrium.
 3. Now suppose the two countries are in a repeated relationship. Once per year, they choose production quantities, and each can observe the amount its rival produced in the previous year. They wish to cooperate to sustain the collusive profit levels you found in part (b). In any one year, one of them can defect from the agreement. If one of them holds the quantity

at the agreed-upon level, what is the best defecting quantity for the other? What are the resulting profits?

4. Write down a matrix that represents this game as a prisoners' dilemma.
5. At what interest rates will collusion be sustainable when the two countries use grim (defect forever) strategies?
8. In the (imaginary) game show *Roll It!!*, two players decide at the end of the game whether to Roll or Steal, given a prize pot worth $\$X$ before them. If both players choose Steal, the game ends, and each gets $\$X/2$. If only one chooses Steal, the game ends, and the person who chose Steal gets $\$X$ while the person who chose Roll gets nothing. Finally, if both choose Roll, an enormous six-sided die is rolled; if it comes up 1, the game ends, and both players lose everything; if it comes up 2 through 6, that many thousands of dollars are added to the prize pot, and the game continues with another choice to Roll or Steal, now with the bigger prize pot.
 1. First suppose that Player 2 always chooses Roll, allowing Player 1 to steal the whole prize pot whenever she wants. Verify that, in order to maximize her expected winnings, player 1 should choose Roll whenever the prize pot is less than \$20,000 and choose Steal whenever it exceeds \$20,000. Hint: Compare Player 1's expected payoff when stealing a pot worth $\$X(+\$X)$ with her expected payoff when rolling *just one more time* and then stealing the pot, if possible.
 2. Suppose that the prize pot has grown to \$20,000 or more. Verify that each player has a dominant strategy to choose Steal, and hence that the game will end with each player getting half of the prize pot.
 3. Suppose that the prize pot has grown to \$18,000 and, as in part (b), that both players are certain to choose Steal once the prize pot reaches \$20,000 or

more. Verify that the resulting game is a prisoners' dilemma.

4. Suppose that the game begins with a prize pot of \$1,000. Does a rollback equilibrium exist in which both players decide to Roll at least one time?
 5. (Optional) What is the highest prize pot the viewers of *Roll It!!* would ever see on the show, assuming that rollback equilibrium choices are always made by the players?
9. People who are kicked off eBay for committing fraud can get back on again by buying a new eBay username on the Aspin Suspensions Forum (www.aspin.com). For example, in April 2019, an “Aged, Verified, High-Limit USA eBay/PayPal Seller Account [with] 500/\$25,000 Limits” could be had for about \$200.²³ But what stops sellers on the Aspin Forum from cheating the cheaters who come to them looking to get back onto eBay? For this question, assume that it costs \$100 to create a new “ready to sell” eBay username, that buyers are willing to pay \$300 for such usernames, that the Aspin Forum is the only place in the world where such usernames can be safely traded, and that the going price is \$200 per username. Assume that sellers discount future payoffs using an annual discount factor of $\delta = 2/3$
1. Suppose that a seller in good standing on the Aspin Forum will be able to sell one eBay username per year forever. What is the present value of the seller's total profits, if the seller never cheats and hence always remains in good standing?
 2. Suppose that the game of purchasing a new username works as follows: First, a buyer decides whether to pay \$200; then, the seller decides whether to deliver (providing the login details for a ready-to-sell eBay username) or cheat the buyer (delivering nothing). Assume that any cheating seller will be detected and permanently barred from the Aspin Forum. Draw the

game tree for this sequential-move game. What are the rollback equilibrium strategies and payoffs?

3. Imagine that eBay takes steps to make it less lucrative to commit fraud on its site. As a result, the going price for a fraudulent eBay username falls to \$120. Verify that a seller with a discount factor of $\delta = 2/3$ now has an incentive to cheat.
4. In the context of part (c), what can the Aspin Forum do so that sellers on the site no longer have an incentive to cheat buyers?

UNSOLVED EXERCISES

1. Two people, Baker and Cutler, play a game in which they choose and divide a prize. Baker decides how large the total prize should be; she can choose either \$10 or \$100. Cutler chooses how to divide the prize chosen by Baker; Cutler can choose either an equal division or a split where she gets 90% and Baker gets 10%. Draw the payoff table of the game and find its equilibria for each of the following situations:
 1. When the moves are simultaneous.
 2. When Baker moves first.
 3. When Cutler moves first.
 4. Is this game a prisoners' dilemma? Why or why not?
2. Sociologist Diego Gambetta begins his book *The Sicilian Mafia*²⁴ by quoting a cattle breeder he interviewed: "When the butcher comes to me to buy an animal, he knows that I want to cheat him [on quality]. But I know that he wants to cheat me [on payment]. So we need Peppe [the mafioso] to make us agree. And we both pay Peppe a percentage of the deal." Relate this situation to the various methods of resolving prisoners' dilemmas in this chapter. Consider the following questions in your answer: (i) Who has a repeated relationship with whom? (ii) Why are both parties in the trade willing to pay Peppe? (iii) What keeps Peppe honest—that is, what stops him from double-crossing a trader who has behaved honestly and extorting him for more payment? Related to this, what determines Peppe's percentage? (iv) Peppe can enforce honesty between the traders either by trashing a cheater's reputation, or by direct physical punishment. In which mode will Peppe's fee be higher?
3. Consider a small town that has a population of dedicated pizza eaters but is able to accommodate only two pizza shops, Donna's Deep Dish and Pierce's Pizza Pies. Each

seller has to choose a price for its pizza, but for simplicity, assume that only two prices are available: high and low. If a high price is set, the sellers can achieve a profit margin of \$12 per pie; the low price yields a profit margin of \$10 per pie. Each store has a loyal captive customer base that will buy 3,000 pies per week, no matter what price is charged by either store. There is also a floating demand of 4,000 pies per week. The people who buy these pies are price conscious and will go to the store with the lower price; if both stores charge the same price, this demand will be split equally between them.

1. Draw the game table for the pizza-pricing game, using each store's profits per week (in thousands of dollars) as payoffs. Find the Nash equilibrium of this game and explain why it is a prisoners' dilemma.
2. Now suppose that Donna's Deep Dish has a much larger loyal clientele that guarantees it the sale of 11,000 (rather than 3,000) pies a week. Profit margins and the size of the floating demand remain the same. Draw the payoff table for this new version of the game and find the Nash equilibrium.
3. How does the existence of the larger loyal clientele for Donna's Deep Dish help solve the pizza stores' dilemma?
4. Six friends stop at a burger joint for lunch. There are two items on the menu: (i) a regular burger that costs \$4 and (ii) a deluxe burger that costs \$8. Each friend feels that eating a regular burger is worth \$5 while eating a deluxe burger is worth \$6.
 1. The friends have agreed to split equally the overall cost of the meal; if the total cost is T , then each friend will pay $T/6$. Verify that the resulting game is a prisoners' dilemma in which each friend has a dominant strategy to order a deluxe burger. Hint: If D is the number of friends who order a deluxe burger,

then the total cost is $T(D) = 8D + 4(6 - D) = 24 + 4D$.

2. One day, the burger joint decides to raise the price of a deluxe burger to \$12. What is the unique Nash equilibrium in the resulting game? Are the friends better off when deluxe burgers cost \$12 or when they cost \$8?
5. A town council consists of three members who vote every year on their own salary increases. Two Yes votes are needed to pass the increase. Each member would like a higher salary but would like to vote against it herself because that looks good to the voters. Specifically, the payoffs of each outcome are as follows:

Raise passes, own vote is No: 10

Raise fails, own vote is No: 5

Raise passes, own vote is Yes: 4

Raise fails, own vote is Yes: 0

Voting is simultaneous. Draw the (three-dimensional) payoff table, and show that in Nash equilibrium the raise fails unanimously. Examine how a repeated relationship among the members can secure them salary increases every year if (1) every member serves a three-year term, (2) every year in rotation, one of them is up for reelection, and (3) the townspeople have short memories, remembering only the members' votes on the salary increase motion of the current year and not those of past years.

6. Consider the following game, which comes from James Andreoni and Hal Varian at the University of Michigan.^{[25](#)} A neutral referee runs the game. There are two players, Row and Column. The referee gives two cards to each: 2 and 7 to Row and 4 and 8 to Column. Who gets what cards is common knowledge. Then, playing simultaneously and

independently, each player is asked to hand over to the referee either his High card or his Low card. The referee hands out payoffs—which come from a central kitty, not from the players’ pockets—that are measured in dollars and depend on the cards that he collects. If Row chooses his Low card, 2, then Row gets \$2; if he chooses his High card, 7, then Column gets \$7. If Column chooses his Low card, 4, then Column gets \$4; if he chooses his High card, 8, then Row gets \$8.

1. Show that the complete payoff table is as follows:

		COLUMN	
		Low	High
ROW	Low	2, 4	10, 0
	High	0, 11	8, 7

2. What is the Nash equilibrium? Verify that this game is a prisoners’ dilemma.

Now suppose the game has the following stages. The referee hands out cards as before; who gets what cards is common knowledge. At stage 1, each player, out of his own pocket, can hand over a sum of money, which the referee is to hold in an escrow account. This amount can be zero but cannot be negative. When both players have made their stage 1 choices of sums to hand over, these choices are publicly disclosed. Then, at stage 2, the two players make their choices of cards, again simultaneously and independently. The referee hands out payoffs from the central kitty in the same way as in the single-stage game before, but in addition, he disposes of the escrow account as follows: If Column chooses his High card, the referee hands over to Column the sum that Row put into the escrow account; if Column chooses his Low card, Row’s sum reverts back to him. The disposition of the

sum that Column deposited depends similarly on Row's card choice. All these rules are common knowledge.

1. (c) Find the rollback (subgame-perfect) equilibrium of this two-stage game. Does it resolve the prisoners' dilemma? What is the role of the escrow account?
7. Glassworks and Clearsmooth compete in the local market for windshield repairs. The market size (total available profits) is \$10 million per year. Each firm can choose whether to advertise on local television. If a firm chooses to advertise in a given year, it costs that firm \$3 million. If one firm advertises and the other doesn't, then the former captures the whole market. If both firms advertise, they split the market 50:50. If both firms choose not to advertise, they also split the market 50:50.
 1. Suppose the two windshield-repair firms know they will compete for just one year. Draw the payoff matrix for this game. Find the Nash equilibrium strategies.
 2. Suppose the firms play this game for five years in a row, and they know that at the end of five years, both firms plan to go out of business. What is the subgame-perfect equilibrium for this five-period game? Explain.
 3. What would be a tit-for-tat strategy in the game described in part (b)?
 4. Suppose the firms play this game repeatedly forever, and suppose that future profits are discounted with an interest rate of 20% per year. Can you find a subgame-perfect equilibrium that involves higher annual payoffs than the equilibrium in part (b)? If so, explain what strategies are involved. If not, explain why not.
8. Consider the pizza stores introduced in Exercise U3, Donna's Deep Dish and Pierce's Pizza Pies. Suppose that they are not constrained to choose from only two possible

prices, but that each can choose a specific price to maximize profits. Suppose further that it costs \$3 to make each pizza (for both stores), and that experience or market surveys have shown that the relation between sales (Q) and price (P) for each firm is as follows:

$$Q_{\text{Pierce}} = 12 - P_{\text{Pierce}} + 0.5 P_{\text{Donna}}.$$

Then profits per week (Y , in thousands of dollars) for each firm are

$$Y_{\text{Pierce}} = (P_{\text{Pierce}} - 3) Q_{\text{Pierce}} = (P_{\text{Pierce}} - 3) (12 - P_{\text{Pierce}} + 0.5 P_{\text{Donna}}),$$

$$Y_{\text{Donna}} = (P_{\text{Donna}} - 3) Q_{\text{Donna}} = (P_{\text{Donna}} - 3) (12 - P_{\text{Donna}} + 0.5 P_{\text{Pierce}}).$$

1. Use these profit functions to determine each firm's best-response rule, as in [Chapter 5](#), and use these best-response rules to find the Nash equilibrium of this pricing game. What prices do the firms choose in equilibrium? How much profit per week does each firm earn?
2. If the firms work together and choose a joint best price, P , then the profit of each will be

$$Y_{\text{Donna}} = Y_{\text{Pierce}} = (P - 3) (12 - P + 0.5P) = (P - 3) (12 - 0.5P).$$

What price do they choose to maximize joint profits?

3. Suppose the two stores are in a repeated relationship, trying to sustain the joint profit-maximizing prices calculated in part (b). They print new menus each month and thereby commit themselves to prices for the whole month. In any one month, one of them can defect from the agreement. If one of them holds the price at the agreed-upon level, what is the

best defecting price for the other? What are its resulting profits? For what interest rates will their collusion be sustainable if both are using the grim strategy?

9. Now we extend the analysis in Exercise S7 to allow for defecting in a collusive triopoly. Exercise S9 in [Chapter 5](#) finds the Nash equilibrium outcome of a VLCC triopoly of Korea, Japan, and China.
 1. Now find the collusive outcome of the triopoly. That is, what total quantity should be set by the three countries collectively in order to maximize their joint profit?
 2. Assume that under the collusive outcome found in part (a), the three countries produce equal quantities of VLCCs, so that each earns an equal share of the collusive profit. How much profit would each country earn? Compare this profit with the amount each earns in the Nash equilibrium outcome.
 3. Now suppose the three countries are in a repeated relationship. Once per year, they choose production quantities, and each can observe the quantity its rivals produced in the previous year. They wish to cooperate to sustain the collusive profit levels found in part (b). In any one year, one of them can defect from the agreement. If the other two countries are expected to produce their share of the collusive outcome found in parts (a) and (b), what is the best defecting quantity for the third to produce? What is the resulting profit for a defecting country when it produces the optimal defecting quantity while the other two produce their collusive quantities?
 4. Of course, the year after one country defects, both of its rivals will also defect. They will all find themselves back at the Nash equilibrium outcome (permanently, if they use the grim strategy). How much does the defecting country stand to gain in one year of defecting from the collusive outcome? How

much will the defecting country then lose in every subsequent year from earning the Nash equilibrium profit instead of the collusive profit?

5. For what interest rates will collusion be sustainable if the three countries are using the grim strategy? Is this set of interest rates larger or smaller than that found in the duopoly case discussed in Exercise S7, part (e)? Why?

Endnotes

- See <https://www.aspkin.com/forums/ebay-accounts-sale/77998-aged-verified-high-limit-usa-ebay-paypal-seller-accounts-sale-delivered-fast.xhtml> (accessed April 18, 2019). [Return to reference 23](#)
- Diego Gambetta, *The Sicilian Mafia* (Cambridge, Mass.: Harvard University Press, 1993), p.15. [Return to reference 24](#)
- James Andreoni and Hal Varian, “Preplay Contracting in the Prisoners’ Dilemma,” *Proceedings of the National Academy of Sciences*, vol. 96, no. 19 (September 14, 1999), pp. 10933 – 38. [Return to reference 25](#)

■ Appendix: Infinite Sums

The computation of present values requires us to determine the current value of a sum of money that is paid to us in the future. As we saw in [Section 2](#), the present value of a sum of money—say, x —that is paid to us n months from now is just $x/(1 + r)^n$, where r is the appropriate monthly rate of return. But the present value of a sum of money that is paid to us next month and every following month in the foreseeable future is more complicated to determine. In that case, the payments continue infinitely, so there is no defined end to the sum of present values that we need to compute. To compute the present value of this flow of payments requires some knowledge of the mathematics of the summation of infinite series.

Consider a player who stands to gain \$36 this month from defecting in a prisoners' dilemma, but who will then lose \$36 every month in the future as a result of her choice to continue defecting while her opponent (using the tit-for-tat, or TFT, strategy) punishes her. In the first of the future months—the first for which there is a loss and the first for which values need to be discounted—the present value of her loss is $36/(1 + r)$; in the second future month, the present value of the loss is $36/(1 + r)^2$; in the third future month, the present value of the loss is $36/(1 + r)^3$. That is, in each of the n future months that she incurs a loss from defecting, that loss equals $36/(1 + r)^n$.

We could write out the total present value of all of her future losses as a large sum with an infinite number of components,

$$\text{PV} = \frac{36}{1 + r} + \frac{36}{(1 + r)^2} + \frac{36}{(1 + r)^3} + \frac{36}{(1 + r)^4} + \frac{36}{(1 + r)^5} + \frac{36}{(1 + r)^6} + \dots,$$

or we could use summation notation as a shorthand device and instead write

$$\text{PV} = \sum_{n=1}^{\infty} \frac{36}{(1+r)^n}.$$

This expression, which is equivalent to the preceding one, is read as “the sum, from n equals 1 to n equals infinity, of 36 over $(1+r)$ to the n th power.” Because 36 is a common factor—it appears in each term of the sum—it can be pulled out to the front of the expression. Thus, we can write the same present value as

$$\text{PV} = 36 \times \sum_{n=1}^{\infty} \frac{1}{(1+r)^n}.$$

We now need to determine the value of the sum within this present-value expression to calculate the actual present value. To do so, we simplify our notation by using the *discount factor* δ in place of $1/(1+r)$. Then the sum that we are interested in evaluating is

$$\sum_{n=1}^{\infty} \delta^n.$$

It is important to note here that $\delta = 1/(1+r) < 1$ because r is strictly positive.

An expert on infinite sums would tell you, after inspecting this last sum, that it converges to the finite value $\delta/(1-\delta)$.^{[26](#)} Convergence is guaranteed because increasingly large powers of a number less than 1, δ in this case, become smaller and smaller, approaching zero as n approaches infinity. The later terms in our present value, then, decrease in size until they get sufficiently small that the series approaches (but technically never exactly reaches) the particular value of the sum. Although a good deal of

more sophisticated mathematics is required to deduce that the convergent value of the sum is $\delta/(1 - \delta)$, proving that this is the correct answer is relatively straightforward.

We use a simple trick to prove our claim. Consider the sum of the first m terms of the series, and denote it by S_m . Thus

$$S_m = \sum_{n=1}^{\infty} \delta^n = \delta + \delta^2 + \delta^3 + \cdots + \delta^{m-1} + \delta^m.$$

Now we multiply this sum by $(1 - \delta)$ to get

$$\begin{aligned} (1 - \delta)S_m &= \delta + \delta^2 + \delta^3 + \cdots + \delta^{m-1} + \delta^m \\ &\quad - \delta^2 - \delta^3 - \delta^4 - \cdots - \delta^m - \delta^{m-1} \\ &= \delta - \delta^{m-1}. \end{aligned}$$

Dividing both sides by $(1 - \delta)$, we have

$$S_m = \frac{\delta - \delta^{m+1}}{1 - \delta}.$$

Finally, we take the limit of this sum as m approaches infinity to evaluate our original infinite sum. As m goes to infinity, the value of δ^{m+1} goes to zero because very large and increasing powers of a number less than 1 get increasingly small, but stay nonnegative. Thus, as m goes to infinity, the right-hand side of the preceding equation goes to $\delta/(1 - \delta)$, which is therefore the limit of S_m as m approaches infinity. This completes the proof.

We need only convert back into r to be able to use our answer in the calculation of present values in our prisoners' dilemma games. Because $\delta = 1/(1 + r)$, it follows that

$$\frac{\delta}{1 - \delta} = \frac{1/(1 + r)}{r/(1 + r)} = \frac{1}{r}.$$

The present value of an infinite stream of \$36s earned each month, starting next month, is then

$$36 \times \sum_{n=1}^{\infty} \frac{1}{(1 + r)^n} = \frac{36}{r}.$$

This is the value that we use in [Section 2](#) to determine whether a player should defect forever. Notice that incorporating a probability of continuation, $p \leq 1$, into the discounting calculations changes nothing in the summation procedure used here. We could easily substitute R for r in the preceding calculations, and $p\delta$ for the discount factor, δ .

Remember that you need to find present values only for losses (or gains) incurred (or accrued) *in the future*. The present value of \$36 lost today is just \$36. So if you wanted the present value of a stream of losses, all of them \$36, that begins *today*, you would take the \$36 lost today and add it to the present value of the stream of losses in the future. We have just calculated that present value as $36/r$. Thus, the present value of the stream of lost \$36s, including the \$36 lost today, would be $36 + 36/r$, or $36[(r + 1)/r]$, which equals $36/(1 - \delta)$. Similarly, if you wanted to look at a player's stream of profits under a particular contingent strategy in a prisoners' dilemma, you would not discount the profit amount earned in the very first period; you would discount only those profit figures that represent money earned in future periods.

Endnotes

- An infinite series *converges* if the sum of the values in the series approaches a specific value, getting closer and closer to that value as additional components of the series are included in the sum. The series *diverges* if the sum of the values in the series gets increasingly larger (more negative) with each addition to the sum. Convergence requires that the components of the series get progressively smaller. [Return to reference 26](#)

11 ■ Collective-Action Games

THE GAMES AND STRATEGIC SITUATIONS considered in the preceding chapters have usually included only two or three players interacting with each other. Such games are common in our own academic, business, political, and personal lives and so are important to understand and analyze. But many social, economic, and political interactions are strategic situations in which numerous players participate at the same time. Strategies for career paths, investment plans, rush-hour commuting routes, and even studying have associated benefits and costs that depend on the actions of many other people. If you have been in any of these situations, you probably thought something was wrong—too many students, investors, and commuters crowding just where you wanted to be, for example. If you have tried to organize fellow students or your community in some worthy cause, you probably faced frustration of the opposite kind—too few willing volunteers. In other words, multiple-person games in a society often seem to produce outcomes that are not deemed satisfactory by many, or even all, of the people in that society. In this chapter, we examine such games from the perspective of the theory that we have already developed. We present an understanding of what goes wrong in such situations and what can be done about it.

In their most general form, such many-player games concern problems of [collective action](#). In these cases, the aims of the whole society, or *collective*, are best served if its members take some particular action or actions, but those actions are not in the best private interests of those individual members. In other words, the socially optimal outcome is not automatically achievable as the Nash equilibrium of the game. Therefore, we must examine how the game can be modified to lead to the optimal outcome, or at least to improve on an unsatisfactory Nash equilibrium. To do

so, we must first understand the nature of such games. We find that they come in three forms, all of them familiar to you by now: prisoners' dilemmas, games of chicken, and assurance games. Although our main focus in this chapter is on situations where numerous individuals play such games at the same time, we build on familiar ground by beginning with games between just two players.

Glossary

collective action

A problem of achieving an outcome that is best for society as a whole, when the interests of some or all individuals will lead them to a different outcome as the equilibrium of a noncooperative game.

1 COLLECTIVE-ACTION GAMES WITH TWO PLAYERS

Imagine that you are a farmer. A neighboring farmer and you can both benefit by constructing an irrigation and flood-control project. The two of you can join together to undertake this project, or one of you can do so on your own. However, after the project has been constructed, both of you will automatically benefit from it. Therefore, each is tempted to leave the work to the other. That is the essence of your strategic interaction and the difficulty of securing collective action.

In [Chapter 4](#), we encountered a game of this kind: Three neighbors were each deciding whether to contribute to a street garden that all of them would enjoy. That game became a prisoners' dilemma in which all three shirked. Our analysis here will include an examination of a broader range of possible payoff structures. Also, in the street-garden game, we rated the possible outcomes on a scale of 1 to 6; when we describe more general collective-action games, we will have to consider more general forms of benefits and costs for each player.

Our irrigation project has two important characteristics. First, its benefits are [nonexcludable](#): A person who has not contributed to paying for it cannot be prevented from enjoying the benefits. Second, its benefits are [nonrival](#): Any one person's benefits are not diminished by the mere fact that someone else is also getting those benefits. Economists call such a project a [pure public good](#); national defense is often given as an example. In contrast, a pure *private* good is fully excludable and rival: Nonpayers can be excluded from its benefits, and if one person gets those benefits, no one

else does. A loaf of bread is a good example of a pure private good. Most goods fall somewhere on the two-dimensional spectrum of varying degrees of excludability and rivalness. We will not go any deeper into this taxonomy, but we mention it to help you relate our discussion to what you may encounter in other courses and books.^{[1](#)}

A. Collective Action as a Prisoners' Dilemma

The costs and the benefits associated with building the irrigation project depend, as do those associated with all collective actions, on which players participate. In turn, the relative sizes of the costs and benefits determine the structure of the game that is played. Suppose either of you acting alone could complete the project in 7 weeks, whereas if the two of you acted together, it would take only 4 weeks of time from each. The two-person project is also of better quality; each of you gets benefits worth 6 weeks of work from a one-person project (whether constructed by you or by your neighbor) and 8 weeks' worth of benefits from a two-person project.

More generally, we can write benefits and costs, and therefore player payoffs, as functions of the number of players participating. So the cost to you of choosing to build the project depends on whether you build it alone or with help. Cost can be written as $C(n)$, where cost, C , depends on the number, n , of players participating in the project. Then $C(1)$ would be the cost to you of building the project alone and $C(2)$ would be the cost *to you* of building the project with your neighbor. Here, $C(1) = 7$ and $C(2) = 4$. Similarly, benefits, B , from the completed project may vary depending on how many players (n) participate in its completion. In our example, $B(1) = 6$ and $B(2) = 8$. Note that these benefits are the same for each of you, regardless of participation, due to the public-good nature of this particular project.

In this game, each of you has to decide whether to work toward the construction of the project or not—that is, to shirk. (Presumably, there is a short window of time in which

the work must be done, and you could pretend to be called away on some very important family matter at the last minute, as could your neighbor.) Figure 11.1 shows the payoffs for this game, measured in units equivalent to the value of a week of work. Payoffs are determined by the difference between the cost and the benefit associated with each action. So the payoff function for choosing Build can be written as $P(n) = B(n) - C(n)$, with payoff for participating, P , depending on the number of participants, n ; $n = 1$ if you build alone and $n = 2$ if your neighbor also chooses Build. The payoff function for choosing Not is $S(n) = B(n)$, with payoff for shirking, S , also depending on the number of participants n ; with $n = 1$ if your neighbor chooses Build. In this case, you incur no cost because you do not participate in the project.

		NEIGHBOR	
		Build	Not
YOU	Build	4, 4	-1, 6
	Not	6, -1	0, 0

FIGURE 11.1 Collective Action as a Prisoners’ Dilemma: Version I

		NEIGHBOR	
		Build	Not
YOU	Build	2.3, 2.3	-1, 6
	Not	6, -1	0, 0

FIGURE 11.2 Collective Action as a Prisoners’ Dilemma: Version II

Given the payoff structure in Figure 11.1, your best response if your neighbor does not participate is not to participate either: Your benefit from completing the project by yourself

(6) is less than your cost (7), for a net payoff of -1 , whereas you can get 0 by not participating. Similarly, if your neighbor does participate, then you can reap the benefit (6) from his work at no cost to yourself, which is better for you than working yourself to get the larger benefit of the two-person project (8) while incurring the cost of the work (4), for a net payoff of 4. The key feature of the game is that it is better for you not to participate no matter what your neighbor does; the same logic holds for him. (In this case, each farmer is said to be a [free rider](#) on his neighbor's effort if he lets the other do all the work and then reaps the benefits all the same.) Thus, not building is the dominant strategy for each of you. But both of you would be better off if you were to work together to build (payoff 4) than if neither of you builds (payoff 0). Therefore, the game is a prisoners' dilemma.

We see in this prisoners' dilemma one of the main difficulties that arises in games of collective action. Individually optimal choices—in this case, a farmer choosing not to build regardless of what the other farmer chooses—may not be optimal from the perspective of society as a whole, even if the society is made up of just two farmers. The [social optimum](#) in a collective-action game is achieved when the sum total of the players' payoffs is maximized; in this prisoners' dilemma, the social optimum is the (Build, Build) outcome. Nash equilibrium behavior by the players does not consistently bring about the socially optimal outcome, however. Hence, the study of collective-action games has focused on methods to improve on observed (generally Nash) equilibrium behavior to move outcomes toward the socially best ones. As we will see, the divergence between Nash equilibrium and socially optimal outcomes appears in every version of collective-action games.

Now consider what the game would look like if the numbers were to change slightly. Suppose the two-person project

yields benefits that are not much better than those from the one-person project: 6.3 weeks’ worth of work to each farmer. Then each of you gets a payoff of $6.3 - 4 = 2.3$ when both of you build. The resulting payoff table is shown in Figure 11.2. The game is still a prisoners’ dilemma and still leads to the equilibrium (Not, Not). However, when both of you build, the total payoff for both of you is only 4.6. The social optimum occurs when one of you builds and the other does not, in which case the total payoff is $6 + (-1) = 5$. There are two possible ways to get this outcome: You build and your neighbor shirks, or your neighbor builds and you shirk. Achieving the social optimum in this case then poses a new problem: Who should build and suffer the payoff of -1 while the other is allowed to be a free rider and enjoy the payoff of 6?

		NEIGHBOR	
		Build	Not
YOU	Build	5, 5	2, 6
	Not	6, 2	0, 0

FIGURE 11.3 Collective Action as Chicken: Version I

B. Collective Action as Chicken

Yet another variation in the numbers of the original prisoners' dilemma game of Figure 11.1 changes the nature of the game. Suppose the cost of the work is reduced so that it becomes better for you to build your own project if your neighbor does not. Specifically, suppose the one-person project requires 4 weeks of work, so $C(1) = 4$, and the two-person project takes 3 weeks from each, so $C(2) = 3$ (to each); the benefits are the same as before. Figure 11.3 shows the payoff matrix resulting from these changes. Now your best response is to shirk when your neighbor works and to work when he shirks. In form, this game is just like a game of chicken, where shirking is the Straight or Tough strategy and working is the Swerve or Weak strategy.

If the outcome of this game is one of its pure-strategy equilibria, your and the neighbor's payoffs sum to 8; this total is less than 10, the total of the payoffs that you could get if both of you build. That is, neither of the Nash equilibria provides as much benefit to the pair of you as a whole as that of your joint optimum where both of you choose to build. If the outcome of the chicken game is its mixed-strategy equilibrium, both of you will fare even worse than in either of the pure-strategy equilibria: The two expected payoffs will add up to something less than 8 (4, to be precise).

The collective-action game of chicken has another possible structure if we make some additional changes to the benefits associated with the project. As with version II of the prisoners' dilemma, suppose the two-person project is not much better than the one-person project. Then each farmer's benefit from the two-person project, $B(2)$, is only 6.3, whereas each still gets a benefit of $B(1) = 6$ from the one-

person project. We ask you to practice your skill by constructing the payoff table for this game. You will find that it is still a game of chicken—call it chicken II. It still has two pure-strategy Nash equilibria, in each of which only one farmer builds, but the sum of the payoffs when both build is only 6.6, whereas the sum when only one farmer builds is 8. The social optimum is for only one farmer to build, but each farmer prefers the equilibrium in which the other builds. This may lead to a new dynamic game in which each waits for the other to build. Or the original game might yield its mixed-strategy equilibrium, with its low expected payoffs.

		NEIGHBOR	
		Build	Not
YOU	Build	4, 4	−4, 3
	Not	3, −4	0, 0

FIGURE 11.4 Collective Action as an Assurance Game

C. Collective Action as Assurance

Finally, let us change the payoffs for version I of the prisoners' dilemma in a different way altogether, leaving the benefits of the two-person project and the costs of building as originally set out, but reducing the benefit of a one-person project to $B(1) = 3$. This change reduces your benefit as a free rider so much that now if your neighbor chooses Build, your best response is Build. Figure 11.4 shows the payoff table for this version of the game. It is now an assurance game with two pure-strategy equilibria: one where both of you participate and the other where neither of you does.

As in the chicken II version of the game, the socially optimal outcome here is one of the two Nash equilibria. But there is a difference. In chicken II, the two players differ in their preferences between the two equilibria, either of which achieves the social optimum. In the assurance game, both of them prefer the same equilibrium, and that is the sole socially optimal outcome. Therefore, achieving the social optimum should be easier in the assurance game than in chicken.

D. Collective Inaction

Many games of collective action have payoff structures that differ somewhat from those in our irrigation project example. Our farmers find themselves in a situation in which the social optimum entails that at least one, if not both, of them participates in the project. Thus the game is one of collective *action*. Other multiplayer games might better be called games of collective *inaction*. In such games, society as a whole prefers that some or all of the individual players do *not* participate or do *not* act. Examples of this type of interaction include choices among rush-hour commuting routes, investment plans, or fishing grounds.

All of these games have the attribute that players must decide whether to take advantage of some common resource, be it a freeway, a high-yielding stock fund, or an abundantly stocked pond. These “collective inaction” games are better known as *common-resource* games; the total payoff to all players reaches its maximum when players refrain from overusing the common resource. The difficulty associated with not being able to reach the social optimum in such games is known as the “tragedy of the commons,” a phrase coined by Garrett Hardin in his paper of the same name.^{[2](#)}

We supposed that the irrigation project yielded equal benefits to both farmers. But what if the outcome of both farmers’ building was that the project used so much water that the farms had too little water for their livestock? Then each player’s payoff could be negative when both choose Build, lower than when both choose Not. This would be yet another variant of the prisoners’ dilemma we encountered in [Section 1.A](#), in which the socially optimal outcome entails neither farmer’s building even though each one still has an individual incentive to do so. Or suppose that one farmer’s

activity causes harm to the other, as would happen if the only way to prevent one farm from being flooded is to divert the water to the other. Then each player's payoffs could be negative if his neighbor chooses Build. Thus, another variant of chicken could also arise. In this variant, each farmer wants to build when the other does not, whereas it would be collectively better if neither of them did.

Just as the problems pointed out in these examples of both collective action and collective inaction are familiar, the various ways of tackling these problems also follow the general principles discussed in earlier chapters. Before turning to solutions, let us see how the problems manifest themselves in the more realistic setting where several players interact simultaneously in such games.

Endnotes

- Public goods are studied in more detail in textbooks on *public economics*, such as Jonathan Gruber, *Public Finance and Public Policy*, 5th ed. (New York: Worth, 2015); Harvey Rosen and Ted Gayer, *Public Finance*, 10th ed. (Chicago: Irwin/McGraw-Hill, 2014); and Joseph Stiglitz and Jay Rosengard, *Economics of the Public Sector*, 4th ed. (New York: W. W. Norton, 2015). [Return to reference 1](#)
- Garrett Hardin, “The Tragedy of the Commons,” *Science*, vol. 162 (1968), pp. 1243 – 48. [Return to reference 2](#)

Glossary

nonexcludable

Benefits that are available to each individual, regardless of whether he has paid the costs that are necessary to secure the benefits.

nonrival

Benefits whose enjoyment by one person does not detract anything from another person's enjoyment of the same benefits.

pure public good

A good or facility that benefits all members of a group, when these benefits cannot be excluded from a member who has not contributed efforts or money to the provision of the good, and the enjoyment of the benefits by one person does not significantly detract from their simultaneous enjoyment by others.

free rider

A player in a collective-action game who intends to benefit from the positive externality generated by others' efforts without contributing any effort of his own.

social optimum

In a collective-action game where payoffs of different players can be meaningfully added together, the social optimum is achieved when the sum total of the players' payoffs is maximized.

2 COLLECTIVE-ACTION PROBLEMS IN LARGE GROUPS

In this section, we extend our irrigation project example to a situation in which a population of N farmers must each decide whether to participate. Here we make use of the notation we introduced above, with $C(n)$ representing the cost each participant incurs when n of the N total farmers have chosen to participate. Similarly, the benefit to each farmer, regardless of participation, is $B(n)$. Each participant's payoff function is $P(n) = B(n) - C(n)$, whereas the payoff function for each nonparticipant, or shirker, is $S(n) = B(n)$.

Suppose you are contemplating whether to participate or to shirk. Your decision will depend on what the other $(N - 1)$ farmers in the population are doing. In general, you will have to make your decision when the other $(N - 1)$ players consist of n participants and $(N - 1 - n)$ shirkers. If you decide to shirk, the number of participants in the project is still n , so you get a payoff of $S(n)$. If you decide to participate, the number of participants becomes $n + 1$, so you get $P(n + 1)$. Therefore, your final decision depends on the comparison of these two payoffs; you will participate if $P(n + 1) > S(n)$, and you will shirk if $P(n + 1) < S(n)$. This comparison holds true for every version of the collective-action game analyzed in [Section 1](#); the differences in player behavior in the different versions arise because the changes in the payoff structure alter the values of $P(n + 1)$ and $S(n)$.

		NEIGHBOR	
		Build	Not
YOU	Build	$P(2), P(2)$	$P(1), S(1)$
	Not	$S(1), P(1)$	$S(0), S(0)$

FIGURE 11.5 General Form of a Two-Person Collective-Action Game

We can relate the two-player examples in [Section 1](#) to this more general framework. If there are just two players, then $P(2)$ is the payoff to one from building when the other also builds, $S(1)$ is the payoff to one from shirking when the other builds, and so on. Therefore, we can generalize the payoff tables of Figures 11.1 through 11.4 into an algebraic form. This general payoff structure is shown in Figure 11.5.

The game illustrated in Figure 11.5 is a prisoners' dilemma if the inequalities

$$P(2) < S(1), P(1) < S(0), P(2) > S(0)$$

all hold at the same time. The first says that the best response to Build is Not, the second says that the best response to Not also is Not, and the third says that (Build, Build) is jointly preferred to (Not, Not). The game is equivalent to version I of the prisoners' dilemma if $2P(2) > P(1) + S(1)$, so that the total payoff is higher when both build than when only one builds. You can establish similar inequalities among these payoffs that yield the other types of games in [Section 1](#).

Return now to the multiplayer version of the game with a general n . Given the payoff functions for the two actions, $P(n+1)$ and $S(n)$, we can use graphs to help us determine which type of game we have encountered and find its Nash equilibrium. We can also then compare the Nash equilibrium with the game's socially optimal outcome.

A. Multiplayer Prisoners' Dilemma

Take a specific version of our irrigation project example in which an entire village of 100 farmers must decide which action to take. Suppose that the project raises the productivity of each farmer's land in proportion to the size of the project; specifically, suppose the benefit to each farmer when n farmers work on the project is $P(n) = 2n$. Suppose also that if you are not working on the project, you can enjoy this benefit and use your extra four weeks of time in some other occupation, so $S(n) = 2n + 4$. Remember that your decision about whether to participate in the project depends on the relative magnitudes of $P(n + 1) = 2(n + 1)$ and $S(n) = 2n + 4$. We draw the two separate graphs of these functions for an individual farmer in Figure 11.6, showing n over its full range from 0 to $(N - 1)$ along the horizontal axis and the payoff to the farmer along the vertical axis. If there are currently very few participants (thus mostly shirkers), your choice will depend on the relative locations of $P(n + 1)$ and $S(n)$ on the left end of the graph. Similarly, if there are already many participants, your choice will depend on the relative locations of $P(n + 1)$ and $S(n)$ on the right end of the graph.

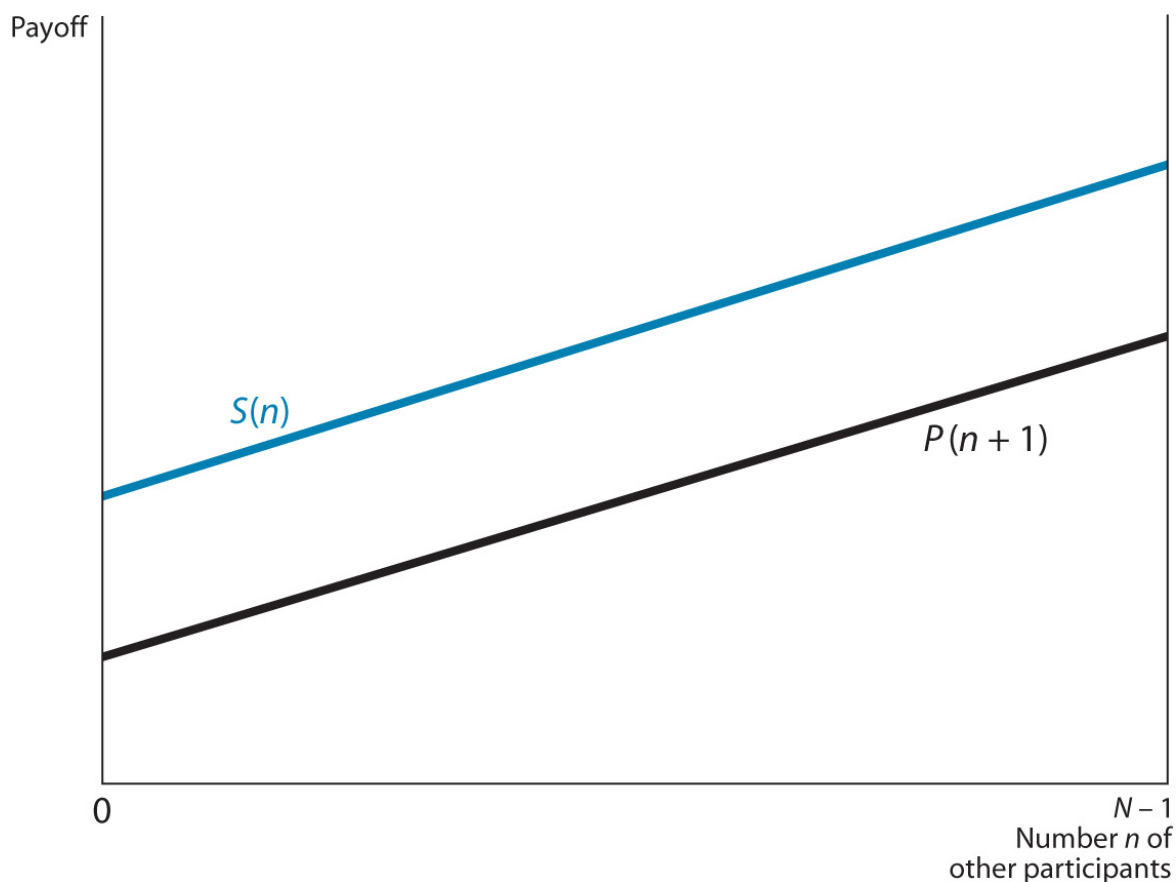


FIGURE 11.6 Multiplayer Prisoners' Dilemma Payoff Graph

Because n can have only integer values, each function $P(n+1)$ and $S(n)$ technically consists only of a discrete set of points, rather than a continuous set as implied by our smooth curves (which in this instance happen to be straight lines). But when N is large, the discrete points are sufficiently close together that we can connect the successive points and show each payoff function as a continuous curve. We use linear $P(n+1)$ and $S(n)$ functions in this section to bring out the basic considerations involved in collective-action games in the simplest possible way; we will discuss more complicated possibilities later.

Recall that you determine your choice of action by considering the number of current participants in the project, n , and the payoffs associated with each action at that n . Figure 11.6 illustrates a case in which the curve $S(n)$ lies entirely above the curve $P(n+1)$. Therefore, no matter how many others participate (that is, no matter how large n gets), your payoff is

higher if you shirk than if you participate; shirking is your dominant strategy. These payoffs are identical for all players, so everyone has a dominant strategy to shirk. Therefore, the Nash equilibrium of the game entails everyone shirking, and the project is not built.

Note that both curves rise as n increases. Thus, for each action you can take, you are better off if more of the others participate in the project. And the left intercept of the $S(n)$ curve is below the right intercept of the $P(n + 1)$ curve, so $S(0) = 4 < P(N) = 102$. This inequality says that if everyone, including you, shirks, your payoff is less than if everyone, including you, participates. Everyone would be better off than they are in the Nash equilibrium of the game if the outcome in which everyone participates could be sustained. This makes the game a prisoners' dilemma.

How does the Nash equilibrium found using the curves in Figure 11.6 compare with the social optimum of this game? To answer this question, we need a way to describe the total social payoff at each value of n ; we do that by using the payoff functions $P(n)$ and $S(n)$ to construct a third function $T(n)$, showing the total payoff to society as a function of n . The total payoff to society when there are n participants consists of the value $P(n)$ for each of the n participants and the value $S(n)$ for each of the $(N - n)$ shirkers:

$$T(n) = nP(n) + (N - n) S(n).$$

The social optimum occurs when the allocation of people between participants and shirkers maximizes the total payoff $T(n)$, or at the number of participants—that is, the value of n —that maximizes $T(n)$. To get a better understanding of where this optimum might be, it is convenient to write $T(n)$ differently, rearranging the expression above to get

$$T(n) = NS(n) - n[S(n) - P(n)].$$

This version of the total social payoff function shows that we can calculate it as if we gave every one of the N farmers the

shirker's payoff, but then removed the shirker's extra benefit $[S(n) - P(n)]$ from each of the n participants.

In collective-action games, as opposed to common-resource games, we normally expect $S(n)$ to increase as n increases. Therefore, the first term in this expression, $NS(n)$, also increases as n increases. If the second term does not increase too fast as n increases—as would be the case if the shirker's extra benefit, $[S(n) - P(n)]$, is small and constant—then the effect of the first term dominates in determining the value of $T(n)$.

This is exactly what happens with the total social payoff function for our current 100-farmer example. Here $T(n) = nP(n) + (N - n)S(n)$ becomes $T(n) = n(2n) + (100 - n)(2n + 4) = 2n^2 + 200n - 2n^2 + 400 - 4n = 400 + 196n$. In this case, $T(n)$ increases steadily with n and is maximized at $n = N$ when no one shirks.

Thus, the large-group version of this game holds the same lesson as our two-person example: Society as a whole would be better off if all of the farmers participated in building the irrigation project (if $n = N$). But payoffs are such that each farmer has an individual incentive to shirk. The Nash equilibrium of the game, at $n = 0$, is not socially optimal. Figuring out how to achieve the social optimum is one of the most important topics in the study of collective action, to which we return later in this chapter.

In other situations, $T(n)$ could be maximized for a different value of n , not just at $n = N$. That is, society's aggregate payoff could be maximized by allowing some shirking. Even in the prisoners' dilemma case, it is not automatic that the total payoff function is maximized when n is as large as possible. If the gap between $S(n)$ and $P(n)$ widens sufficiently fast as n increases, then the negative effect of the second term in the expression for $T(n)$ outweighs the positive effect of the first term as n approaches N . Then it may be best to let some people shirk—that is, the socially optimal value for n may be less than N . This result mirrors that of our prisoners' dilemma version II in [Section 1](#).

This type of outcome would arise in our village if $S(n)$ were $4n + 4$, rather than $2n + 4$. Then $T(n) = -2n^2 + 396n + 400$, which is no longer a linear function of n . In fact, a graphing calculator or some basic calculus shows that this $T(n)$ is maximized at $n = 99$, rather than at $n = 100$ as was true before. The change to the payoff structure has created an inequality in the payoffs—the shirkers fare better than the participants—which adds another dimension of difficulty to society's attempts to resolve the dilemma. How, for example, would the village designate exactly one farmer to be the shirker?

B. Multiplayer Chicken

Now we consider some of the other configurations that can arise in the payoffs. For example, when $P(n) = 4n + 36$, so that $P(n + 1) = 4n + 40$, and $S(n) = 5n$, the two payoff curves will cross in the graph. This case is illustrated in Figure 11.7. Here, for small values of n , $P(n + 1) > S(n)$, so if few others are participating, your optimal choice is to participate. For large values of n , $P(n + 1) < S(n)$, so if many others are participating, your optimal choice is to shirk. Note the equivalence of these two statements to the idea in the two-person chicken game that you shirk if your neighbor works and you work if he shirks. This case is indeed a game of chicken. More generally, a collective-action game of chicken occurs when you are given a choice between two actions, and you prefer to do the one that most others are *not* doing.

We can also use Figure 11.7 to determine the location of the Nash equilibrium of this version of the game. Because you choose to participate when n is small and to shirk when n is large, the equilibrium must be some intermediate value of n . Only at that n where the two curves intersect are you indifferent between your two choices. This location represents the equilibrium value of n . In our graph, $P(n + 1) = S(n)$ when $4n + 40 = 5n$ or when $n = 40$; that is the Nash equilibrium number of farmers from the village who will participate in the irrigation project.

If the two curves intersect at a point corresponding to an integer value of n , then that is the Nash equilibrium number of participants. If that is not the case, then strictly speaking, the game has no Nash equilibrium. But in practice, if the current value of n in the population is the integer just to the left of the point of intersection, then one more farmer will just want to participate, whereas if the current value of n is the integer just to the right of the point of intersection, one farmer will want to switch to shirking. Therefore, the number of participants will stay within a small range around the point of intersection,

and we can justifiably speak of the intersection as the equilibrium in some approximate sense.

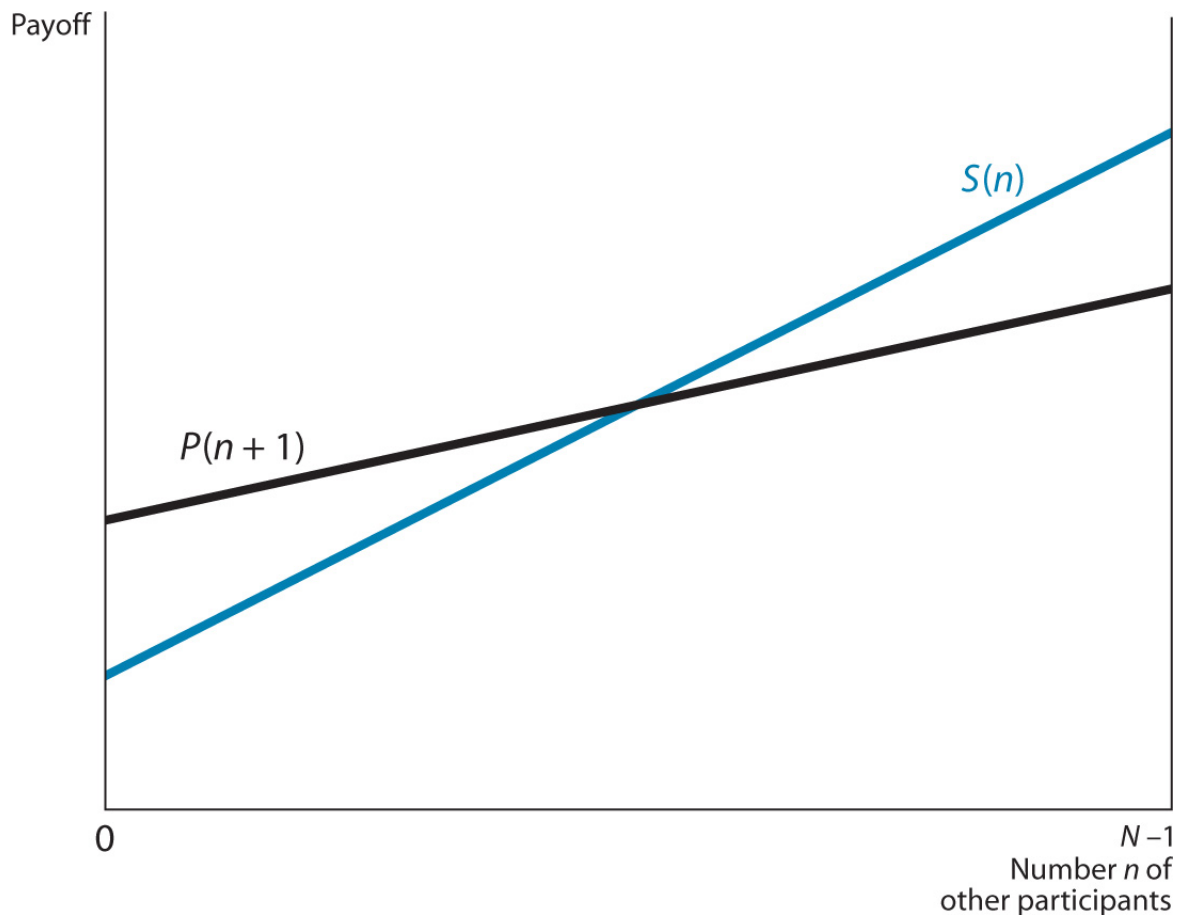


FIGURE 11.7 Multiplayer Chicken Payoff Graph

The payoff structure illustrated in Figure 11.7 shows both curves positively sloped, although they don't have to be. It is conceivable that the benefit for each person is smaller when more people participate, so the curves could be negatively sloped instead. The important feature of the collective-action game of chicken is that when few people are taking one action, it is better for any one person to take that action; and when many people are taking one action, it is better for any one person to take the other action.

What is the socially optimal outcome in the collective-action game of chicken? If each participant's payoff $P(n)$ increases as the number of participants increases, and if each shirker's

payoff $S(n)$ does not become too much greater than the $P(n)$ of each participant, then the total social payoff is maximized when everyone participates. This is the outcome in our example, where $T(n) = 536n - n^2$; total social payoff continues to increase for values of n beyond N (100 here), so $n = N$ is the social optimum.

But more generally, some cases of collective-action chicken entail social optima in which it is better to let some shirk. If our group of farmers numbered 300 instead of 100, our example here would yield such an outcome. The socially optimal number of participants, found on a graphing calculator or using calculus, would be 268. The idea that the social optimum can sometimes be attained when only some people work and the rest shirk emerged in the difference between versions I and II of chicken in our examples in [Section 1](#); here, we see an example of the result in a larger population. For an exercise, you may try generating a payoff structure that leads to such an outcome for our village of 100 farmers. By changing the payoff functions, $P(n)$ and $S(n)$, one can find games of chicken where the socially optimal number of participants could even be smaller than that in the Nash equilibrium. We return to examine the question of the social optima of all of these versions of the game in greater detail in [Section 3](#).

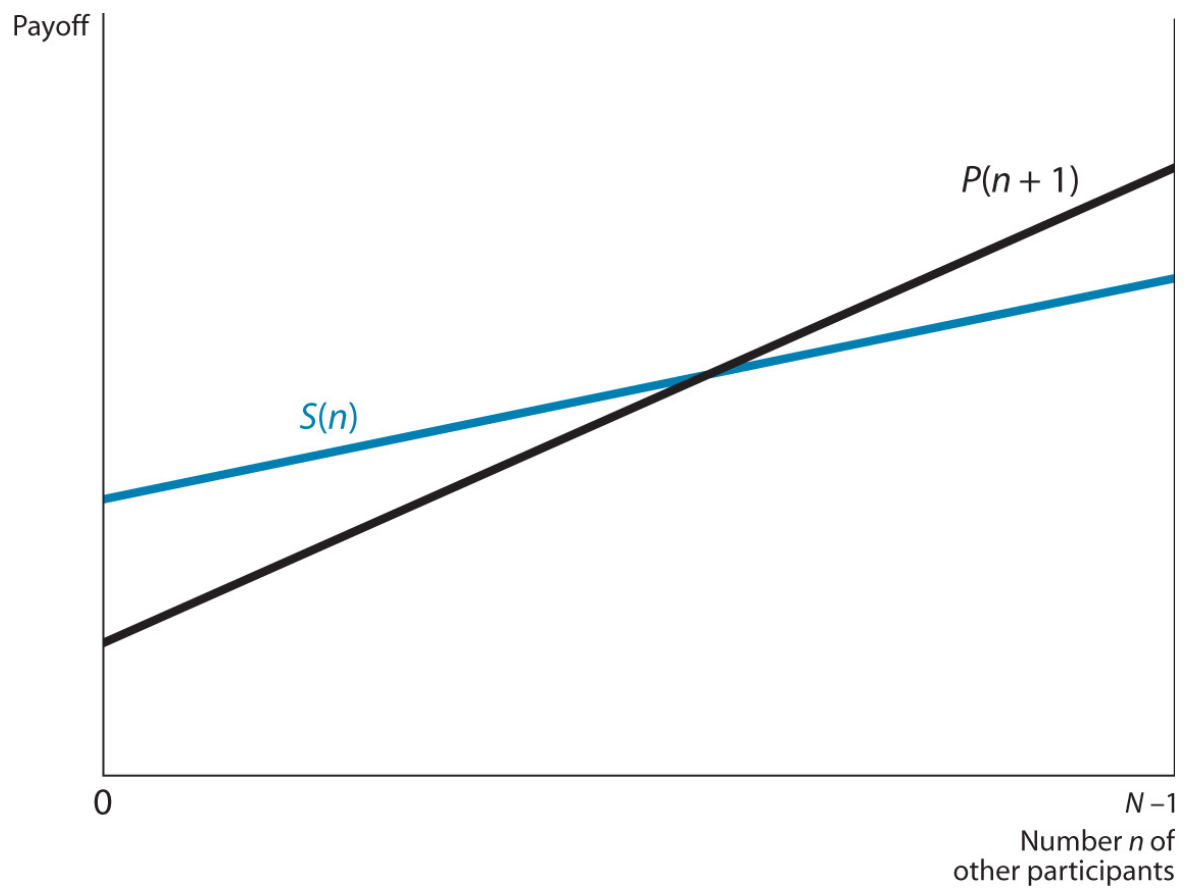


FIGURE 11.8 Multiplayer Assurance Payoff Graph

C. Multiplayer Assurance

Finally, we consider the third possible type of collective-action game: assurance. Figure 11.8 shows the payoff curves for the assurance case, where we suppose that the village farmers get $P(n+1) = 4n + 4$ and $S(n) = 2n + 100$. Here, $S(n) > P(n+1)$ for small values of n , so if few others are participating, then you want to shirk, too. But $P(n+1) > S(n)$ for large values of n , so if many others are participating, then you want to participate, too. In other words, unlike chicken, assurance is a collective-action game in which you want to make the choice that the others are making.

Figure 11.8 looks nearly identical to Figure 11.7, except that the S and P curves have swapped positions. In Figure 11.7, the S curve starts lower and is steeper than the P curve, but in Figure 11.8, the S curve starts higher and is flatter than the P curve. These details of curve placement might seem inconsequential, but they illustrate the essence of the difference between chicken games and assurance games. In a chicken game, each farmer wants to participate when others shirk (graphically, the P curve starts out higher than the S curve), but each has less to gain from participating as more participate (graphically, the P curve rises more slowly than the S curve). Eventually, when the number of other participants in a chicken game is sufficiently large, each farmer prefers to shirk (graphically, the S curve is higher than the P curve once n becomes sufficiently large). On the other hand, in an assurance game, each farmer wants to shirk when others shirk also (graphically, the S curve starts out higher than the P curve), but each has more to gain from participating as more participate (graphically, the P curve rises more quickly than the S curve). Eventually, when the number of other participants in an assurance game is sufficiently large, each farmer prefers to participate (graphically, the P curve is higher than the S curve once n becomes sufficiently large). More simply put, in a multiplayer chicken game, players prefer to “go against the crowd,” while in a multiplayer assurance game, they prefer to “go with the crowd.”

What are the Nash equilibria of the assurance game in Figure 11.8? For any initial value of n to the left of the intersection, each farmer will want to shirk; so there is a Nash equilibrium at $n = 0$, where everyone shirks. But the opposite is true to the right of the intersection. In that portion of the graph, each farmer will want to participate, yielding a second Nash equilibrium at $n = N$, where everyone participates.

Technically, there is also a third Nash equilibrium of this game if the value of n at the intersection is an integer value, as it is in our example. There we find that $P(n + 1) = 4n + 4 = 2n + 100 = S(n)$ when $n = 48$. Thus, if n were exactly 48, we would see an outcome in which there were some participants and some shirkers. This situation could be an equilibrium only if the value of n were exactly right. Even then, it would be a highly unstable situation. If any one farmer accidentally joined the wrong group, his choice would alter the incentives for everyone else, driving the game to one of the (stable) endpoint equilibria. (We will discuss the concept of equilibrium “stability” in detail in [Chapter 12](#), in the context of evolutionary games.)

The social optimum in this game is fairly easy to see in Figure 11.8. Because both curves are rising—so that each person is better off if more people participate—the right-hand extreme equilibrium is clearly the better one for society. This is confirmed in our example by noting that $T(n) = 2n^2 + 100n + 10,000$, which is an increasing function of n for all positive values of n ; thus, the socially optimal value of n is the largest one possible, or $n = N$. In the assurance case, then, the socially optimal outcome is actually one of the stable Nash equilibria of the game. As such, it may be easier to achieve than in some of the other cases. The critical question regarding the social optimum, regardless of whether it represents a Nash equilibrium of the underlying game, is how to bring it about.

So far, our examples have focused on relatively small groups of 2 or 100 persons. When the total number of people in the population, N , is very large, however, and any one person makes only a very small difference, then $P(n + 1)$ is almost the same as

$P(n)$. Thus, the condition under which any one person chooses to shirk is $P(n) < S(n)$. Expressing this inequality in terms of the benefits and costs of the common project in our example—namely, $P(n) = B(n) - C(n)$ and $S(n) = B(n)$ —we see that $P(n)$ [unlike $P(n + 1)$ in our preceding calculations] is *always* less than $S(n)$; individual persons will *always* want to shirk when N is very large. That is why problems of collective provision of public projects in a large group almost always manifest themselves as prisoners' dilemmas. But as we have seen, this result is not necessarily true for smaller groups. Neither is it true for large groups in other contexts such as traffic congestion, a case we discuss later in this chapter.

In general, we must allow for a broader interpretation of the payoffs $P(n)$ and $S(n)$ than we did in the specific case involving the benefits and the costs of a project. We cannot assume, for example, that the payoff functions will be linear. In fact, in the most general case, $P(n)$ and $S(n)$ can be any functions of n and can intersect many times. Then there can be multiple equilibria, although each can be thought of as representing one of the types described so far.³ And some games will be of the common-resource type as well, so when we allow for completely general games, we will speak of two actions labeled P and S , which have no necessary connotation of “participation” and “shirking,” but allow us to continue with the same symbols for the payoffs. Thus, when n players are taking the action P , $P(n)$ becomes the payoff of each player taking the action P , and $S(n)$ becomes that of each player taking the action S .

Endnotes

- Several exercises at the end of this chapter present examples of simple situations with nonlinear payoff curves and multiple equilibria. For a more general analysis and classification of such diagrams, see Thomas Schelling, *Micromotives and Macrobehavior* (New York: W. W. Norton & Company, 1978), Chapter. 7. The theory can be taken further by allowing each player a continuous choice (for example, the number of hours of participation) instead of just a binary choice of whether to participate. Many such situations are discussed in more specialized books on collective action, for example, Todd Sandler, *Collective Action: Theory and Applications* (Ann Arbor: University of Michigan Press, 1993), and Richard Cornes and Todd Sandler, *The Theory of Externalities, Public Goods, and Club Goods*, 2nd ed. (New York: Cambridge University Press, 1996). [Return to reference](#)

3 SPILLOVER EFFECTS, OR EXTERNALITIES

So far, we have seen that collective-action games occur in prisoners' dilemma, chicken, and assurance forms. We have also seen that the Nash equilibria in such games rarely yield the socially optimal level of participation (or nonparticipation). And even when the social optimum is a Nash equilibrium, it is usually only one of several equilibria that may arise. Here we delve further into the differences between the individual (or private) incentives in such games and the group (or social) incentives. We also describe more carefully the effects of each individual's decision on other individuals as well as on the collective. This analysis makes explicit why differences in incentives exist, how they are manifested, and how one might go about achieving socially better outcomes than those that arise in Nash equilibrium.

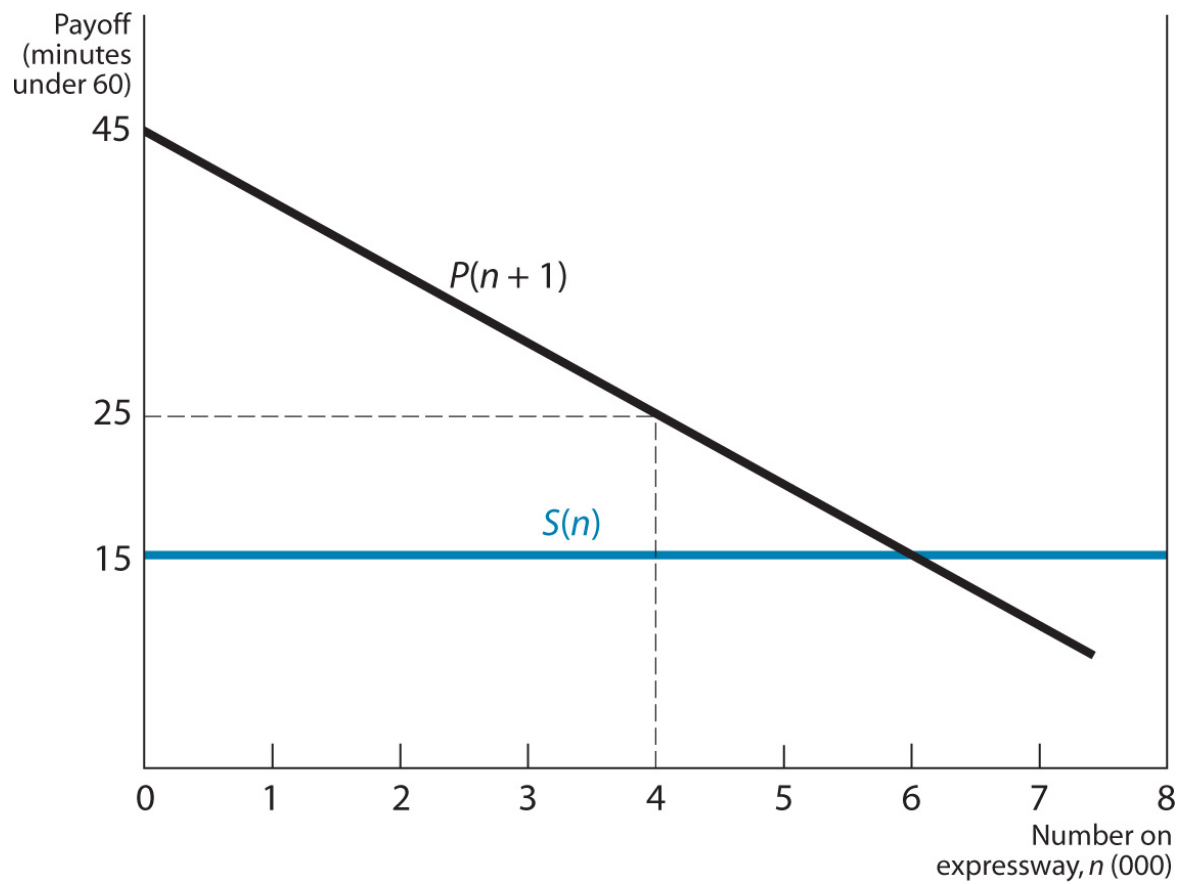


FIGURE 11.9 Commuting Route-Choice Game

A. Commuting and Spillover Effects

We start by thinking about a large population of 8,000 commuters who drive every day from a suburb to the city and back. As one of these commuters, you can take either the expressway (action P) or a network of local roads (action S). The local-roads route takes a constant 45 minutes, no matter how many cars are going that way. The expressway takes only 15 minutes when uncongested. But every driver who chooses the expressway increases the time for every other driver on the expressway by 0.005 minutes (about one-quarter of a second).

We can measure the payoffs in minutes of time saved—by how much the commute time is less than 1 hour, for instance. Then the payoff to drivers on the local roads, $S(n)$, is a constant $60 - 45 = 15$, regardless of the value of n . But the payoff to drivers on the expressway, $P(n)$, depends on n ; in particular, $P(n) = 60 - 15 = 45$ for $n = 0$, but $P(n)$ decreases by $5/1,000$ (or $1/200$) for every commuter on the expressway. Thus, $P(n) = 45 - 0.005n$. We graph the two payoff curves in Figure 11.9.

Suppose that initially 4,000 cars are on the expressway; $n = 4,000$. With so many cars on that road, it takes each of them $15 + 4,000 \times 0.005 = 15 + 20 = 35$ minutes to commute to work; each gets a payoff of $P(n) = 25$ [which is $60 - 35$, or $P(4,000)$]. As shown in Figure 11.9, that payoff is better than what local-road drivers obtain. You, a local-road driver, might therefore decide to switch from driving the local roads to driving on the expressway. Your switch would increase by 1 the value of n and would thereby affect the payoffs of all the other commuters. There would now be 4,001 drivers (including you) on the expressway, and the commute time for each would be 35 and $1/200$, or 35.005, minutes; each would now get a payoff of $P(n + 1) = P(4,001) = 24.995$. This payoff is still higher than the 15 from driving on the local roads. Thus, you have a private incentive to make the switch, because for you, $P(n + 1) > S(n)$ ($24.995 > 15$).

Your switch yields you a *private* gain—because it is privately enjoyed by you—equal to the difference between your payoffs before and after the switch; this private gain is $P(n + 1) - S(n) = 9.995$ minutes. Because you are only one person and therefore a small part of the whole group, the gain in payoff that you receive in relation to the total group payoff is small, or *marginal*. Thus, we call your gain the marginal private gain associated with your switch.

But now the 4,000 other drivers on the expressway each take 0.005 minutes more as a result of your decision to switch; the payoff to each changes by $P(4,001) - P(4,000) = -0.005$. Similarly, the drivers on the local roads face a payoff change of $S(4,001) - S(4,000)$, but this is zero in our example. The cumulative effect on all these other drivers is $4,000 \times -0.005 = -20$ (minutes). Your action, switching from local roads to expressway, has caused this effect on the others' payoffs. Such an effect of one person's action on others' payoffs is called a spillover effect, external effect, or externality. Again, because you are only a very small part of the whole group, we should actually call your effect on others the *marginal spillover effect*.

Taken together, the marginal private gain and the marginal spillover effect are the full effect of your switch on the population of commuters, or the overall marginal change in the whole society's payoff. We call this the marginal social gain associated with your switch. This “gain” may actually be positive or negative, so the use of the word *gain* is not meant to imply that all switches will benefit society as a whole. In fact, in our commuting example, the marginal social gain is $9.995 - 20 = -10.005$ (minutes). Thus, the overall social effect of your switch is negative; the total social payoff is reduced by just over 10 minutes.

B. Spillover Effects: The General Case

We can describe the effects we observe in the commuting example more generally by returning to our total social payoff function, $T(n)$, where n represents the number of people choosing P , so $N - n$ is the number of people choosing S . Suppose that initially n people have chosen P and that one person switches from S to P . Then the number choosing P increases by 1 to $(n + 1)$, and the number choosing S decreases by 1 to $(N - n - 1)$, so the total social payoff becomes

$$T(n + 1) = (n + 1) P(n + 1) + [N - (n + 1)] S(n + 1).$$

The increase in the total social payoff is the difference between $T(n)$ and $T(n + 1)$:

$$\begin{aligned} T(n + 1) - T(n) &= (n + 1) P(n + 1) + [N - (n + 1)] S(n + 1) - \\ & nP(n) + (N - n) S(n) = [P(n + 1) - S(n)] + n[P(n + 1) - P(n)] + \\ & [N - (n + 1)] [S(n + 1) - S(n)] \end{aligned}$$

(11.1)

after collecting and rearranging terms.

Equation (11.1) describes mathematically the various different effects of one person's switch from S to P that we saw earlier in the commuting example. The equation shows how the marginal social gain is divided into the marginal changes in payoffs for the subgroups of the population.

The first of the three terms in equation (11.1)—namely, $[P(n + 1) - S(n)]$ —is the marginal private gain enjoyed by the person who switches. As we saw above, this term is what drives a person's choice, and all such individual choices then determine the Nash equilibrium.

The second and third terms in equation (11.1) are just the quantifications of the spillover effects of one person's switch on the others in the population. For the n other people choosing

P , each sees her payoff change by the amount $[P(n+1) - P(n)]$ when one more person switches to P ; this spillover effect is seen in the second group of terms in equation (11.1). There are also $N - (n+1)$ (or $N - n - 1$) others still choosing S after the one person switches, and each of these players sees her payoff change by $[S(n+1) - S(n)]$; this spillover effect is shown in the third group of terms in the equation. Of course, the effect that one driver's switch has on commuting time for any other driver on either route is very small, but when there are numerous other drivers (that is, when N is large), the full spillover effect can be substantial.

Thus, we can rewrite equation (11.1) for a switch by one person either from S to P or from P to S as

Marginal social gain = marginal private gain + marginal spillover effect.

For an example in which one person switches from S to P , we have

Marginal social gain = $T(n+1) - T(n)$,

Marginal private gain = $P(n+1) - S(n)$, and

Marginal spillover effect = $n[P(n+1) - P(n)] + [N - (n+1)][S(n+1) - S(n)]$.

USING CALCULUS FOR THE GENERAL CASE Before examining some spillover situations in more detail to see what can be done to achieve socially better outcomes, we restate the general concepts of the analysis of collective-action games in the language of calculus. If you do not know this language, you can omit the remainder of this section without loss of continuity; if you do know it, you will find the alternative statement much simpler to grasp and to use than the algebra employed earlier.

If the total number N of people in the group is very large—say, in the hundreds or thousands—then one person can be regarded as a very small, or infinitesimal, part of this whole. This allows us to treat the number n as a continuous variable. If $T(n)$ is the total social payoff, we calculate the effect of changing n by

considering an increase of an infinitesimal marginal quantity dn , instead of a full unit increase from n to $(n + 1)$. To the first order, the change in payoff is $T'(n) dn$, where $T'(n)$ is the derivative of $T(n)$ with respect to n . Using the expression for the total social payoff,

$$T(n) = nP(n) + (N - n) S(n),$$

and differentiating, we have

$$\begin{aligned} T'(n) &= P(n) + nP'(n) - S(n) + (N - n) S'(n) \\ &= [P(n) - S(n)] + nP'(n) + (N - n) S'(n). \end{aligned} \quad (11.2)$$

This equation is the calculus equivalent of equation (11.1). $T'(n)$ represents the marginal social gain. The marginal private gain is $P(n) - S(n)$, which is just the change in the payoff of the person making the switch from S to P . In equation (11.1), we had $P(n + 1) - S(n)$ for this change in payoff; now we have $P(n) - S(n)$. This is because the infinitesimal addition of dn to the group of the n people choosing P does not change the payoff to any one of them by a significant amount. However, the total change in their payoff, $nP'(n)$, is sizable, and is recognized in the calculation of the spillover effect [it is the second term in equation (11.2)], as is the change in the payoff of the $(N - n)$ people choosing S [namely, $(N - n) S'(n)$], the third term in equation (11.2). These last two terms constitute the marginal spillover effect in equation (11.2).

In the commuting example, we had $P(n) = 45 - 0.005n$, and $S(n) = 15$. Then, with the use of calculus, we see that the private marginal gain for each driver who switches to the expressway when n drivers are already using it is $P(n) - S(n) = 30 - 0.005n$. Because $P'(n) = -0.005$ and $S'(n) = 0$, the spillover effect is $n \times (-0.005) + (N - n) \times 0 = -0.005n$, which equals -20 when $n = 4,000$. The answer is the same as before, but calculus simplifies the derivation and helps us find the optimum directly.

C. Commuting Revisited: Negative Externalities

A negative externality exists when the action of one person *lowers* others' payoffs—that is, when it imposes some extra costs on the rest of society. We saw this in our commuting example, where the marginal spillover effect of one person's switch to the expressway was negative, entailing an extra 20 minutes of drive time for other commuters. But the individual who changes her route to work does not take the spillover—the externality—into account when making her choice. She is motivated only by her own payoffs. (Remember that any guilt that she may suffer from harming others should already be reflected in her payoffs.) She will change her action from *S* to *P* as long as this change has a positive marginal *private* gain. She is then made better off by the change.

But society would be better off if the commuter's decision were governed by the marginal *social* gain. In our example, the marginal social gain is negative (-10.005), but the marginal private gain is positive (9.995), so the individual driver makes the switch, even though society as a whole would be better off if she did not do so. More generally, in situations with negative externalities, the marginal social gain will be smaller than the marginal private gain due to the existence of the negative spillover effect. Individuals will make decisions based on a cost-benefit calculation that is the wrong one from society's perspective. As a result, individuals will choose actions with negative spillover effects more often than society would like them to do.

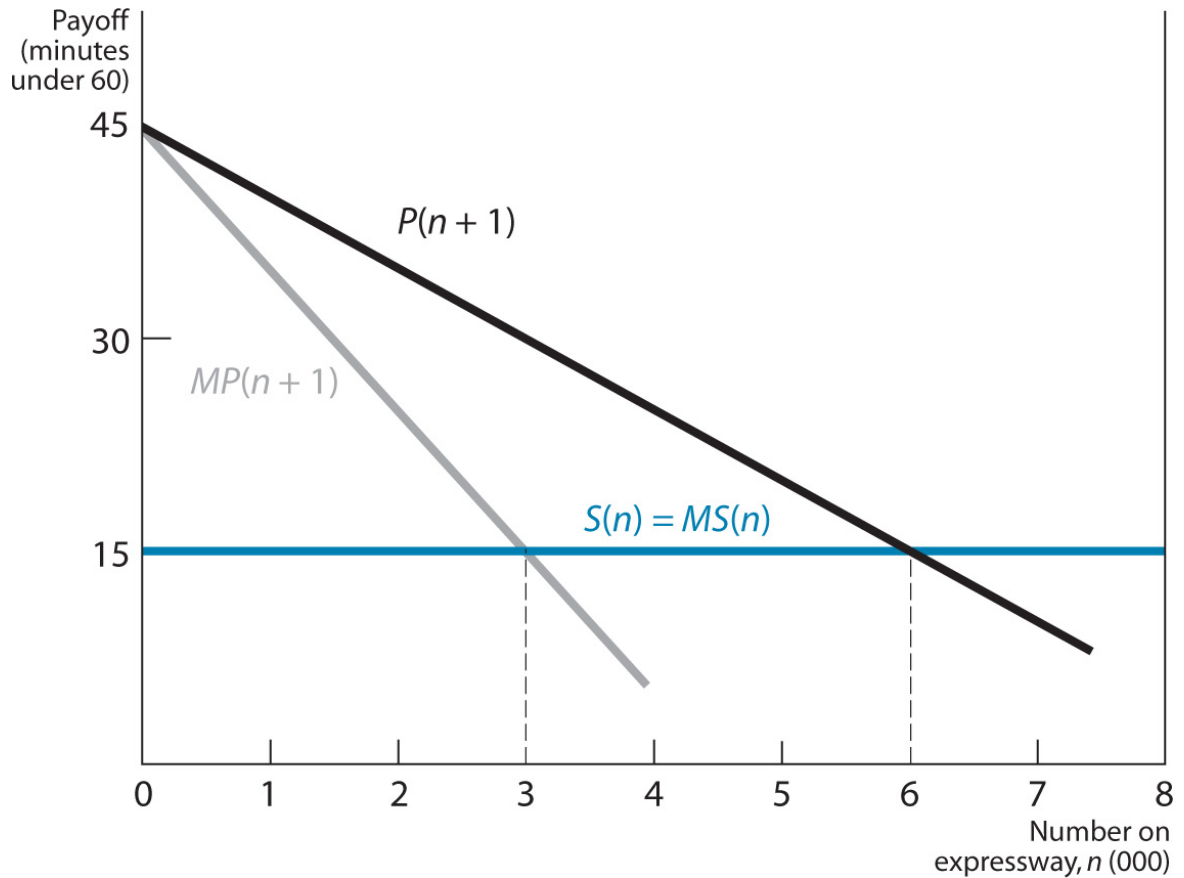


FIGURE 11.10 Equilibrium and Optimum in Commuting Game

We can use equation (11.1) to calculate the precise conditions under which a switch will be beneficial for a particular person versus for society as a whole. Recall that if n people are already using the expressway and another driver is contemplating switching from the local roads to the expressway, she stands to gain from this switch if $P(n+1) > S(n)$, whereas the total social payoff increases if $T(n+1) - T(n) > 0$. The private gain is positive if

$$45 - (n+1) \times 0.005 > 15$$

$$44.995 - 0.005n > 15$$

$$n < 200 (44.995 - 15) = 5,999,$$

whereas the condition for the social gain to be positive is

$$45 - (n + 1) \times 0.005 - 15 - 0.005n > 0$$

$$29.995 - 0.01n > 0$$

$$n < 2,999.5.$$

Thus, if given the free choice, commuters will crowd onto the expressway until there are almost 6,000 of them, but all crowding beyond 3,000 reduces the total social payoff. Society as a whole would be best off if the number of commuters on the expressway were kept down to 3,000.

We show this result graphically in Figure 11.10; this figure replicates Figure 11.9, with the addition of marginal private and social gain curves. The two curves indicating $P(n + 1)$ and $S(n)$ meet at $n = 5,999$; that is, at the value of n for which $P(n + 1) = S(n)$ or for which the marginal private gain is just zero. Everywhere to the left of this intersection point, any one driver on the local roads calculates that she will get a positive gain by switching to the expressway. As some drivers make this switch, the numbers on the expressway increase—the value of n in society rises, as was the case in our example in [Section 3.A](#). Conversely, to the right of the intersection point (that is, for $n > 5,999$), $S(n) > P(n + 1)$; so each of the $(n + 1)$ drivers on the expressway stands to gain by switching to the local road. As some do so, the numbers on the expressway decrease, and n falls. From the left of the intersection point, this process converges to $n = 5,999$, and from the right, it converges to $n = 6,000$.

If we had used the calculus approach, we would have regarded 1 as a very small increment in relation to n and graphed $P(n)$ instead of $P(n + 1)$. Then the intersection point would have been at $n = 6,000$ instead of at $n = 5,999$. As you can see, it makes very little difference in practice. What this means is that we can call $n = 6,000$ the Nash equilibrium of the route choice game when choices are governed by purely individual considerations. Given a free choice, 6,000 of the 8,000 commuters will choose the expressway, and only 2,000 will drive on the local roads.

But we can also interpret the outcome in this game from the perspective of the whole society of commuters. Society gains from

an increase in the number of commuters, n , on the expressway when $T(n + 1) - T(n) > 0$ and loses from an increase in n when $T(n + 1) - T(n) < 0$. To figure out how to show this on the graph, we express the idea somewhat differently. We rearrange equation (11.1) into two pieces, one depending only on P and the other depending only on S :

$$\begin{aligned} T(n + 1) - T(n) &= (n + 1) P(n + 1) + [N - (n + 1)] S(n + 1) - \\ &\quad nP(n) - [N - n] S(n) \\ &= \{P(n + 1) + n[P(n + 1) - P(n)]\} - \{S(n) + [N - (n + 1)][S(n + 1) - S(n)]\}. \end{aligned}$$

The expression in the first set of braces is the effect on the payoffs of the set of commuters who choose P ; this expression includes the $P(n + 1)$ of the switcher and the spillover effect, $n[P(n + 1) - P(n)]$, on all the other n commuters who choose P . We call this expression the marginal social payoff for the P -choosing subgroup when their number increases from n to $n + 1$, or $MP(n + 1)$ for short. Similarly, the expression in the second set of braces is the marginal social payoff for the S -choosing subgroup, or $MS(n)$ for short. Then, the full expression for $T(n + 1) - T(n)$ tells us that the total social payoff increases when one person switches from S to P (or decreases if the switch is from P to S) if $MP(n + 1) > MS(n)$. The total social payoff decreases when one person switches from S to P (or increases when the switch is from P to S) if $MP(n + 1) < MS(n)$.

Using our expressions for $P(n + 1)$ and $S(n)$ in the commuting example, we have

$$MP(n + 1) = 45 - (n + 1) \times 0.005 + n \times (-0.005) = 44.995 - 0.01n$$

while $MS(n) = 15$ for all values of n . Figure 11.10 illustrates $MP(n + 1)$ and $MS(n)$. Note that $MS(n)$ coincides with $S(n)$ everywhere because the local roads are never congested. But the $MP(n + 1)$ curve lies below the $P(n + 1)$ curve. Because of the negative spillover effect, the social gain from one person's

switching to the expressway is less than the private gain to the switcher.

The $MP(n + 1)$ and $MS(n)$ curves meet at $n = 2,999$, or approximately 3,000. To the left of this intersection, $MP(n + 1) > MS(n)$, and society stands to gain by allowing one more person on the expressway. To the right, the opposite is true, and society stands to gain by shifting one person from the expressway to the local roads. Thus, the socially optimal allocation of drivers is 3,000 on the expressway and 3,000 on the local roads.

If you wish to use calculus, you can write the total payoff for the expressway drivers as $nP(n) = n(45 - 0.005n) = 45n - 0.005n^2$. Then $MP(n + 1)$ is the derivative of this with respect to n —namely, $45 - 0.005 \times 2n = 45 - 0.01n$. The rest of the analysis can proceed as before.

How might this society achieve this optimum allocation of its drivers? Different cultures and political groups approach such problems in different ways, each with its own merits and drawbacks. The society could simply restrict access to the expressway to 3,000 drivers. But how would it choose those 3,000? It could adopt a first-come, first-served rule, but then drivers would race one another to get there early and waste a lot of time. A bureaucratic society could set up criteria based on complex calculations of needs and merits as defined by civil servants; then everyone would undertake some costly activities to meet these criteria. In a politicized society, the important “swing voters” or organized pressure groups or contributors might be favored. In a corrupt society, those who bribe the officials or the politicians might get the preference. A more egalitarian society could allocate the rights to drive on the expressway by lottery or rotate them from one month to the next. A scheme that lets you drive only on certain days, depending on the last digit of your car’s license plate, is an example. But such a scheme is not so egalitarian as it seems, because the rich can have two cars and choose license-plate numbers that will allow them to drive every day.

Many economists prefer a more open system of charges. Suppose each driver on the expressway is made to pay a tax t , measured in units of time. Then the private benefit from using the expressway becomes $P(n) - t$, and the number n in the Nash equilibrium will be determined by $P(n) - t = S(n)$. [Here, we are ignoring the tiny difference between $P(n)$ and $P(n + 1)$, which is possible when N is very large.] We know that the socially optimal value of n is 3,000. Using the expressions $P(n) = 45 - 0.005n$ and $S(n) = 15$, and plugging in 3,000 for n , we find that $P(n) - t = S(n)$ —that is, drivers are indifferent between the expressway and the local roads—when $45 - 15 - t = 15$, or $t = 15$. If we value time at a minimum wage of about \$10 an hour, 15 minutes comes to \$2.50. This is the tax or toll that, when charged, will keep the numbers on the expressway down to what is socially optimal.

Note that when 3,000 drivers are on the expressway, the addition of one more increases the time spent by each of them by 0.005 minutes, for a total of 15 minutes. This is exactly the tax that each driver is being asked to pay. In other words, each driver is made to pay the cost of the negative spillover effect that she imposes on the rest of society. This tax brings home to each driver the extra cost of her action and therefore induces her to take the socially optimal action; economists say that she is being made to [internalize the externality](#). This idea, that people whose actions hurt others should pay for the harm that they cause, adds to the appeal of this approach. But the proceeds from the tax are not used to compensate the others directly. If they were, then each expressway user would count on receiving from others just what she pays, and the whole purpose would be defeated. Instead, the proceeds of the tax go into general government revenues, where they may or may not be used in a socially beneficial manner.

Those economists who prefer to rely on markets argue that if the expressway has a private owner, his profit motive will induce him to charge just enough for its use to reduce the number of users to the socially optimal level. An owner knows that if he charges a tax t for each user, the number of users n will be determined by $P(n) - t = S(n)$. His revenue will be $t n = n[P(n) - S(n)]$, and he will act in such a way as to maximize this revenue. In our

example, the revenue is $n[45 - 0.005n - 15] = n[30 - 0.005n] = 30n - 0.005n^2$. It is easy to see that this revenue is maximized when $n = 3,000$. But in this case, the revenue goes into the owner's pocket; most people regard that as a bad solution.

D. Positive Spillover Effects

Many aspects of positive spillover effects or positive externalities can be understood simply as mirror images of those of negative spillover effects. A person's marginal private gain from undertaking activities with positive spillover effects is less than society's marginal gain from such activities. Therefore, such actions will be underused and their benefits underprovided in the Nash equilibrium. A better outcome can be achieved by augmenting people's incentives; providing those persons whose actions create positive spillover effects with a reward just equal to the spillover benefit will achieve the social optimum.

Indeed, the distinction between positive and negative spillover effects is to some extent a matter of semantics. Whether a spillover effect is positive or negative depends on which choice you call P and which you call S . In the commuting example, suppose we called choosing the local roads P and choosing the expressway S . Then one commuter's switch from S to P will reduce the time taken by all the others who choose S , so this action will have a positive spillover effect on them. In another example, consider vaccination against some infectious disease. Each person getting vaccinated reduces his own risk of catching the disease (marginal private gain) and reduces the risk of others' getting the disease through him (spillover effect). If being unvaccinated is called the S action, then getting vaccinated has a positive spillover effect. If remaining unvaccinated is called the P action, then the act of remaining unvaccinated has a negative spillover effect. These two ways of labeling actions here suggest two sorts of policies that might be used to bring individual actions into conformity with the social optimum: Society can either reward those who get vaccinated or penalize those who fail to do so.

But actions with positive spillover effects can have one very important new feature that distinguishes them from actions with negative spillover effects—namely, [positive feedback](#). Suppose

the spillover effect of your choosing P is to increase the payoff to the others who are also choosing P . Then your choice increases the attractiveness of that action (P) and may induce some others to take it also, setting in train a process that culminates in everyone's taking that action. Conversely, if very few people are choosing P , then it may be so unattractive that they, too, give it up, leading to a situation in which everyone chooses S . In other words, positive feedback can give rise to multiple Nash equilibria, which we now illustrate by using a very real example.

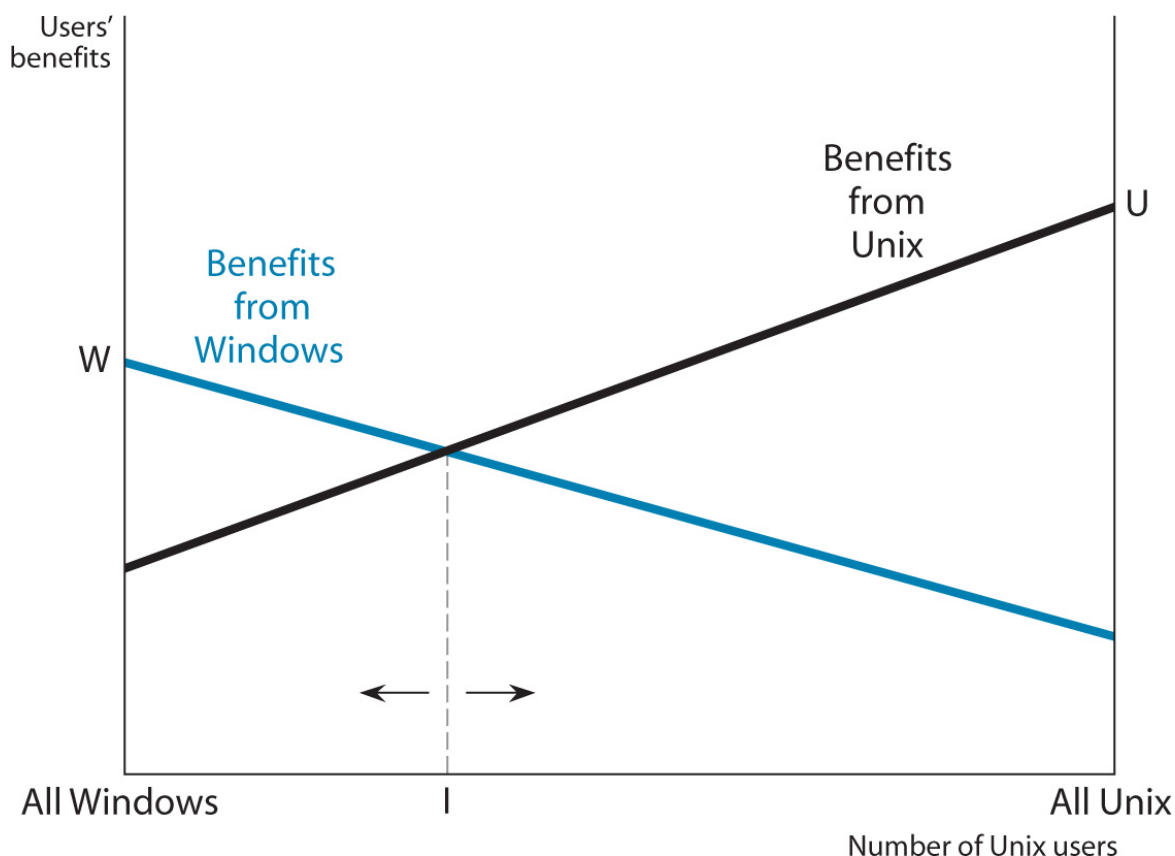


FIGURE 11.11 Payoffs in Operating System - Choice Game

When you buy a computer, you have to choose between one with a Windows operating system and one with an operating system based on Unix, such as Linux. The higher the number of Unix users rises, the better it will be to purchase such a computer. The system will have fewer bugs because more users will have detected those that exist, more application software will be available, and more experts will be available to help with any problems that

arise. Similarly, a Windows-based computer will be more attractive the more Windows users there are. In addition, many computing aficionados would argue that the Unix system is superior. Without necessarily taking a position on that matter, we show what will happen if that is the case. Will individual choice lead to the socially best outcome?

A graph similar to Figures 11.6 through 11.8 can be used to show the payoffs to an individual computer purchaser of the two strategies, Unix and Windows. As shown in Figure 11.11, the Unix payoff rises as the number of Unix users rises, and the Windows payoff rises as the number of Unix owners falls (and the number of Windows users rises). As already explained, the graph is drawn assuming that the payoff to Unix users when everyone in the population is a Unix user (at the point labeled U) is higher than the payoff to Windows users when everyone in the population is a Windows user (at W).

If the current population has only a small number of Unix users, then the situation is represented by a point to the left of the intersection of the two payoff curves at I, and each individual user finds it better to choose Windows. When there is a larger number of Unix users in the population, placing the society to the right of I, it is better for each user to choose Unix. Thus, a mixed population of Unix and Windows users is sustainable as an equilibrium only when the current population has exactly I Unix users; only then will no member of the population have any incentive to switch operating systems. And even that situation is unstable. Suppose just one person accidentally makes a different decision. If he switches to Windows, his choice will push the population to the left of I, in which case others will have an incentive to switch to Windows, too. If he switches to Unix, the population moves to the right of I, creating an incentive for more people to switch to Unix. The cumulative effect of these switches will eventually push the society to an all-Unix or an all-Windows outcome; these are the two stable equilibria of the game.⁴

But which of the two stable equilibria will be achieved in this game? The answer depends on where the game starts. If you look at

the configuration of today's computer users, you will see a heavily Windows-oriented population. Thus, it seems that because there are so few Unix users (or so many Windows users), the world is moving toward the all-Windows equilibrium. Schools, businesses, and private users have become [locked in](#) to this particular equilibrium as a result of an accident of history. If it is indeed true that Unix provides more benefits to society when used by everyone, then the all-Unix equilibrium should be preferred over the all-Windows one that we are approaching. Unfortunately, although society as a whole might be better off with the change, no individual computer user has an incentive to make a change. Only coordinated action can swing the pendulum toward Unix. A critical mass of individual users, more than I in Figure 11.11, must use Unix before it becomes individually rational for others to choose the same operating system.

There are many examples of similar choices of convention being made by different groups of people. The most famous cases are those in which it has been argued, in retrospect, that a wrong choice was made. Advocates claim that electric cars could have been developed for greater efficiency than gasoline; it certainly would have been cleaner. Proponents of the Dvorak typewriter/computer keyboard configuration claim that it would be better than the QWERTY keyboard if used everywhere. Many engineers agree that Betamax had more going for it than VHS in the video recorder market. In such cases, the whims of the public or the genius of advertisers help determine the ultimate equilibrium and may lead to a "bad" or "wrong" outcome from society's perspective. Other situations do not suffer from such difficulties. Few people concern themselves with fighting for a reconfiguration of traffic-light colors, for example.⁵

The ideas of positive feedback and lock-in find an important application in macroeconomics. Production is more profitable the higher the level of demand in the economy, which happens when national income is higher. In turn, income is higher when firms are producing more and are therefore hiring more workers. This positive feedback creates the possibility of multiple equilibria, of which the high-production, high-income one is better for society, but individual decisions may lock the economy into the

low-production, low-income equilibrium. The better equilibrium could be turned into a focal point by public declaration—such as “The only thing we have to fear is fear itself” —but the government can also inject demand into the economy to the extent necessary to move it to the better equilibrium. In other words, the possibility of unemployment due to a deficiency of aggregate demand—as discussed in the supply-and-demand language of economic theory by the British economist John Maynard Keynes in his well-known 1936 book titled *Employment, Interest, and Money*—can be seen from a game-theoretic perspective as the result of a failure to solve a collective-action problem.^{[6](#)}

Endnotes

- The term *positive feedback* may create the impression that it is a good thing, but in technical language “positive” merely connotes “reinforcing” and includes no general value judgment about the outcome. In this example, the same positive feedback mechanism could lead to either an all-Unix outcome or an all-Windows outcome; one outcome could be worse than the other. [Return to reference 4](#)
- Not everyone agrees that the Dvorak keyboard and the Betamax video recorder were clearly superior alternatives. See two articles by S. J. Liebowitz and Stephen E. Margolis, “Network Externality: An Uncommon Tragedy,” *Journal of Economic Perspectives*, vol. 8 (Spring 1994), pp. 146 – 49, and “The Fable of the Keys,” *Journal of Law and Economics*, vol. 33 (April 1990), pp. 1 – 25. [Return to reference 5](#)
- John Maynard Keynes, *Employment, Interest, and Money* (London: Macmillan, 1936). See also John Bryant, “A Simple Rational-Expectations Keynes-type Model,” *Quarterly Journal of Economics*, vol. 98 (1983), pp. 525 – 28, and Russell Cooper and Andrew John, “Coordination Failures in a Keynesian Model,” *Quarterly Journal of Economics*, vol. 103 (1988), pp. 441 – 63, for formal game-theoretic models of unemployment equilibria. [Return to reference 6](#)

Glossary

marginal private gain

The change in an individual's own payoff as a result of a small change in a continuous-strategy variable that is at his disposal.

spillover effect

Same as external effect.

external effect

When one person's action alters the payoff of another person or persons. The effect or spillover is *positive* if one's action raises others' payoffs (for example, network effects) and *negative* if it lowers others' payoffs (for example, pollution or congestion). Also called **externality** or spillover effect.

externality

Same as external effect.

marginal social gain

The change in the aggregate social payoff as a result of a small change in a continuous-strategy variable chosen by one player.

internalize the externality

To offer an individual a reward for the external benefits he conveys on the rest of society, or to inflict a penalty for the external costs he imposes on the rest, so as to bring his private incentives in line with social optimality.

positive feedback

When one person's action increases the payoff of another person or persons taking the same action, thus increasing their incentive to take that action too.

locked in

A situation where the players persist in a Nash equilibrium that is worse for everyone than another Nash equilibrium.

4 A BRIEF HISTORY OF IDEAS

A. The Classics

The problem of collective action has been recognized by social philosophers and economists for a very long time. The seventeenth-century British philosopher Thomas Hobbes argued that society would break down in a “war of all against all” unless it was ruled by a dictatorial monarch, or *Leviathan* (the title of his book). One hundred years later, the French philosopher Jean-Jacques Rousseau described the problem of a prisoners’ dilemma in his *Discourse on Inequality*. A stag hunt needs the cooperation of the whole group of hunters to encircle and kill the stag, but any individual hunter who sees a hare may find it better for himself to leave the circle to chase the hare. But Rousseau thought that such problems were the product of civilization and that people in the natural state lived harmoniously as “noble savages.” At about the same time, two Scots pointed out some dramatic solutions to such problems. David Hume, in his *Treatise on Human Nature*, argued that the expectations of future returns of favors can sustain cooperation. Adam Smith’s *Wealth of Nations* developed a grand vision of an economy in which the production of goods and services motivated purely by private profit could result in an outcome that was best for society as a whole.⁷

The optimistic interpretation persisted, especially among many economists and even several political scientists, to the point where it was automatically assumed that if an outcome was beneficial to a group as a whole, the actions of its members would bring that outcome about. This belief received a necessary rude shock in the mid-1960s, when Mancur Olson published *The Logic of Collective Action*. He pointed out that the best collective outcome would not prevail unless it was in each individual’s private interest to perform the assigned action—that is, unless it was a Nash equilibrium.

However, Olson did not specify the collective-action game very precisely. Although it looked like a prisoners' dilemma, Olson insisted that it was not necessarily so, and we have already seen that the problem can also take the form of a game of chicken or an assurance game.⁸

Another major class of collective-action problems—namely, those concerning the depletion of common-access resources—received attention at about the same time. If a resource such as a fishery or a meadow is open to all, each user will exploit it as much as he can, because any self-restraint on his part will merely make more available for the others to exploit. As we mentioned earlier, Garrett Hardin wrote a well-known article on this subject titled “The Tragedy of the Commons.” Common-resource problems are unlike our irrigation project game, in which each person has a strong private incentive to free ride on the efforts of others. In regard to a common resource, each person has a strong private incentive to exploit it to the full, making everyone else pay the cost that results from the degradation of the resource.

B. Modern Approaches and Solutions

Until recently, many social scientists and most physical scientists took a Hobbesian line on the common-resource problem, arguing that it can be solved only by a government that forces everyone to behave cooperatively. Others, especially economists, retained their Smithian optimism. They argued that placing the resource in proper private ownership, where its benefits can be captured in the form of profit by the owner, will induce the owner to restrain its use in a socially optimal manner. He will realize that the value of the resource (fish or grass, for example) may be higher in the future because less will be available, and therefore he can make more profit by saving some of it for that future.

Nowadays, thinkers on all sides have begun to recognize that collective-action problems come in diverse forms and that there is no uniquely best solution to all of them. They also understand that groups or societies do not stand helpless in the face of such problems, but devise various ways to cope with them. Much of this work has been informed by game-theoretic analysis of repeated prisoners' dilemmas and similar games.⁹

Solutions to collective-action problems of all types must induce individuals to act cooperatively or in a manner that would be best for the group, even though their interests may best be served by doing something else—in particular, taking advantage of the others' cooperative behavior.¹⁰ Humans exhibit much in the way of cooperative behavior. The act of reciprocating gifts and the skill of detecting cheating, for example, are so common in all societies and throughout history that there is reason to argue that they may be instincts.¹¹ But human societies generally rely heavily on purposive social and cultural customs, norms, and sanctions

in inducing cooperative behavior from their individual members. These methods are conscious, deliberate attempts to design the game in order to solve the collective-action problem. [12](#)

We approach the matter of solution methods from the perspective of the type of game being played.

A solution is easiest if the collective-action problem takes the form of an assurance game. Then it is in every person's private interest to take the socially best action if he expects all other persons to do likewise. In other words, the socially optimal outcome is a Nash equilibrium. The only problem is that the same game has other, socially worse, Nash equilibria. Then all that is needed to achieve the best Nash equilibrium, and thereby the social optimum, is to make it a focal point—that is, to ensure the convergence of the players' expectations on it. Such a convergence can result from a social [custom](#) or [convention](#)—namely, a mode of behavior that finds automatic acceptance because it is in everyone's interest to follow it so long as others are expected to do likewise. For example, if all the farmers, herders, weavers, and other producers in an area want to get together to trade their wares, all they need is the assurance of finding others with whom to trade. Then the custom that the market is held in village X on day Y of every week makes it optimal for everyone to be there on that day. [13](#)

One complication remains. For the desired outcome to be a focal point, each person must have confidence that all others understand it, which in turn requires that those others have confidence that all others understand it . . . In other words, the point must be common knowledge. Usually, some prior social action is necessary to ensure that this is true. Publication in a medium that is known by everyone to be sufficiently widely read, and discussion in an inward-facing circle so that everyone knows that everyone else was present

and paying attention, are some methods used for this purpose. [14](#)

Our analysis in [Section 2](#) suggested that individual payoffs are often configured in such a way that collective-action problems, particularly of large groups, take the form of a prisoners' dilemma. Not surprisingly, the methods for coping with such problems have received the most attention.

The simplest solution method attempts to change people's preferences so that the game is no longer a prisoners' dilemma. If individuals get sufficient pleasure from cooperating, or suffer enough guilt or shame when they cheat, they will cooperate to maximize their own payoffs. If the extra payoff from cooperation is conditional—one gets pleasure from cooperating or guilt or shame from cheating if, but only if, many others are cooperating—then the game can turn into an assurance game. In one of its equilibria, everyone cooperates because everyone else does, and in the other, no one cooperates because no one else does. Then the collective-action problem is the simpler one of making the better equilibrium the focal point. If the extra payoff from cooperation is unconditional—one gets pleasure from cooperating or guilt or shame from cheating regardless of what the others do—then the game can have a unique equilibrium where everyone cooperates. In many situations, it is not even necessary for everyone to have such payoffs. If a substantial proportion of the population does, that may suffice for the desired collective outcome.

Some such prosocial preferences may be innate, hardwired by a biological evolutionary process. But they are more likely to be social or cultural products. Most societies make deliberate efforts to instill prosocial thinking in children during the process of socialization in families and schools. The growth of prosocial preferences is seen in experiments using ultimatum and dictator games of the kind we discussed

in [Chapter 3](#). When these experiments are conducted on children of different ages, very young children behave selfishly. By age eight, however, they develop a significant sense of equality. True prosocial preferences develop gradually thereafter, with some relapses, into an adult fair-mindedness. Thus, a long process of education and experience instills *internalized norms* into people's preferences. [15](#)

However, people do differ in the extent to which they internalize prosocial preferences, and the process may not go far enough to solve many collective-action problems. Most people have sufficiently broad understanding of what the socially cooperative action is in most situations, but individuals retain the personal temptation to cheat. Therefore, a system of external sanctions or punishments is needed to sustain cooperative action. We call these widely understood, but not automatically followed, rules of behavior *enforced norms*.

In [Chapter 10](#), we described in detail several methods for achieving a cooperative outcome in prisoners' dilemma games, including repetition, penalties (or rewards), and leadership. In that discussion, we were mainly concerned with two-person dilemmas. The same methods apply to enforcement of norms in collective-action problems in large groups, with some important modifications or innovations.

We saw in [Chapter 10](#) that repetition was the most prominent of these methods, so we focus the most attention on it here. Repetition can achieve cooperative outcomes as equilibria of individual actions in a repeated two-person prisoners' dilemma by holding up the prospect that cheating will lead to a breakdown of cooperation. More generally, what is needed to maintain cooperation is the expectation in the mind of each player that his personal benefits from cheating will be transitory and that they will quickly be replaced by a payoff lower than that associated with cooperative behavior. If

players are to believe that cheating is not beneficial from a long-term perspective, cheating should be detected quickly, and the punishment that follows (reduction in future payoffs) should be sufficiently swift, sure, and painful.

A group has one advantage in this respect over a pair of individual persons. The same pair may not have occasion to interact all that frequently, but each of them is likely to interact with *someone* in the group all the time. Therefore, B's temptation to cheat A can be countered by his fear that others, such as C, D, and so on, whom he meets in the future will punish him for this action. An extreme case where bilateral interactions are not repeated and punishment must be inflicted on one's behalf by a third party is, in Yogi Berra's well-known saying, "Always go to other people's funerals. Otherwise they won't go to yours."

But a group has some offsetting disadvantages over direct bilateral interaction when it comes to sustaining good behavior in repeated interactions. The required speed and certainty of detection and punishment of cheating suffer as the numbers in the group increase. One sees many instances of successful cooperation in small village communities that would be unimaginable in a large city or state.

Start with the detection of cheating, which is never easy. In most real situations, payoffs are not completely determined by the players' actions, but are subject to some random fluctuations. Even with two players, if one gets a low payoff, he cannot be sure that the other cheated; it may have been that his low payoff resulted from some random event. With more people, an additional question enters the picture: If someone cheated, who was it? Punishing someone without being sure of his guilt beyond a reasonable doubt is not only morally repulsive, but also counterproductive. The incentive to cooperate gets blunted if even cooperative actions are susceptible to punishment by mistake.

Next, with multiple players, even when cheating is detected and the cheater identified, this information has to be conveyed sufficiently quickly and accurately to others. For this, the group must be small, or else it must have a good communication or gossip network. Also, members should not have much reason to accuse others falsely.

Finally, even after cheating is detected and the information spread to the whole group, the cheater's punishment—enforcement of the social norm—has to be arranged. A third person often has to incur some personal cost to inflict such punishment. For example, if C is called on to punish B, who had previously cheated A, C may have to forgo some profitable business that he could have transacted with B. Then the inflicting of punishment is itself a collective-action game and suffers from the same temptation to shirk—that is, not to participate in the punishment. A society could construct a second-round system of punishments for shirking, but that in turn may be yet another collective-action problem! However, humans seem to have evolved an instinct whereby people get some personal pleasure from punishing cheaters even when they have not themselves been the victims of a particular act of cheating.¹⁶ Interestingly, the notion that “one should impose sanctions, even at personal cost, on violators of enforced social norms” seems itself to have become an internalized norm.¹⁷

Norms are reinforced by observation of society's general adherence to them, and they lose their force if they are frequently seen to be violated. Before the advent of the welfare state, when those who fell on hard economic times had to rely on help from family or friends or their immediate small social group, the work ethic constituted a norm that held in check the temptation to slacken one's own efforts and become a free rider on the support of others. As government took over the supporting role and unemployment compensation or welfare became an entitlement, this norm of

the work ethic weakened. After the sharp increases in unemployment in Europe in the late 1980s and early 1990s, a significant fraction of the population became users of the official support system, and the norm weakened even further. [18](#)

Different societies or cultural groups may develop different conventions and norms to achieve the same purpose. At the trivial level, each culture has its own set of good manners—ways of greeting strangers, indicating approval of food, and so on. When two people from different cultures meet, misunderstandings can arise. More importantly, each company or office has its own ways of getting things done. The differences among these customs and norms are subtle and difficult to pin down, but many mergers fail because of a clash of these “corporate cultures.”

Next, consider the chicken form of collective-action games. Here, the nature of the remedy depends on whether the largest total social payoff is attained when everyone participates (what we called chicken version I in [Section 1.B](#)) or when some cooperate and others are allowed to shirk (chicken II). For chicken I, where everyone has the individual temptation to shirk, the problem is much like that of sustaining cooperation in the prisoners’ dilemma, and all the earlier remarks on that game apply here, too. Chicken II is different—easier in one respect and harder in another. Once an assignment of roles between participants and shirkers is made, no one has the private incentive to switch: If the other driver is assigned the role of going straight, then you are better off swerving, and the other way around. Therefore, if a custom creates the expectation of an equilibrium, it can be maintained without further social intervention such as sanctions. However, in this equilibrium, the shirkers get higher payoffs than the participants do, and this inequality can create its own problems for the game; the conflicts and tensions, if they are major, can threaten the whole fabric of

the society. If the interaction is repeated, the inequality problem can be solved by rotating the roles of participants and shirkers to equalize payoffs over time.

Sometimes the problem of differential payoffs in version II of the prisoners' dilemma or chicken is solved not by restoring equality, but by [oppression](#) or [coercion](#), which forces a dominated subset of society to accept the lower payoff and allows the dominant subgroup to enjoy the higher payoff. In many societies throughout history, the work of handling animal carcasses was forced on particular groups or castes in this way. And the history of the maltreatment of racial and ethnic minorities and of women provides vivid examples of such practices. Once such a system becomes established, no one member of the oppressed group can do anything to change the situation. The oppressed must get together as a group and act to change the whole system, though how to do so is itself another problem of collective action.

Finally, consider the role of leadership in solving collective-action problems. In [Chapter 10](#), we pointed out that if the players are of very unequal "size," the prisoners' dilemma may disappear because it may be in the private interest of the larger player to continue cooperation and to accept the cheating of the smaller player. Here we recognize the possibility of a different kind of bigness—namely, having a "big heart." People in most groups differ in their preferences, and many groups have one or a few who take genuine pleasure in expending personal effort to benefit the whole. If there are enough such people for the task at hand, then the collective-action problem disappears. Most schools, churches, local hospitals, and other worthy causes rely on the work of such willing volunteers. This solution, like others before it, is more likely to work in small groups, where the fruits of their actions are more closely

and immediately visible to the benefactors, who are therefore encouraged to continue.

C. Applications

In her book *Governing the Commons*, Elinor Ostrom describes several examples of resolution of common-resource problems at local levels. Most of them require taking advantage of features specific to the context in order to set up systems of detection and punishment. A fishing community on the Turkish coast, for example, assigns and rotates locations to its members; the person who is assigned a good location on any given day will naturally observe and report any intruder who tries to usurp his place. Many other users of common resources, including the grazing commons in medieval England, actually restricted access and controlled overexploitation by allocating complex, tacit, but well-understood rights to individual persons. In one sense, this solution bypasses the common-resource problem by dividing up the resource into a number of privately owned subunits.

The most striking feature of Ostrom's range of cases is their immense variety. Some of the prisoners' dilemmas involving the exploitation of common-property resources that she examined were solved by private initiative by the group of people actually in the dilemma; others were solved by external public or governmental intervention. In some instances, the dilemma was not resolved at all, and the group remained trapped in the all-shirk outcome. Despite this variety, Ostrom identifies several common features that make prisoners' dilemmas of collective action easier to solve: (1) it is essential to have an identifiable and stable group of potential participants; (2) the benefits of cooperation have to be large enough to make it worth paying all the costs of monitoring and enforcing the rules of cooperation; and (3) it is very important that the members of the group can communicate with one another. This last feature accomplishes several things. First, it makes the norms clear—everyone

knows what behavior is expected, what kind of cheating will not be tolerated, and what sanctions will be imposed on cheaters. Next, it spreads information about the efficacy of the cheating-detection mechanism, thereby building trust and removing the suspicion that each participant might hold that he is abiding by the rules while others are getting away with breaking them. Finally, it enables the group to monitor the effectiveness of the existing arrangements and to improve on them as necessary. All these requirements look remarkably like those identified in [Chapter 10](#) from our theoretical analysis of the prisoners' dilemma and from the observations of Axelrod's tournaments.

Ostrom's study of the fishing village also illustrates what can be done if the social optimum requires different persons to do different things, in which case some get higher payoffs than others. In a repeated relationship, the advantageous position can rotate among the participants, thereby maintaining some sense of equality over time.

Ostrom finds that an external enforcer of cooperation may not be able to detect cheating or impose punishment with sufficient clarity and swiftness. Thus, the frequent call for centralized or government policy to solve collective-action problems is often proved wrong. Another example comes from village communities or "communes" in late-nineteenth-century Russia. These communities solved many collective-action problems of irrigation, crop rotation, management of woods and pastures, and road and bridge construction and repair in just this way. "The village . . . was not the haven of communal harmony. . . . It was simply that the individual interests of the peasants were often best served by collective activity." Reforms by early-twentieth-century czarist governments and by Soviet revolutionaries of the 1920s alike failed, partly because the old system had such a hold on the peasants' minds that they resisted anything new, but also because the reformers failed to understand the role

that some of the prevailing practices played in solving collective-action problems and thus failed to replace them with equally effective alternatives.^{[19](#)}

The difference between small and large groups is well illustrated by Avner Greif's comparison of two groups of traders in countries around the Mediterranean Sea in medieval times. The Maghribis were Jewish traders who relied on extended family and social ties. If one member of this group cheated another, the victim informed all the others by writing letters. When guilt was convincingly proved, no one in the group would deal with the cheater. This system worked well on a small scale of trade. But as trade expanded around the Mediterranean, the group could not find sufficiently close or reliable insiders to go to the countries with the new trading opportunities.

In contrast, the Genoese traders established a more official legal system. A contract had to be registered with the central authorities in Genoa. The victim of any cheating or violation of the contract had to take a complaint to the authorities, who carried out the investigation and imposed the appropriate fines on the cheater. This system, with all its difficulties of detection, could be more easily expanded with the expansion of trade.^{[20](#)} As economies have grown and world trade has expanded, we have seen a similar shift from tightly linked groups to more arm's-length trading relationships, and from enforcement based on repeated interactions to official law.

The idea that smaller groups are more successful at solving collective-action problems, which forms the major theme of Olson's *Logic of Collective Action* (see footnote 8), has led to an insight important in political science. In a democracy, all voters have equal political rights, and the majority's preference should prevail. But we see many instances in which this does not happen. The effects of policies are generally

good for some groups and bad for others. To get its preferred policy adopted, a group has to take political action—lobbying, publicity, campaign contributions, and so on. To do these things, the group must solve a collective-action problem, because each member of the group may hope to shirk and enjoy the benefits that the others' efforts have secured. If small groups are better able to solve this problem, then the policies resulting from the political process will reflect *their* preferences, even if other groups who fail to organize are more numerous and suffer losses greater than the successful groups' gains.

The most dramatic example of policies reflecting the preferences of an organized group comes from the arena of trade policy. A country's restrictions on imports help domestic producers whose goods compete with these imports, but they hurt the consumers of the imported goods and the domestic competing goods alike, because prices for these goods are higher than they would be otherwise. The domestic producers are few in number, and the consumers are almost the whole population of the country; thus, the total dollar amount of the consumers' losses is typically far bigger than the total dollar amount of the producers' gains. Political considerations based on constituency membership numbers and economic considerations of dollar gains and losses alike would lead us to expect a consumer victory in this policy arena; we would expect to see at least a push for the idea that import restrictions should be abolished, but we don't. The smaller and more tightly knit associations of producers are better able to organize for political action than are the numerous, dispersed consumers.

More than 70 years ago, the American political scientist E. E. Schattschneider provided the first extensive documentation and discussion of how pressure politics drives trade policy. He recognized that "the capacity of a group for organization has a great influence on its activity," but he did not

develop any systematic theory of what determines this capacity.²¹ The analysis of Olson and others has improved our understanding of the issue, but the triumph of pressure politics over economics persists in trade policy to this day. For example, in the late 1980s, the U.S. sugar trade policy that severely limited the amount of imported sugar cost each of the 240 million people in the United States about \$11.50 per year, for a total of about \$2.75 billion per year, while it increased the annual incomes of about 10,000 sugar-beet farmers by about \$50,000 each, and those of 1,000 sugarcane farmers by as much as \$500,000 each, for a total of about \$1 billion. The net loss to the U.S. economy was \$1.75 billion.²² Each of the unorganized consumers continues to bear his small share of the cost in silence; many of them are not even aware that each is paying \$11.50 a year too much for his sweet tooth.

If this overview of the theory and practice of solving collective-action problems seems diverse and lacking a neat summary statement, that is because the problems are equally diverse, and the solutions depend on the specifics of each problem. The one general lesson that we can provide is the importance of letting the participants themselves devise solutions by using their local knowledge of the situation, their advantage of proximity in monitoring the cooperative or shirking actions of others in the community, and their ability to impose sanctions on shirkers by exploiting various ongoing relationships within the social group.

Finally, a word of caution. You might be tempted to come away from this discussion of collective-action problems with the impression that individual freedom always leads to harmful outcomes that can and must be improved by social norms and sanctions. Remember, however, that societies face problems other than those of collective action, some of which are better solved by individual initiative than by joint efforts. Societies can often get hidebound and autocratic, becoming

trapped in their norms and customs and stifling the innovation that is so often the key to economic growth. Collective action can become collective inaction.[23](#)

Endnotes

- The great old books cited in this paragraph have been reprinted many times in many different versions. For each, we list the year of original publication and the details of one relatively easily accessible reprint. In each case, the editor of the reprinted version provides an introduction that conveniently summarizes the main ideas. Thomas Hobbes, *Leviathan; or the Matter, Form, and Power of Commonwealth Ecclesiastical and Civil*, 1651 (Everyman Edition, London: J. M. Dent, 1973); David Hume, *A Treatise of Human Nature*, 1739 (Oxford: Clarendon Press, 1976); Jean-Jacques Rousseau, *A Discourse on Inequality*, 1755 (New York: Penguin Books, 1984); Adam Smith, *An Inquiry into the Nature and Causes of the Wealth of Nations*, 1776 (Oxford: Clarendon Press, 1976). [Return to reference 7](#)
- Mancur Olson, *The Logic of Collective Action* (Cambridge, Mass.: Harvard University Press, 1965). [Return to reference 8](#)
- Prominent in this literature are Michael Taylor, *The Possibility of Cooperation* (New York: Cambridge University Press, 1987); Elinor Ostrom, *Governing the Commons* (New York: Cambridge University Press, 1990); and Matt Ridley, *The Origins of Virtue* (New York: Viking Penguin, 1996). [Return to reference 9](#)
- The problem of the need to attain cooperation and its solutions are not unique to human societies. Examples of cooperative behavior in the animal kingdom have been explained by biologists in terms of the advantage of the gene and of the evolution of instincts. For more, see Chapter 12 and Ridley, *Origins of Virtue*. [Return to reference 10](#)
- See Ridley, *Origins of Virtue*, Chapters 6 and 7. [Return to reference 11](#)

- The social sciences do not have precise and widely accepted definitions of terms such as *custom* and *norm*; nor are the distinctions among such terms always clear and unambiguous. We set out some definitions in this section, but be aware that you may find different usage in other books. Our approach is similar to those found in Richard Posner and Eric Rasmusen, “Creating and Enforcing Norms, with Special Reference to Sanctions,” *International Review of Law and Economics*, vol. 19, no. 3 (September 1999), pp. 369 – 82, and in David Kreps, “Intrinsic Motivation and Extrinsic Incentives,” *American Economic Review*, Papers and Proceedings, vol. 87, no. 2 (May 1997), pp. 359 – 64; Kreps uses the term *norm* for all the concepts that we classify under different names.

Sociologists have a taxonomy of norms that is different from that of economists. It is based on the importance of the norms (those pertaining to trivial matters such as table manners are called *folkways*, and those pertaining to weightier matters are called *mores*), and on whether the norms are formally codified as *laws*. They also maintain a distinction between *values* and norms, recognizing that some norms may run counter to persons’ values and therefore require sanctions to enforce them. This distinction corresponds to ours between customs, internalized norms, and enforced norms. The conflict between individual values and social goals arises for enforced norms, but not for customs or *conventions*, as we label them, or for internalized norms. See Donald Light and Suzanne Keller, *Sociology*, 4th ed. (New York: Knopf, 1987), pp. 57 – 60.

[Return to reference 12](#)

- In his study of the emergence of cooperation, *Cheating Monkeys and Citizen Bees* (New York: Free Press, 1999), the evolutionary biologist Lee Dugatkin labels this case

“selfish teamwork.” He argues that such behavior is likelier to arise in times of crisis, because each person is pivotal at those times. In a crisis, the outcome of the group interaction is likely to be disastrous for everyone if even one person fails to contribute to the group’s effort to get out of the dire situation. Thus, each person is willing to contribute so long as the others do. We will mention Dugatkin’s full classification of alternative approaches to cooperation in Chapter 12 on evolutionary games. [Return to reference 13](#)

- See Michael Chwe, *Rational Ritual: Culture, Coordination, and Common Knowledge* (Princeton, N.J.: Princeton University Press, 2001), for a discussion of this issue and numerous examples and applications of it. [Return to reference 14](#)
- Colin Camerer, *Behavioral Game Theory* (Princeton, N.J.: Princeton University Press, 2003), pp. 65–67. See also pp. 63–75 for an account of differences in prosocial behavior along different dimensions of demographic characteristics and across different cultures. [Return to reference 15](#)
- For evidence of such altruistic punishment instinct, see Ernst Fehr and Simon Gächter, “Altruistic Punishment in Humans,” *Nature*, vol. 415 (January 10, 2002), pp. 137–40. [Return to reference 16](#)
- Our distinction between internalized norms and enforced norms is similar to Kreps’s distinction between functions (iii) and (iv) of norms (Kreps, “Intrinsic Motivation and Extrinsic Incentives,” p. 359). Society can also reward desirable actions just as it can punish undesirable ones. Again, the rewards, financial or otherwise, can be given externally, or players’ payoffs can be changed so that they take pleasure in doing the right thing. The two types of rewards can interact; for example, the peerages and knighthoods given to British philanthropists and others who do good deeds for British

society are external rewards, but individual persons value them only because respect for knights and peers is a British social norm. [Return to reference 17](#)

- Assar Lindbeck, “Incentives and Social Norms in Household Behavior,” *American Economic Review*, Papers and Proceedings, vol. 87, no. 2 (May 1997), pp. 370 – 77. [Return to reference 18](#)
- Orlando Figes, *A People’s Tragedy: The Russian Revolution 1891 – 1924* (New York: Viking Penguin, 1997), pp. 89 – 90, 240 – 41, 729 – 30. See also Ostrom, *Governing the Commons*, p. 23, for other instances where external, government-enforced attempts to solve common-resource problems actually made them worse. [Return to reference 19](#)
- Avner Greif, “Cultural Beliefs and the Organization of Society: A Historical and Theoretical Reflection on Collectivist and Individualist Societies,” *Journal of Political Economy*, vol. 102, no. 5 (October 1994), pp. 912 – 50. [Return to reference 20](#)
- E. E. Schattschneider, *Politics, Pressures, and the Tariff* (New York: Prentice-Hall, 1935); see especially pp. 285 – 86. [Return to reference 21](#)
- Stephen V. Marks, “A Reassessment of the Empirical Evidence on the U.S. Sugar Program,” in *The Economics and Politics of World Sugar Policies*, ed. Stephen V. Marks and Keith E. Maskus (Ann Arbor: University of Michigan Press, 1993), pp. 79 – 108. [Return to reference 22](#)
- David Landes, *The Wealth and Poverty of Nations* (New York: W. W. Norton, 1998), Chapters 3 and 4, makes a spirited case for this effect. [Return to reference 23](#)

Glossary

norm

A pattern of behavior that is established in society by a process of education or culture, to the point that a person who behaves differently experiences a negative psychic payoff.

sanction

Punishment approved by society and inflicted by others on a member who violates an accepted pattern of behavior.

custom

Same as **convention**.

convention

A mode of behavior that finds automatic acceptance as a focal point, because it is in each individual's interest to follow it when others are expected to follow it too (so the game is of the assurance type). Also called **custom**.

oppression

In this context, same as **coercion**.

coercion

In this context, forcing a player to accept a lower payoff in an asymmetric equilibrium in a collective action game, while other favored players are enjoying higher payoffs. Also called **oppression** in this context.

5 “HELP!” : A GAME OF CHICKEN WITH MIXED STRATEGIES

In our discussion of the chicken variant of collective-action games in earlier sections of this chapter, we looked only at the pure-strategy equilibria. But we know from [Chapter 7](#) that such games have mixed-strategy equilibria, too. In collective-action problems, where each participant is thinking, “It is better if I wait for enough others to participate so that I can shirk; but then again, maybe they won’t, in which case I should participate,” mixed strategies nicely capture the spirit of such vacillation. Our last story is a dramatic, even chilling, application of such a mixed-strategy equilibrium.

In 1964 in New York City (in Kew Gardens, Queens), a woman named Kitty Genovese was killed in a brutal attack that lasted more than half an hour. She screamed through it all, and although her screams were heard by many people, and at least three actually witnessed some part of the attack, no one went to help her, or even called the police.

The story created a sensation, and commentators found several ready theories to explain it. The press and most of the public saw this episode as a confirmation of their belief that New Yorkers—or big-city dwellers, or Americans, or people more generally—were apathetic or didn’t care about their fellow human beings.

However, even a little introspection or observation will convince you that people do care about the well-being of other humans, even strangers. Social scientists offered a different explanation for what happened, which they labeled [pluralistic ignorance](#). The idea behind this explanation is that no one can be sure about what is happening or whether help is really needed and how much. People look to one another for clues or guidance about these matters and try to interpret other people’s behavior in this light. If they see that no one else is doing anything to help,

they interpret it as meaning that help is probably not needed, and so they don't do anything either. This explanation has some intuitive appeal, but is unsatisfactory in the Kitty Genovese context. There is a very strong presumption that a screaming woman needs help. What did the onlookers think—that a movie was being shot in their obscure neighborhood? If so, where were the lights, the cameras, the director, other crew?

A better explanation would recognize that although each onlooker might experience strong personal distress from Kitty's suffering and might get genuine personal pleasure if she were saved, each must balance that against the cost of getting involved. You might have to identify yourself if you call the police; you might then have to appear in court as a witness, and so on. Thus, we see that each person may prefer to wait for someone else to call and hope to get for himself the free rider's benefit of the pleasure of a successful rescue.

Social psychologists have a slightly different version of this idea of free riding, which they label [diffusion of responsibility](#). In this version, the idea is that everyone might agree that help is needed, but they are not in direct communication with one another and so cannot coordinate on who should help. Each person may believe that help is someone else's responsibility. And the larger the group, the more likely it is that each person will think that someone else will probably help, and that therefore he can save himself the trouble and the cost of getting involved.

Social psychologists conducted some experiments to test this hypothesis. They staged situations in which someone needed help of different kinds in different places and with different-sized crowds present. Among other things, they found that the larger the size of the crowd, the less likely help was to come forth.

The concept of diffusion of responsibility seems to explain this finding, but not completely. It claims that the larger the crowd, the less likely any one person is to help. But there are more people in a larger crowd, and only one person is needed to act and call the police to secure help. To make it less likely that

even one person helps, the chance of any one person helping has to decrease sufficiently fast with the increase in the total number of potential helpers to offset that increase. To find out whether it does so requires game-theoretic analysis, which we now supply.²⁴

We consider only the aspect of diffusion of responsibility in which action is not consciously coordinated, and we leave aside all other complications of information and inference. Thus, we assume that everyone believes that taking action is necessary and is worth the cost.

Suppose N people are in the group. The action brings each of them a benefit B . Only one person is needed to take the action; more are redundant. Anyone who acts bears the cost C . We assume that $B > C$; so it is worth any one person's while to act even if no one else is acting. Thus, the action is justified in a very strong sense.

The problem is that anyone who takes the action gets the benefit B and pays the cost C for a net payoff of $(B - C)$, whereas he would get the higher payoff B if someone else took the action. Thus, each person has the temptation to let someone else take the action and to become a free rider on another's effort. When all N people are thinking thus, what will be the equilibrium outcome?

If $N = 1$, the single person has a simple decision problem rather than a game. He gets $B - C > 0$ if he takes the action and 0 if he does not. Therefore, he goes ahead and helps.

If $N > 1$, we have a game of strategic interaction with several equilibria. Let us begin by ruling out some possibilities. With $N > 1$, there cannot be a pure-strategy Nash equilibrium in which all people act, because then any one of them would do better by switching to free riding. Likewise, there cannot be a pure-strategy Nash equilibrium in which no one acts, because *given that no one else is acting* (remember that in Nash equilibrium, each player takes the others' strategies as given), it pays any one person to act.

There *are* Nash equilibria in which exactly one person acts; in fact, there are N such equilibria, one corresponding to each member of the group. But when everyone is making the decision individually in isolation, there is no way to coordinate and designate who is to act. Even if members of the group were to attempt such coordination, they might try to negotiate over the responsibility and not reach a conclusion, at least not in time to be of help. Therefore, it is of interest to examine equilibria in which all members have identical strategies.

We already saw that there cannot be an equilibrium in which all N people follow the same pure strategy. Therefore, we should see whether there can be an equilibrium in which they all follow the same mixed strategy. Actually, mixed strategies are quite appealing in this context. The people are isolated, and each is trying to guess what the others will do. Each is thinking, Perhaps I should call the police . . . but maybe someone else will . . . but what if they don't . . . ? Each breaks off this process at some point and does the last thing that he thought of in this chain, but we have no good way of predicting what that last thing is. A mixed strategy carries the flavor of this idea of a chain of guesswork being broken at a random point.

So suppose P is the probability that any one person will not act. If one particular person is willing to mix strategies, he must be indifferent between the two pure strategies of acting and not acting. Acting gets him $(B - C)$ for sure. Not acting will get him 0 if none of the other $(N - 1)$ people act and B if at least one of them does act. Because the probability that any one person fails to act is P , and because they are deciding independently, the probability that none of the $(N - 1)$ others acts is P^{N-1} , and the probability that at least one does act is $(1 - P^{N-1})$. Therefore, the expected payoff to the one person when he does not act is

$$0 \times P^{N-1} + B(1 - P^{N-1}) = B(1 - P^{N-1}).$$

And that one person is indifferent between acting and not acting when

$$B - C = B(1 - P^{N-1}) \quad \text{or when} \quad P^{N-1} = \frac{C}{B} \quad \text{or} \quad P = \left(\frac{C}{B} \right)^{1/(N-1)}.$$

Note how this indifference condition for *one* selected player determines the probability with which the *other* players mix their strategies.

Having obtained the equilibrium mixture probabilities, we can now see how they change as the group size N changes. Remember that $C/B < 1$. As N increases from 2 to infinity, the power $1/(N - 1)$ decreases from 1 to 0. Then C/B raised to this power—namely, P —increases from C/B to 1. Remember that P is the probability that any one person does not take the action. Therefore, the probability of action by any one person—namely, $(1 - P)$ —falls from $1 - C/B = (B - C)/B$ to 0. [25](#)

In other words, the more people there are, the less likely any one of them is to act. This is intuitively true, and in good conformity with the idea of diffusion of responsibility. But it does not yet give us the conclusion that help is less likely to be forthcoming in a larger group. As we said before, help requires action by only one person. Just because there are more and more people, each of whom is less and less likely to act, we cannot conclude immediately that the probability of *at least one* of them acting gets smaller. More calculation is needed to see whether this is the case.

Because the N persons are randomizing independently in the Nash equilibrium, the probability Q that *not even one* of them helps is

$$Q = P^N = \left(\frac{C}{B} \right)^{N/(N-1)}.$$

As N increases from 2 to infinity, $N/(N - 1)$ decreases from 2 to 1, and then Q increases from $(C/B)^2$ to C/B . Correspondingly, the probability that *at least one* person helps—namely, $(1 - Q)$ —decreases from $1 - (C/B)^2$ to $1 - C/B$.²⁶

So our exact calculation does bear out the hypothesis: The larger the group, the *less* likely help is to be given at all. The probability of help does not, however, reduce to zero even in very large groups; instead, it levels off at a positive value—namely, $(B - C)/B$ —which depends on the benefit and cost of action to each individual.

We see how game-theoretic analysis sharpens the ideas from social psychology with which we started. The diffusion of responsibility theory takes us part of the way—namely, to the conclusion that any one person is less likely to act when he is part of a larger group. But the desired conclusion—that larger groups are less likely to provide help at all—needs further and more precise probability calculation based on the analysis of individual mixing and the resulting interactive (game) equilibrium.

And now we ask, did Kitty Genovese die in vain? Do the theories of pluralistic ignorance, diffusion of responsibility, and free riding still play out in the decreased likelihood of individual action within increasingly large cities? Perhaps not. John Tierney of the *New York Times* has publicly extolled the virtues of “urban cranks,”²⁷ people who encourage the civility of the group through prompt punishment of those who exhibit unacceptable behavior—including litterers, noise polluters, and the generally obnoxious boors of society. Such cranks are essentially enforcers of a cooperative norm for society. And as Tierney surveys the actions of known cranks, he reminds the rest of us that “new cranks must be mobilized! At this very instant, people are wasting time reading while norms are being flouted out on the street. . . . You don’t live alone in this world! Have you enforced a norm today?” In other words, we need social norms, as well as some people who have internalized the norm of enforcing norms.

Endnotes

- For a fuller account of the Kitty Genovese story and for the analysis of such situations from the perspective of social psychology, see John Sabini, *Social Psychology*, 2nd ed. (New York: W. W. Norton, 1995), pp. 39 - 44. Our game-theoretic model is based on Thomas Palfrey and Howard Rosenthal, "Participation and the Provision of Discrete Public Goods," *Journal of Public Economics*, vol. 24 (1984), pp. 171 - 93. Many purported facts of the story have been recently challenged in *Kitty Genovese: The Murder, the Bystanders, and the Crime that Changed America* by Kevin Cook (New York: W. W. Norton, 2014), but the power and impact of the originally reported story on American thinking about urban crime remains, and it is still a good example for game-theoretic analysis. [Return to reference 24](#)
- Consider the case in which $B = 10$ and $C = 8$. Then P equals 0.8 when $n = 2$, rises to 0.998 when $n = 100$, and approaches 1 as N continues to rise. The probability of action by any one person is $1 - P$, which falls from 0.2 to 0 as N rises from 2 toward infinity. [Return to reference 25](#)
- With the same sample values for B (10) and C (8), this result implies that increasing N from 2 to infinity increases the probability that not even one person helps from 0.64 to 0.8. And the probability that at least one person helps falls from 0.36 to 0.2. [Return to reference 26](#)
- John Tierney, "The Boor War: Urban Cranks, Unite—Against All Uncivil Behavior. Eggs Are a Last Resort," *New York Times Magazine*, January 5, 1997. [Return to reference 27](#)

Glossary

pluralistic ignorance

A situation of collective action where no individual knows for sure what action is needed, so everyone takes the cue from other people's actions or inaction, possibly resulting in persistence of wrong choices.

diffusion of responsibility

A situation where action by one or a few members of a large group would suffice to bring about an outcome that all regard as desirable, but each thinks it is someone else's responsibility to take this action.

SUMMARY

Multiplayer games generally concern problems of *collective action*. The general structure of collective-action games may be manifested as a prisoners' dilemma, a game of chicken, or an assurance game. The critical difficulty with such games in any form is that the Nash equilibrium arising from individually rational choices may not be the *socially optimal* outcome—the outcome that maximizes the sum of the payoffs of all the players.

In collective-action games, when a person's action has some effect on the payoffs of all the other players, we say that there are *spillover effects*, or *externalities*. They can be positive or negative and can lead to outcomes driven by private interests that are not socially optimal. When actions create negative spillover effects, they are overused from the perspective of society; when actions create positive spillover effects, they are underused. The additional possibility of *positive feedback* exists when there are positive spillover effects; in such a case, the game may have multiple Nash equilibria.

Problems of collective action have been recognized for many centuries and discussed by scholars from diverse fields. Several early works professed no hope for the situation, but others offered up dramatic solutions. The most recent treatments of the subject acknowledge that collective-action problems arise in diverse areas and that there is no single optimal solution. Social scientific analysis suggests that social *custom*, or *convention*, can lead to cooperative behavior. Other possibilities for solutions come from the creation of *norms* of acceptable behavior. Some of these norms are *internalized* in individuals' payoffs; others must be *enforced* by the use of *sanctions* in response to the

uncooperative behavior. Much of the literature agrees that small groups are more successful at solving collective-action problems than large ones.

In large-group games, *diffusion of responsibility* can lead to behavior in which individual persons wait for others to take action and *free ride* on the benefits of that action. If help is needed, it is less likely to be given at all as the size of the group available to provide it grows.

KEY TERMS

[coercion](#) ([452](#))

[collective action](#) ([420](#))

[convention](#) ([448](#))

[custom](#) ([448](#))

[diffusion of responsibility](#) ([457](#))

[external effect](#) ([436](#))

[externality](#) ([436](#))

[free rider](#) ([423](#))

[internalize the externality](#) ([441](#))

[locked in](#) ([444](#))

[marginal private gain](#) ([436](#))

[marginal social gain](#) ([436](#))

[nonexcludable](#) ([421](#))

[nonrival](#) ([421](#))

[norm](#) ([447](#))

[oppression](#) ([452](#))

[pluralistic ignorance](#) ([456](#))

[positive feedback](#) ([443](#))

[pure public good](#) ([421](#))

[sanction](#) ([447](#))

[social optimum](#) ([423](#))

[spillover effect](#) ([436](#))

Glossary

collective action

A problem of achieving an outcome that is best for society as a whole, when the interests of some or all individuals will lead them to a different outcome as the equilibrium of a noncooperative game.

nonexcludable

Benefits that are available to each individual, regardless of whether he has paid the costs that are necessary to secure the benefits.

nonrival

Benefits whose enjoyment by one person does not detract anything from another person's enjoyment of the same benefits.

pure public good

A good or facility that benefits all members of a group, when these benefits cannot be excluded from a member who has not contributed efforts or money to the provision of the good, and the enjoyment of the benefits by one person does not significantly detract from their simultaneous enjoyment by others.

free rider

A player in a collective-action game who intends to benefit from the positive externality generated by others' efforts without contributing any effort of his own.

social optimum

In a collective-action game where payoffs of different players can be meaningfully added together, the social optimum is achieved when the sum total of the players' payoffs is maximized.

marginal private gain

The change in an individual's own payoff as a result of a small change in a continuous-strategy variable that is

at his disposal.

spillover effect

Same as **external effect**.

external effect

When one person's action alters the payoff of another person or persons. The effect or spillover is *positive* if one's action raises others' payoffs (for example, network effects) and *negative* if it lowers others' payoffs (for example, pollution or congestion). Also called **externality** or **spillover effect**.

externality

Same as **external effect**.

marginal social gain

The change in the aggregate social payoff as a result of a small change in a continuous-strategy variable chosen by one player.

internalize the externality

To offer an individual a reward for the external benefits he conveys on the rest of society, or to inflict a penalty for the external costs he imposes on the rest, so as to bring his private incentives in line with social optimality.

positive feedback

When one person's action increases the payoff of another person or persons taking the same action, thus increasing their incentive to take that action too.

locked in

A situation where the players persist in a Nash equilibrium that is worse for everyone than another Nash equilibrium.

norm

A pattern of behavior that is established in society by a process of education or culture, to the point that a person who behaves differently experiences a negative psychic payoff.

sanction

Punishment approved by society and inflicted by others on a member who violates an accepted pattern of behavior.

custom

Same as **convention**.

convention

A mode of behavior that finds automatic acceptance as a focal point, because it is in each individual's interest to follow it when others are expected to follow it too (so the game is of the assurance type). Also called **custom**.

oppression

In this context, same as **coercion**.

coercion

In this context, forcing a player to accept a lower payoff in an asymmetric equilibrium in a collective action game, while other favored players are enjoying higher payoffs. Also called **oppression** in this context.

pluralistic ignorance

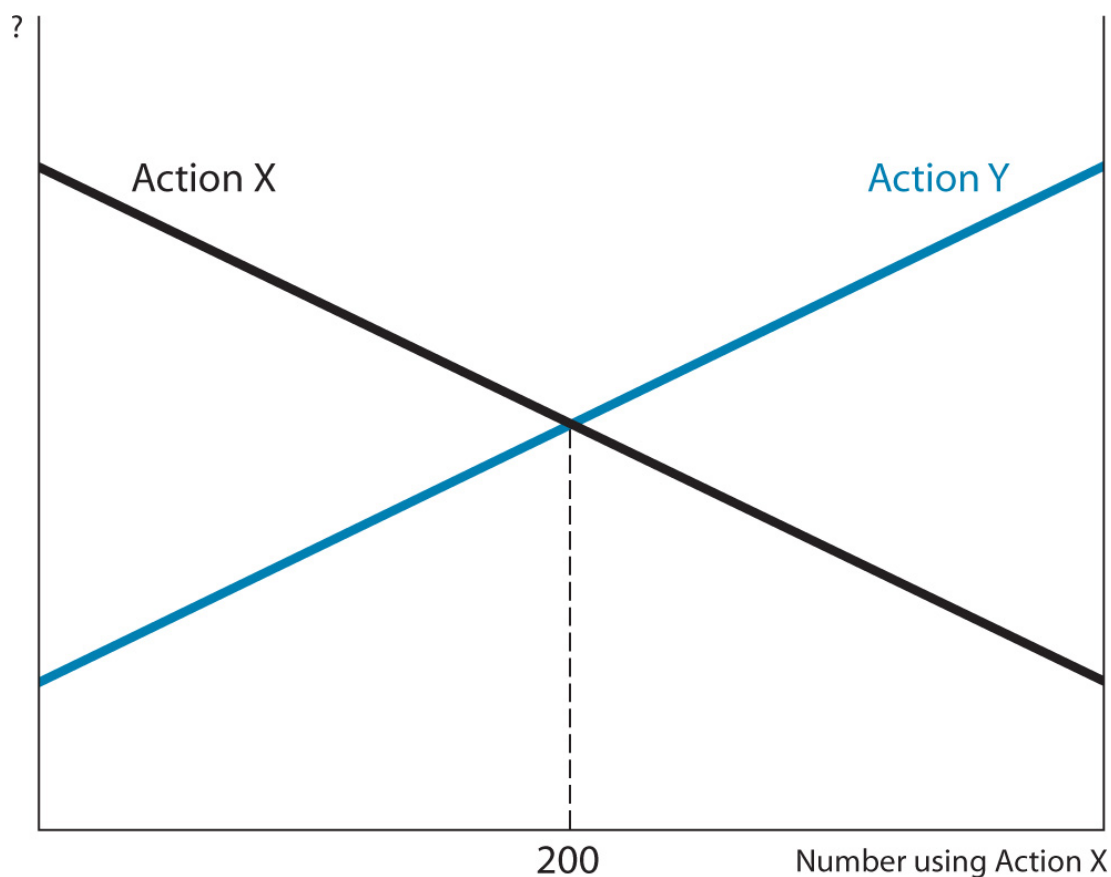
A situation of collective action where no individual knows for sure what action is needed, so everyone takes the cue from other people's actions or inaction, possibly resulting in persistence of wrong choices.

diffusion of responsibility

A situation where action by one or a few members of a large group would suffice to bring about an outcome that all regard as desirable, but each thinks it is someone else's responsibility to take this action.

SOLVED EXERCISES

1. Suppose that 400 people are choosing between Action X and Action Y. The relative payoffs of the two actions depend on how many of the 400 people choose Action X and how many choose Action Y. The payoffs are as shown in the following graph, but the vertical axis is not labeled, so you do not know whether the curves show the benefits or the costs of the two actions.



-
1. You are told that the outcome in which 200 people choose Action X is an *unstable* equilibrium. If 100 people are currently choosing Action X, would you expect the number of people choosing Action X to increase or decrease over time? Why?
 2. For the graph to be consistent with the behavior that you described in part (a), should the curves be labeled as

indicating the *costs* or *benefits* of Action X and Action Y?
Explain your answer.

2. A group has 100 members. Each person can choose to participate or not participate in a common project. If n of them participate in the project, then each participant derives the benefit $p(n) = n$, and each of the $(100 - n)$ shirkers derives the benefit $s(n) = 4 + 3n$.
 1. Is this an example of a prisoners' dilemma, a game of chicken, or an assurance game?
 2. Write the expression for the total social payoff for the group.
 3. Show, either graphically or mathematically, that the maximum total social payoff for the group occurs when $n = 74$.
 4. What difficulties will arise in trying to get exactly 74 participants and allowing the remaining 26 to shirk?
 5. How might the group try to overcome the difficulties identified in part (d)?
3. Consider a small geographic region with a total population of 1 million. There are two towns, Alphaville and Betaville, in which each person can choose to live. For each person, the benefit of living in a town increases for a while with the size of the town (because larger towns have more amenities and so on), but after a point it decreases (because of congestion and so on). If x is the fraction of the population that lives in the same town as you do, your payoff is given by

$$x \text{ if } 0 \leq x \leq 0.4$$

$$0.6 - 0.5x \text{ if } 0.4 < x \leq 1.$$

1. Draw a graph like Figure 11.11, showing the benefits of living in the two towns, as the fraction of the population living in one versus the other varies continuously from 0 to 1.
2. Equilibrium is reached either when both towns are populated and their residents have equal payoffs or when one town—say, Betaville—is totally depopulated, and the residents of the other town (Alphaville) get a higher payoff than would the very first person who seeks to populate Betaville. Use your graph to find all such equilibria.
3. Now consider a dynamic process of adjustment whereby people gradually move toward the town whose residents currently enjoy a larger payoff than do the residents of the other

town. Which of the equilibria identified in part (b) will be stable with these dynamics? Which ones will be unstable?

4. Suppose an amusement park is being built in a city with a population of 100. Voluntary contributions are being solicited to cover the cost. Each citizen is being asked to give \$100. The more people who contribute, the larger the park will be, and the greater the benefit to each citizen. But it is not possible to keep out the noncontributors; they get their share of this benefit anyway. Suppose that when there are n contributors in the population, where n can be any whole number between 0 and 100, the benefit to each citizen in monetary unit equivalents is n^2 dollars.
 1. Suppose that initially no one is contributing. You are the mayor of the city. You would like everyone to contribute, and you can use persuasion on some people. What is the minimum number whom you need to persuade before everyone else will join in voluntarily?
 2. Find the Nash equilibria of the game where each citizen is deciding whether to contribute.
5. In the Italian province of Tuscany, enterprising “hunters” search through the woods each day for the most expensive food in the world, truffles. These delicate subterranean fungi are beloved²⁸ for their strong pungency and are literally worth their weight in gold, at a cost up to \$200 per ounce. The most prized truffles grow only in the wild and must be found by smell (by dog or pig), making them very hard to find (and harder still if many others are also out looking for them). Consider a game played by 99 residents of San Miniato, a Tuscan town famous for its white truffles, deciding each day whether to hunt for truffles or do some other work. If H people go hunting, each will, on average, find $10 - (H/10)$ ounces of truffles, which sell in the San Miniato market for \$200/ounce. (For simplicity, assume that there is unlimited demand at this price.) Those who do not go hunting will do work that earns them \$50 for the day.
 1. Confirm that each resident prefers to go truffle hunting if no one else does, but prefers not to go truffle hunting if everyone else does.
 2. Suppose that the residents decide sequentially whether to go truffle hunting. How many residents go truffle hunting in the resulting rollback equilibrium? How much daily income is generated in the town (adding up the income of all residents)?

3. One day, the mayor of San Miniato decides to limit the number of people who are allowed to go truffle hunting each day. How many people should be allowed to go truffle hunting each day in order to maximize total town income?
4. Residents who are not selected to go truffle hunting under the mayor's plan will have an incentive to try to sneak away and go hunting anyway. How might the townspeople address this challenge?
6. In the early 1950s, Henry O. Bakken struck oil on his farm near Tioga, North Dakota, but it wasn't until the 1990s that geologists appreciated the massive scale of his discovery: an enormous formation underlying parts of Montana, North Dakota, Saskatchewan, and Manitoba and holding an estimated 413 billion barrels of oil. Little of that oil could be recovered prior to recent advances in drilling technology, including hydraulic fracturing (known as "fracking") and horizontal drilling. In 2015, the North Dakota Department of Natural Resources estimated that the break-even price for drilling "the Bakken" was \$40 per barrel, meaning that landowners can earn a profit from drilling their land if oil is selling for more than \$40/barrel, but not otherwise. Because the Bakken formation is so large, drilling activity in the area can affect global oil prices. For this exercise, suppose that the global oil price will equal $$(60 - 25X)$ per barrel, where $0 \leq X \leq 1$ is the fraction of landowners over the Bakken who actively drill their land. Suppose that the number of landowners is extremely large, so that any fraction X is possible. Assume also that each landowner who decides to drill will extract exactly 100,000 barrels of oil. Finally, for simplicity, assume that drilling is a once-and-for-all decision and that all landowners decide simultaneously whether to drill.²⁹
 1. Express each landowner's profit (or loss) when drilling as a function of X . Confirm that each landowner prefers to drill if no one else does (when $X = 0$), but prefers not to drill if everyone else does (when $X = 1$).
 2. Describe the Nash equilibrium of the (simultaneous-move) drilling game. What fraction of landowners choose to drill in the Nash equilibrium?

One day, leaders in the United States and Canada reach an agreement to impose a new "drilling tax" on landowners who drill anywhere over the Bakken. The tax system collects \$10 in tax revenue per barrel of oil extracted and distributes

that money back to all owners of land over the formation (equally) in the form of a “carbon dividend.”

3. Express each landowner's profit (or loss), when drilling and when not drilling, as functions of X , taking into account the drilling tax and the carbon dividend. Confirm that each landowner prefers to drill if no one else does (when $X = 0$), but prefers not to drill if everyone else does (when $X = 1$).
4. Describe the Nash equilibrium of the drilling game after the drilling tax and carbon dividend are introduced. What fraction of landowners choose to drill in this Nash equilibrium?
5. Do landowners over the Bakken earn more or less profit after the drilling tax and carbon dividend are introduced? Explain your answer in terms of how this policy impacts the collective-action problem.
7. Put the Keynesian idea of unemployment described at the end of [Section 3.D](#) into a game with a properly specified set of payoffs, and show the multiple equilibria in a graph. Show the level of production (national product) on the vertical axis as a function of a measure of the level of demand (national income) on the horizontal axis. Equilibrium is reached when national product equals national income—that is, when the function relating the two cuts the 45° line. For what shapes of the function can there be multiple equilibria? Why might you expect such shapes in reality? Suppose that income increases when current production exceeds current income, and that income decreases when current production is less than current income. In this dynamic process, which equilibria are stable and which ones unstable?
8. Write a brief description of a strategic game that you have witnessed or participated in that includes a large number of players and in which individual players' payoffs depend on the number of other players and their actions. Try to illustrate your game with a graph if possible. Discuss the outcome of the actual game in light of the fact that many such games do not have socially optimal outcomes. Do you see evidence of such an outcome in your game?

UNSOLVED EXERCISES

1. Figure 11.5 illustrates the payoffs in a general, two-person collective-action game. There we showed various inequalities in the algebraic payoffs [$p(1)$, etc.] that made the game a prisoners' dilemma. Now you are asked to find similar inequalities corresponding to other kinds of games.
 1. Under what type of payoff structure(s) is the two-person game a chicken game? What further condition(s) on payoffs make the game version I of chicken (as in Figure 11.3)?
 2. Under what type of payoff structure(s) is the two-person game an assurance game?
2. A group has 100 members. Each person can choose to participate or not participate in a common project. If n of them participate in the project, then each participant derives the benefit $p(n) = n$, and each of the $(100 - n)$ shirkers derives the benefit $s(n) = 3n - 180$.
 1. Is this collective-action problem an example of a prisoners' dilemma, a game of chicken, or an assurance game?
 2. Write the expression for the total social payoff of the group.
 3. Show, either graphically or mathematically, that the maximum total social payoff for the group occurs when $n = 100$.
 4. What difficulties will arise in trying to get all 100 group members to participate?
 5. How might the group try to overcome the difficulties identified in part (d)?
3. A class with 30 students enrolled is given a homework assignment with five questions. The first four are the usual kinds of problems, worth a total of 90 points. But the fifth is an interactive game for the class. The question reads: "You can choose whether to answer this question. If you choose to do so, you merely write, 'I hereby answer Question 5.' If you choose not to answer Question 5, your score for the assignment will be based on your performance on the first four questions. If you choose to answer Question 5, then your scoring will be as follows: If fewer than half of the students in the class answer Question 5, you get 10 points for Question 5; that is, 10 points will be added to your score on the other four questions to get your total score for the assignment. If half or more than half of

the students in the class answer Question 5, you get -10 points; that is, 10 points will be subtracted from your score on the other questions.”

1. Draw a graph illustrating the payoffs from the two possible strategies, Answer Question 5 and Don't Answer Question 5, in relation to the number of other students who answer it. Find the Nash equilibrium of the game.
2. What would you expect to see happen in this game if it were actually played in a college classroom? Why? Consider two cases: (i) The students make their choices individually, with no communication. (ii) The students make their choices individually, but can discuss these choices ahead of time in a discussion forum available on the class Web site.
4. There are two routes for driving from A to B. One is a freeway, and the other consists of local roads. The benefit of using the freeway is constant and equal to 1.8, irrespective of the number of people using it. Local roads get congested when too many people use this alternative, but if not enough people use it, the few isolated drivers run the risk of becoming victims of crimes. Suppose that when a fraction x of the population is using the local roads, the benefit of this mode to each driver is given by

$$1 + 9x - 10x^2.$$

1. Draw a graph showing the benefits of the two driving routes as functions of x , regarding x as a continuous variable that can range from 0 to 1.
2. Identify all possible equilibrium traffic patterns from your graph in part (a). Which equilibria are stable? Which ones are unstable? Why?
3. What value of x maximizes the total social payoff to the whole population?
5. Suppose a class of 100 students is comparing two careers—lawyer and engineer. An engineer gets take-home pay of \$100,000 per year, irrespective of the numbers who choose this career. Lawyers make work for one another, so as the total number of lawyers increases, the income of each lawyer increases—up to a point. Ultimately, the competition among them drives down the income of each. Specifically, if there are N lawyers, each will get $100N - N^2$ thousand dollars a year. The annual cost of running a legal practice (office space, secretary, paralegals, access to online reference services, and so forth) is \$800,000. Therefore, each

lawyer takes home $100N - N^2 - 800$ thousand dollars a year when there are N of them.

1. Draw a graph showing the take-home income of each lawyer on the vertical axis and the number of lawyers on the horizontal axis. (Plot a few points—say, for 0, 10, 20, . . . , 90, 100 lawyers. Fit a curve to the points, or use a computer graphics program if you have access to one.)
2. When career choices are made in an uncoordinated way, what are the possible equilibrium outcomes?
3. Now suppose the whole class decides how many should become lawyers, aiming to maximize the total take-home income of the whole class. What will be the number of lawyers? (If you can, use calculus, regarding N as a continuous variable. Otherwise, you can use graphical methods or a spreadsheet.)
6. A group of 12 countries is considering whether to form a monetary union. They differ in their assessments of the costs and benefits of this move, but each stands to gain more from joining, and lose more from staying out, when more of the other countries choose to join. The countries are ranked in order of their liking for joining, 1 having the highest preference for joining and 12 the least. Each country has two actions, IN and OUT. Let

$$B(i, n) = 2.2 + n - i$$

be the payoff to a country with ranking i when it chooses IN and n others have chosen IN. Let

$$S(i, n) = i - n$$

be the payoff to a country with ranking i when it chooses OUT and n others have chosen IN.

1. Show that for country 1, IN is the dominant strategy.
2. Having eliminated OUT for country 1, show that IN becomes the dominant strategy for country 2.
3. Continuing in this way, show that all countries will choose IN.
4. Contrast the payoffs in this outcome with those where all choose OUT. How many countries are made worse off by the formation of the union?

Endnotes

- The French gastronome J.A. Brillat-Savarin famously referred to truffles as “the diamond of the kitchen” in his 1825 treatise on food and cooking. (Brillat-Savarin, one of the world’s first gastronomic essayists, also wrote, “Tell me what you eat, and I will tell you what you are.” Certainly, truffle lovers are a distinct breed.) See Jean Anthelme Brillat-Savarin, *The Physiology of Taste; or, Meditations on Transcendental Gastronomy* (Urbana, IL: Project Gutenberg, 2004), Meditation 6, Section 7 and Aphorism 4. Retrieved April 30, 2019, from www.gutenberg.org/ebooks/5434. [Return to reference 28](#)
- In reality, landowners have a valuable option: to wait and drill later. For more on the value of options, see Avinash Dixit and Robert Pindyck, *Investment under Uncertainty* (Princeton, N.J.: Princeton University Press, 1994). [Return to reference 29](#)

12 ■ Evolutionary Games

WE HAVE SO FAR STUDIED GAMES with many different features—simultaneous and sequential moves, zero-sum and non-zero-sum payoffs, strategic moves to manipulate rules of games to come, one-shot and repeated play, and even games of collective action in which a large number of people play simultaneously. In all of these games, we maintained the ground rules of conventional game theory—namely, that every player in these games has an internally consistent value system, can calculate the consequences of her strategic choices, and makes the choices that best favor her interests. We recognized the possibility that players’ value systems include regard for others, and occasionally—for example, in our discussion of quantal-response equilibrium in [Chapter 5](#)—we allowed that the players recognize the possibility of errors. But we maintained the assumption that each player makes a conscious and calculated choice from her available strategies.

In our presentation of the empirical evidence on strategic choice in several of the earlier chapters, we pointed out several “behavioral” departures from the theory of rational decision making. The most cogent and best-developed theory of such behavior comes from the psychologist and 2002 Nobel laureate in economics Daniel Kahneman.¹ He argues that people have two different systems of decision making: System 1 is instinctive and fast, System 2 is calculating and slow. The fast, instinctive system may be partly hardwired into the brain by evolution, but it is also the result of extensive experience and practice, which build intuition. This system is valuable because it saves much mental effort and time, and it is often the first to be deployed when making a decision. Given enough time and attention, it may be supplemented or supplanted by the more calculating and slower System 2. When the instinctive System 1 is used on any one occasion, the

outcome constitutes an addition to the stock of experience and may lead to a gradual modification of the instinct.

This theory suggests a very different mode of game playing and analysis of games. Players come to a game with the instinctive System 1 and play the strategy it indicates. This strategy may or may not be optimal for the occasion. Its outcome, if good, reinforces the instinct; otherwise, it contributes to gradual change in the instinct. Of course, the outcome depends on what strategies the other player or players deploy, which depends on the state of their instinctive systems, which in turn depends on their experience, and so forth. We need to find out where this process of interactive dynamics of instincts goes. In particular, we need to determine whether it converges to some fixed strategy choices and, if so, whether those choices correspond to what the calculating slow system would have dictated. The biological theory of evolution and evolutionary dynamics offers one approach to this analysis, which we develop in this chapter.

Endnotes

- Daniel Kahneman, *Thinking, Fast and Slow* (New York: Farrar, Straus and Giroux, 2011). [Return to reference 1](#)

1 THE FRAMEWORK

The biological theory of evolution rests on three fundamental concepts: heterogeneity, fitness, and selection. The starting assumption is that a significant part of animal behavior is genetically determined; a complex of one or more genes (a [genotype](#)) governs a particular pattern of behavior, called a behavioral [phenotype](#). Natural diversity in the gene pool ensures heterogeneity of phenotypes in the population. Some behaviors are better suited than others to the prevailing conditions, and the success of a phenotype is given a quantitative measure called its [fitness](#). People are used to thinking of this success in terms of the common but misleading phrase “survival of the fittest” ; however, the ultimate test of biological fitness is not mere survival, but reproductive success. That is what enables an animal to pass on its genes to the next generation and perpetuate its phenotype. The fitter phenotypes then become relatively more numerous in the next generation compared with the less fit phenotypes. This process of [selection](#) is the dynamic that changes the mix of genotypes and phenotypes and may lead eventually to a stable state.

From time to time, chance produces new genetic [mutations](#). Many of these mutations produce behaviors (that is, phenotypes) that are ill suited to the environment, and they die out. But occasionally a mutation leads to a new phenotype that is fitter. Then such a mutant gene can successfully [invade](#) a population—that is, spread to become a significant proportion of the population.

At any time, a population may contain some or all of its biologically conceivable phenotypes. Those that are fitter than others will increase in proportion, some unfit phenotypes may die out, and other phenotypes not currently in

the population may try to invade it. Biologists call a configuration of a population and its current phenotypes evolutionarily stable if the population cannot be invaded successfully by any mutant. This criterion is a static test of evolutionary stability, but often a more dynamic criterion is applied: A configuration is evolutionarily stable if the dynamic process of evolution, starting from any arbitrary mixture of phenotypes in the population, converges to that configuration.²

The fitness of a phenotype depends on the relationship of the individual organism to its environment; for example, the fitness of a particular bird depends on the aerodynamic characteristics of its wings. It also depends on the whole complex of the proportions of different phenotypes that exist in the environment—how aerodynamic the bird's wings are relative to those of the rest of its species. Thus, the fitness of a particular animal—with its particular behavioral traits, such as aggression or sociability—depends on whether other members of its species are predominantly aggressive or passive, crowded or dispersed, and so on. For our purposes, this interaction among phenotypes within a species is the most interesting aspect of the story. Sometimes an individual member of a species interacts with members of another species; then the fitness of a particular type of sheep, for example, may depend on the traits that prevail in the local population of wolves. We consider this type of interaction as well.

All of this evolutionary theory finds a ready parallel in game theory. The behavior of a phenotype can be thought of as a *strategy* of the animal in its interactions with others—for example, to fight or to retreat. The difference is that the choice of strategy is not a purposive calculation, as it would be in standard game theory; rather, it is a genetically predetermined fixture of the phenotype. These interactions lead to *payoffs* to the phenotypes. In biology, the payoffs to

an animal measure its fitness; when we apply these ideas outside of biology, they can have other connotations of success in the social, political, or economic games in question.

The players' payoffs or fitness numbers can be shown in a payoff table just like that for a standard game, with all conceivable phenotypes of one animal arrayed along the rows of the matrix and those of the other along the columns. If more than two animals interact simultaneously—which is called [playing the field](#) in biology—the payoffs can be shown by functions like those for collective-action games described in [Chapter 11](#). We will consider pair-by-pair matches for most of this chapter and will look at many-player interactions briefly in [Section 2.F](#).

Because the population contains a mix of phenotypes, different pairs selected from it will bring to their interactions different combinations of strategies. The actual quantitative measure of the fitness of a phenotype is the average payoff that it gets in all its interactions with others in the population. Those animals with higher fitness will have greater evolutionary success. The eventual outcome of the population dynamics will be an evolutionarily stable configuration of the population.

Biologists have used this approach very successfully. Combinations of aggressive and cooperative behavior, locations of nesting sites, and many more phenomena that elude more conventional explanations can be understood as the stable outcomes of an evolutionary process of selection of fitter strategies. Interestingly, biologists developed the idea of evolutionary games by using the preexisting body of game theory, drawing from its language but modifying the assumption of conscious maximizing to suit their needs. Now game theorists, in turn, are using insights from the research

on biological evolutionary games to enrich their own subject.³

Indeed, the theory of evolutionary games provides a ready-made framework for studying Kahneman's two systems of decision making.⁴ The idea that animals play genetically fixed strategies can be interpreted more broadly in applications of the theory other than in biology. In human interactions, a strategy may be embedded in a player's mind for a variety of reasons—not only by genetics, but also (and probably more importantly) by socialization, cultural background, education, or a rule of thumb based on past experience. All of these can be captured in Kahneman's instinctive, fast System 1. The population can consist of a mixture of different people with different backgrounds or experiences that embed different System 1 strategies in their minds. Thus, for example, some politicians may be motivated to adhere to certain moral or ethical codes even at the cost of electoral success, whereas others are mainly concerned with their own reelection; similarly, some firms may pursue profit alone, whereas others are motivated by social or ecological objectives. We can call each logically conceivable strategy that can be embedded in this way a phenotype for the population of players in the context being studied.

From a population, with its heterogeneity of embedded strategies, pairs of phenotypes are repeatedly randomly selected to interact (play the game) with others of the same or different "species." In each interaction, the payoff of each player depends on the strategies of both; this dependence is governed by the usual rules of the game and can be illustrated in a game table or tree. The *fitness* of a particular strategy is defined as its aggregate or average payoff in its pairings with all the strategies in the population. Some strategies have higher level of fitness than others; in the next generation—that is, the next round of play—these higher-fitness strategies will be used by more

players and will proliferate. Strategies with lower fitness will be used by fewer players and will decay or die out. Occasionally, someone may experiment with or adopt a previously unused strategy from the collection of those that are logically conceivable. This corresponds to the emergence of a mutant.

Although we use the biological analogy, the reason that the fitter strategies proliferate and the less fit ones die out in socioeconomic games differs from the strict genetic mechanism of biology: Players who fared well in the last round will transmit information to their friends and colleagues playing the next round, those who fared poorly in the last round will observe which strategies succeeded better and will try to imitate them, and some purposive thinking and revision of previous rules of thumb will take place between successive rounds. Such “social” and “educational” mechanisms of transmission are far more important in most strategic games than any biological genetics; indeed, they are responsible for reinforcing the reelection orientation of legislators and the profit-maximization motive of firms. Finally, conscious experimentation with new strategies substitutes for the accidental mutation in biological games. Gradual modification in the light of outcomes, experience, observation, and experiment constitute the dynamics of Kahneman’s calculating, slower System 2.

Evolutionarily stable configurations of biological games can be of two kinds. First, a single phenotype may prove fitter than any other, and the population may come to consist of it alone. Such an evolutionarily stable outcome is called [monomorphism](#)—that is, the population contains a single (mono) form (morph). In that case, the unique prevailing strategy is called an [evolutionarily stable strategy \(ESS\)](#). The other possibility is that two or more phenotypes are equally fit (and fitter than some others not played), so they may be able to coexist in certain proportions. Then the

population is said to exhibit polymorphism—that is, a multiplicity (poly) of forms (morph). Such a state will be stable if no new phenotype or feasible mutant can achieve a fitness higher than the fitnesses of the types that are already present in the polymorphic population. Polymorphism comes close to the game-theoretic notion of a mixed strategy. However, there is an important difference: To get polymorphism, no individual player need follow a mixed strategy. Each can follow a pure strategy, but the population exhibits a mixture because different individual players pursue different pure strategies.

The whole setup—the population, its conceivable collection of phenotypes, the payoff matrix in the interactions of the phenotypes, and the rule for the evolution of proportions of the phenotypes in relation to their fitness—constitutes an evolutionary game. An evolutionarily stable configuration of the population can be called an *equilibrium* of the evolutionary game.

Some evolutionary games are symmetric, with the two players on similar footing—for example, two members of the same species competing with each other for food or mates; in a social science interpretation, they could be two elected officials competing for the right to continue in public office. In the payoff table for the game, each can be designated as the row player or the column player with no difference in outcome. Other evolutionary games are asymmetric; such games involve two species, such as a predator and a prey in biology, or a firm and a customer in economics. We develop the analysis of evolutionary games and their stable equilibria, as usual, through a series of illustrative examples.

Endnotes

- The dynamics of phenotypes is driven by an underlying dynamics of genotypes, but, at least at the elementary level, evolutionary biology focuses its analysis at the phenotype level and conceals the genetic aspects of evolution. We will do likewise in our exposition of evolutionary games. Some theories at the genotype level can be found in the materials cited in footnote 3. [Return to reference 2](#)
- Robert Pool, “Putting Game Theory to the Test,” *Science*, vol. 267 (March 17, 1995), pp. 1591 – 93, is a good article for general readers and has many examples from biology. John Maynard Smith deals with such games in biology in his *Evolutionary Genetics*, 2nd ed. (Oxford: Oxford University Press, 1998), [Chapter 7](#), and *Evolution and the Theory of Games* (Cambridge: Cambridge University Press, 1982); the former also gives much background on evolution. Recommended for advanced readers are Peter Hammerstein and Reinhard Selten, “Game Theory and Evolutionary Biology,” in *Handbook of Game Theory*, vol. 2, ed. R. J. Aumann and S. Hart (Amsterdam: North Holland, 1994), pp. 929 – 93; and Jorgen Weibull, *Evolutionary Game Theory* (Cambridge, Mass.: MIT Press, 1995). [Return to reference 3](#)
- Indeed, applications of the evolutionary perspective need not stop with game theory. The following joke offers an “evolutionary theory of gravitation” as an alternative to Newton’ s or Einstein’ s physical theories:

Question: Why does an apple fall from the tree to the earth?

Answer: Originally, apples that came loose from trees went in all directions. But only those that were

genetically predisposed to fall to the earth could reproduce.

[Return to reference 4](#)

Glossary

[evolutionarily stable](#)

A population is evolutionarily stable if it cannot be successfully invaded by a new mutant phenotype.

[evolutionarily stable strategy \(ESS\)](#)

A phenotype or strategy that can persist in a population, in the sense that all the members of a population or species are of that type; the population is evolutionarily stable (static criterion). Or, starting from an arbitrary distribution of phenotypes in the population, the process of selection will converge to this strategy (dynamic criterion).

[fitness](#)

The expected payoff of a phenotype in its games against randomly chosen opponents from the population.

[genotype](#)

A gene or a complex of genes, which give rise to a phenotype and which can breed true from one generation to another. (In social or economic games, the process of breeding can be interpreted in the more general sense of teaching or learning.)

[phenotype](#)

A specific behavior or strategy, determined by one or more genes. (In social or economic games, this can be interpreted more generally as a customary strategy or a rule of thumb.)

[selection](#)

The dynamic process by which the proportion of fitter phenotypes in a population increases from one generation to the next.

[invasion](#)

The appearance of a small proportion of mutants in the population.

[mutation](#)

Emergence of a new genotype.

playing the field

A many-player evolutionary game where all animals in the group are playing simultaneously, instead of being matched in pairs for two-player games.

monomorphism

All members of a given species or population exhibit the same behavior pattern.

polymorphism

An evolutionarily stable equilibrium in which different behavior forms or phenotypes are exhibited by subsets of members of an otherwise identical population.

2 SOME CLASSIC GAMES IN AN EVOLUTIONARY SETTING

In earlier chapters, especially [Chapters 4](#) and [7](#), we introduced and analyzed several games that have become classics of the theory and have been given memorable stories and names—prisoners’ dilemma, chicken, and so on. What happens if we replace the assumption of calculated rational choice in these games by specifying that players come from populations of phenotypes with given strategies, and that the population evolves by selection of the fitter types? Here, we reexamine those games one by one from this evolutionary perspective.

A. Prisoners' Dilemma

Suppose the population is made up of two phenotypes. One type consists of players who are natural-born cooperators; they always work toward the outcome that is jointly best for all players. The other type consists of defectors; they work only for themselves. As an example, we use the restaurant pricing game described in [Chapter 5](#) and presented in a simplified version in [Chapter 10](#). Here, we use the simpler version in which only two pricing choices are available: the jointly best price of \$26 and the Nash equilibrium price of \$20. A cooperator restaurateur would always choose \$26, whereas a defector would always choose \$20. The payoffs (profits) for each phenotype in a single play of this discrete dilemma are shown in Figure 12.1. This figure is the same as Figure 10.2, except that here we call the players simply Row and Column, because each can be any individual restaurateur in the population who is chosen at random to compete against a random rival. Remember that under the evolutionary scenario, no one has the choice between defecting and cooperating; each player is “born” with one trait or the other. Which is the more successful (fitter) trait in the population?

A defecting restaurateur gets a payoff of 288 (\$28,800 a month) if matched against another defecting type, and a payoff of 360 (\$36,000 a month) if matched against a cooperating type. A cooperating type gets 216 (\$21,600 a month) if matched against a defecting type, and 324 (\$32,400 a month) if matched against another cooperating type. No matter what the type of the matched rival, the defecting type does better than the cooperating type. Therefore, the defecting type has a better expected payoff (and is thus fitter) than the cooperating type, irrespective of the proportions of the two types in the population.

		COLUMN	
		\$20 (Defect)	\$26 (Cooperate)
ROW	\$20 (Defect)	288, 288	360, 216
	\$26 (Cooperate)	216, 360	324, 324

FIGURE 12.1 Payoff Table for Restaurant Prisoners' Dilemma (in Hundreds of Dollars per Month)

A little more formally, let x be the proportion of cooperators in the population. Consider any one particular cooperator. In a random pairing, the probability that she will meet another cooperator (and get 324) is x , and the probability that she will meet a defector (and get 216) is $(1 - x)$. Therefore, a typical cooperator's expected payoff is $324x + 216(1 - x)$. For a defector, the probability of meeting a cooperator (and getting 360) is x , and that of meeting another defector (and getting 288) is $(1 - x)$. Therefore, a typical defector's expected profit is $360x + 288(1 - x)$. Now it is immediately apparent that

$360x + 288(1 - x) > 324x + 216(1 - x)$ for all x between 0 and 1.

Therefore, a defector has a higher expected payoff, and is fitter, than a cooperator. This outcome will lead to an increase in the proportion of defectors (a decrease in x) from one "generation" of players to the next, until the population consists entirely of defectors.

Moreover, once a population consists entirely of defectors, it cannot be invaded by mutant cooperators. To see why, consider any very small value of x , meaning that the proportion of cooperators in the population is very small. The cooperators will be less fit than the prevailing defectors, and their proportion in the population will not increase, but will be driven to zero; the mutant strain will die out.

B. Comparing the Evolutionary and Rational-Player Models

The preceding analysis shows that in the prisoners' dilemma, defectors have higher fitness than cooperators, and that an all-defector population cannot be invaded by mutant cooperators. Thus, the evolutionarily stable configuration of the population is monomorphic, consisting of the single strategy or phenotype Defect. We therefore call Defect the *evolutionarily stable strategy* for this population. Of course, Defect is also a strictly dominant strategy when this game is played by rational players. This is not a coincidence. If a game has a strictly dominant strategy, that strategy will also be the (unique) ESS.

More generally, even in games where players do not have a strictly dominant strategy, every ESS must correspond to a Nash equilibrium. To see why, suppose the contrary for the moment. If the use of some strategy (call it S) by all players is not a Nash equilibrium, then some other strategy (call it R) must yield a higher payoff for one player when played against S. A mutant playing R would achieve greater fitness in a population playing S and so would invade successfully. Thus, S cannot be an ESS. In other words, if the use of S by all players is not a Nash equilibrium, then S cannot be an ESS. This is the same as saying that if S is an ESS, it must be a Nash equilibrium for all players to use S.

The evolutionary approach therefore provides a backdoor justification for the rational approach. Even when players are not consciously maximizing, if the more successful strategies get played more often and the less successful ones die out, and if this process converges eventually to a stable strategy, then the outcome must be the same as that resulting from consciously rational play.

Although an ESS must be a Nash equilibrium of the corresponding rational-play game, the converse is not true. There may be multiple Nash equilibria, not all of which are ESS. The concept of an ESS therefore gives us a justification, based on a stability argument, for selecting among multiple Nash equilibria. The examples considered next (chicken, assurance, and penalty kicks) are especially useful

for developing your understanding of and intuition for when and why Nash equilibria fail to be evolutionarily stable.

C. Chicken

Remember our 1950s youths racing their cars toward each other and seeing who will be the first to swerve to avoid a collision? Now suppose the players have no choice in the matter: Each is genetically hardwired to be either a Wimp (who always swerves) or a Macho (who always goes straight). The population consists of a mixture of the two types. Pairs are picked at random every week to play the game. Figure 12.2 shows the payoff table for any two such players—say, A and B. (The numbers replicate those we used in Figure 4.16.)

How will the two types fare? The answer depends on the initial proportions in the population. If the population is almost all Wimps, then a Macho mutant will win and score 1 lots of times, whereas the Wimps will get mostly zeroes because they will mostly meet their own kind. But if the population is mostly Macho, then a Wimp mutant scores -1 , which may look bad but is better than the -2 that all the Machos get. You can think of this outcome appropriately in terms of the biological context and the sexism of the 1950s: In a population of Wimps, a Macho newcomer will show all the rest to be chickens and so will impress all the girls. But if the population consists mostly of Machos, they will be in the hospital most of the time and the girls will have to go for the few Wimps who are healthy.

		B	
		Wimp	Macho
A	Wimp	0, 0	-1, 1
	Macho	1, -1	-2, -2

FIGURE 12.2 Payoff Table for Chicken

In other words, each type is fitter when it is relatively rare in the population. Therefore, each can successfully invade a population consisting of the other type. We should expect to see both types in the population in equilibrium; that is, we should expect an ESS with a mixture of phenotypes, or polymorphism.

To find the proportions of Wimps and Machos in such an ESS, let us calculate the fitness of each type in a general mixed population. Let

x be the fraction of Machos and $(1 - x)$ be the proportion of Wimps. A Wimp meets another Wimp and gets 0 for a fraction $(1 - x)$ of pairings, and meets a Macho and gets -1 for a fraction x of pairings. Therefore, the fitness of a Wimp is $0 \times (1 - x) - 1 \times x = -x$. Similarly, the fitness of a Macho is $1 \times (1 - x) - 2x = 1 - 3x$. The Macho type is fitter if

$$1 - 3x > -x$$

$$2x < 1$$

$$x < \frac{1}{2}.$$

If the population is less than half Macho, then the Machos will be fitter, and their proportion will increase. In contrast, if the population is more than half Macho, then the Wimps will be fitter, and the Macho proportion will fall. Either way, the population proportion of Machos will tend toward $\frac{1}{2}$, and this 50-50 mix will be the stable polymorphic ESS.

Figure 12.3 shows this outcome graphically. Each straight line shows the fitness (the expected payoff in a match against a random member of the population) for one type in relation to the proportion x of Machos.⁵ The fitness of the Wimp type as a function of the proportion of the Machos is $-x$, as we saw two paragraphs ago; it is shown by the gently falling line that starts at 0 where $x = 0$ and goes to -1 where $x = 1$. The corresponding function for the Macho type is $1 - 3x$; its fitness is shown by the rapidly falling line that starts at 1 where $x = 0$ and falls to -2 where $x = 1$. The Macho line lies above the Wimp line for $x < \frac{1}{2}$ and below it for $x > \frac{1}{2}$, showing that the Macho type is fitter when the value of x is small, and the Wimp type is fitter when x is large.

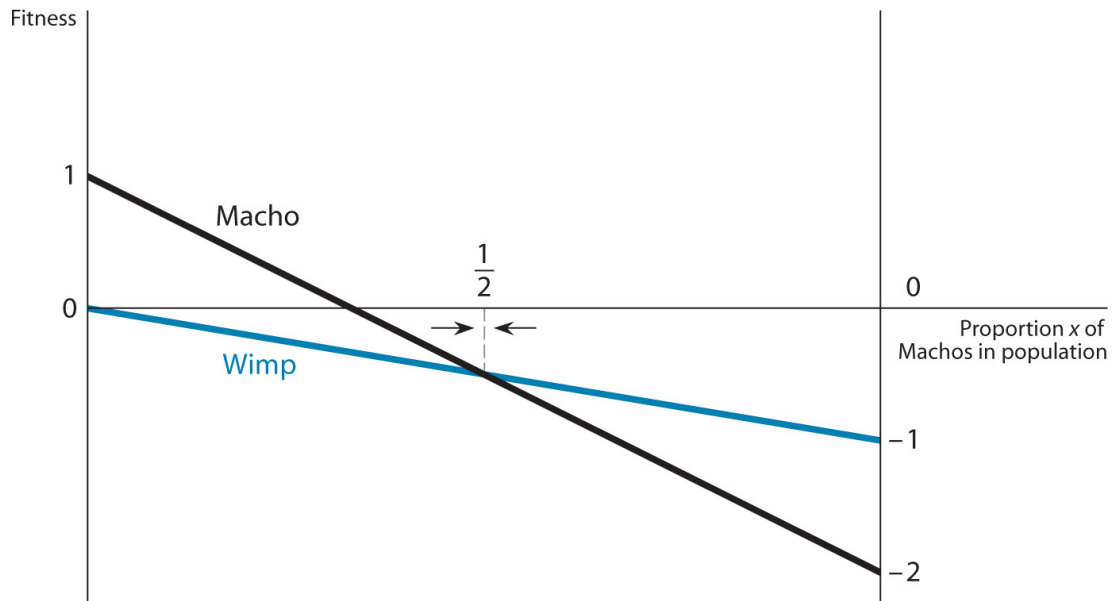


FIGURE 12.3 Fitness Graphs and Polymorphic Equilibrium for Chicken

Now we can compare and contrast the evolutionary theory of this game with our earlier theory of [Chapters 4](#) and [7](#), which was based on the assumption that the players were conscious, rational calculators of strategies. There, we found three Nash equilibria: two in pure strategies, where one player goes straight and the other swerves, and one in mixed strategies, where each player goes straight with a probability of $\frac{1}{2}$ and swerves with a probability of $\frac{1}{2}$.

If the population is truly 100% Macho, then all players are equally fit (or equally unfit, in their hospital beds!). Similarly, in a population of nothing but Wimps, all are equally fit. But these monomorphic configurations are unstable. In an all-Macho population, a Wimp mutant will outscore them and invade successfully.⁶ Once some Wimps get established, no matter how few, our analysis shows that their proportion will rise inexorably toward $\frac{1}{2}$. Similarly, an all-Wimp population is vulnerable to a successful invasion of mutant Machos, and the process again goes to the same polymorphism. Thus, the polymorphic configuration is the only evolutionarily stable outcome.

Most interesting is the connection between the mixed-strategy equilibrium of the rationally played game and the polymorphic equilibrium of the evolutionary game. The mixture proportions in the equilibrium strategy of the former game are *exactly the same* as the

population proportions in the latter game, a 50 - 50 mixture of Wimp and Macho. But the interpretations of these outcomes differ. In the rational framework, each player mixes his own strategies; in the evolutionary framework, every member of the population uses a pure strategy, but different members use different strategies, and so we see a mixture of strategies in the population.⁷

Once again we see a correspondence between Nash equilibria in a rationally played game and stable outcomes in a game with the same payoff structure when played according to the evolutionary rules. We also get better understanding of the mixed-strategy equilibrium, which seemed puzzling when we looked at chicken from the rational perspective. Rational chicken left open the possibility of costly mistakes. Each player went straight one time in two, so one time in four, the players collided. The pure-strategy equilibria avoided the collisions. At that time, this may have led you to think that there was something undesirable about the mixed-strategy equilibrium, and you may have wondered why we were spending time on it. Now you see the reason. That seemingly strange equilibrium emerges as the stable outcome of a natural dynamic process in which each player tries to improve his payoff against the population that he confronts.

D. Assurance

We illustrated assurance games in [Chapter 4](#) with the story of Holmes and Watson deciding where to meet to compare notes at the end of a busy day. We could imagine the same payoff structure applying to a pair of friends deciding where to meet for coffee, at Starbucks or the local diner, for example. In the evolutionary context, we assume that each player in such a coffeehouse-choice game is born liking either Starbucks or the local diner, and that the population includes some of each type. We also assume that pairs of players are chosen at random each day to play the game. We denote the strategies here with S (for Starbucks) and L (for local diner). Figure 12.4 shows the payoff table for a random pairing in this game. The payoffs are the same as those illustrated earlier in Figure 4.14; only the names of the players and the actions have changed.

If we reframe this game as one played by rational strategy-choosing players, we find two equilibria in pure strategies, (S, S) and (L, L), of which the latter is better for both players. If they communicate and coordinate explicitly, they can settle on it quite easily. But if they are making the choices independently, they need to coordinate through a convergence of expectations—that is, by finding a focal point.

The rationally played assurance game has a third equilibrium, in mixed strategies, that we found in [Chapter 7](#). In that equilibrium, each player chooses Starbucks with a probability of $\frac{2}{3}$ and the local diner with a probability of $\frac{1}{3}$; the expected payoff for each player is $\frac{2}{3}$. As we showed in [Chapter 7](#), this payoff is worse than the one associated with the less attractive of the two pure-strategy equilibria, (S, S), because independent mixing leads the players to make clashing or bad choices quite a lot of the time. Here, the bad outcome (a payoff of 0) has a probability of $\frac{4}{9}$: The two players go to different meeting places almost half the time.

FRIEND 2			
		S	L
FRIEND 1	S	1, 1	0, 0
	L	0, 0	2, 2

FIGURE 12.4 Payoff Matrix for the Assurance Game

What happens when this is an evolutionary game? In the population at large, each member is hardwired, either to choose S or to choose L. Randomly chosen pairs of such people are assigned to attempt a meeting. Suppose x is the proportion of S types in the population and $(1 - x)$ is that of L types. Then the fitness of a particular S type individual—her expected payoff in a random encounter of this kind—is $x \times 1 + (1 - x) \times 0 = x$. Similarly, the fitness of each L type is $x \times 0 + (1 - x) \times 2 = 2(1 - x)$. Therefore, the S type is fitter when $x > 2(1 - x)$, or when $x > \frac{2}{3}$. The L type is fitter when $x < \frac{2}{3}$. At the balancing point $x = \frac{2}{3}$, the two types are equally fit.

Once again, as in chicken, the probabilities associated with the mixed-strategy equilibrium that would obtain under rational choice seem to reappear under evolutionary rules as the population proportions in a polymorphic equilibrium. But now this mixed equilibrium is not stable. The slightest chance departure of the proportion x from the balancing point $\frac{2}{3}$ will set in motion a cumulative process that takes the population mix farther away from the balancing point. If x increases from $\frac{2}{3}$, the S type becomes fitter and propagates faster, increasing x even more. If x falls from $\frac{2}{3}$, the L type becomes fitter and propagates faster, lowering x even more. Eventually x will either rise all the way to 1 or fall all the way to 0, depending on which disturbance occurs. The difference relative to the game of chicken is that in chicken, each type is fitter when it is rarer, so the population proportions tend to move away from the extremes and toward a mid-range balancing point. In contrast, in the assurance game, each type is fitter when it is more numerous, because the more of the population that is the same type as you, the lower the risk of failing to meet, so population proportions tend to move toward the extremes.

Figure 12.5 illustrates the fitness graphs and equilibria for the assurance game. This diagram is very similar to Figure 12.3 in that the two lines show the fitnesses of the two types in relation to their proportion in the population, and the intersection of the lines gives the balancing point. The only difference is that, away from the balancing point, the more numerous type is the fitter, whereas in Figure 12.3, it was the less numerous type.

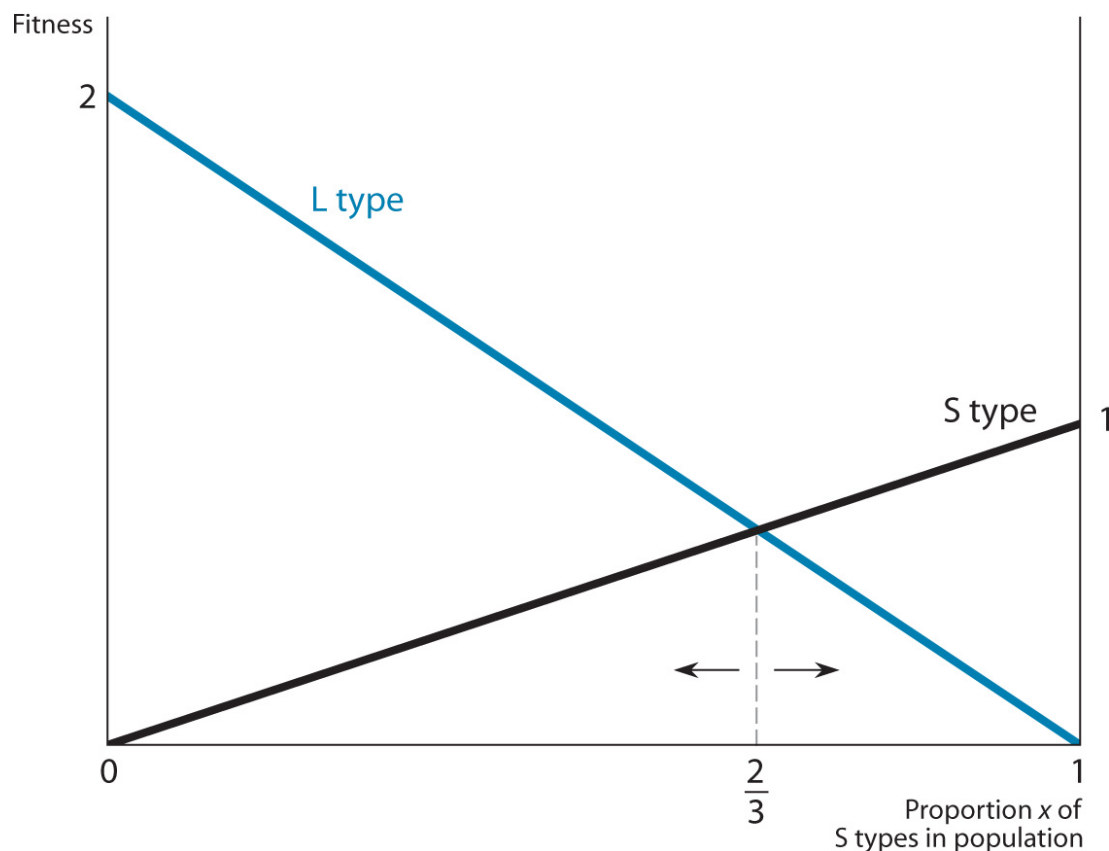


FIGURE 12.5 Fitness Graphs and Equilibria for the Assurance Game

Because each type is less fit when it is rare, the two extreme monomorphic configurations of the population are the only possible evolutionarily stable states. It is easy to check that both outcomes are ESSs according to the static test: An invasion by a small mutant population of the other type will die out because the mutants, being rare, will be less fit. Thus, in assurance or coordination games, in contrast to chicken, the evolutionary process does not preserve the bad equilibrium where there is a positive probability that the players choose clashing strategies. Finally, evolutionary dynamics do not guarantee convergence to the better of the two equilibria when starting from an arbitrary initial mixture of phenotypes—where the population ends up depends on where it starts.

E. Soccer Penalty Kicks

In [Chapter 7](#), we studied mixed-strategy equilibria in the game of penalty kicking in soccer. Some new features of this game emerge when we consider it through the lens of evolutionary games. Most important is its asymmetry: Kickers and goalies have very different roles. In the terminology of this chapter, we may as well regard them as different *species* playing an asymmetric game. We have a population of goalies and a population of kickers, and for each interaction one representative from each population is picked to play. Instead of the goalie making a deliberate choice of whether to go left or right as each kick is taken, and the kicker deciding on each occasion whether to kick to the goalie's left or right, each goalie is either a left-leaner or a right-leaner, and likewise each kicker is either a left-sider or a right-sider. The success of each type of kicker will depend on what type of goalie he faces, and vice versa.

Figure 12.6 is an illustrative matrix of payoffs for such a game. The kicker's payoff is the percentage of times he succeeds in scoring a goal, and the goalie's payoff is the percentage of times he saves the kick (or the kick misses). Remember that for the kicker, left means the goalie's left.

		Goalie	
		Left	Right
Kicker	Left	30, 70	90, 10
	Right	90, 10	50, 50

FIGURE 12.6 Payoff Matrix for the Soccer Penalty-Kick Game

Suppose the kicker population consists of a fraction x of left-siders (and $1 - x$ of right-siders), and the goalie population has a fraction y of left-leaners (and $1 - y$ of right-leaners). Facing a goalie chosen randomly from this population, a left-side kicker will have greater fitness than a right-side kicker if

$30y + 90(1 - y) > 90y + 50(1 - y)$, or $90 - 60y > 40y + 50$, or $y < 0.4$.

If y is small—that is, if the population of goalies is predominantly right-leaning—then left-side kickers are more successful, which agrees with intuition. Conversely, facing a kicker chosen randomly from the kicker population, a left-leaning goalie will have greater fitness than a right-leaning goalie if

$$70x + 10(1 - x) > 10x + 50(1 - x), \text{ or } 60x + 10 > 50 - 40x, \text{ or } x > 0.4.$$

If x is large—that is, if the population of kickers has predominantly left-siders—then a left-leaning goalie will be more successful, as is again quite intuitive.

Experience based on historical success rates will gradually change behaviors of the current populations and newcomers. Will the dynamics converge to a stable configuration of population proportions?

If $x = 0.4$ and $y = 0.4$, both types of each species will be equally fit, and there will be no inducement to change behavior. This suggests $(0.4, 0.4)$ as the obvious candidate for an ESS, and, confirming the correspondence we found earlier in other games, it is also the unique mixed-strategy Nash equilibrium of the conventional, rationally played game.⁸ Let us examine its stability.

If $y < 0.4$, then left-side kickers are more successful, so x will tend to increase; if $y > 0.4$, x will tend to decrease. Conversely, if $x < 0.4$, then right-leaning goalies are more successful, so y will tend to decrease; if $x > 0.4$, y will tend to increase. We show these dynamics graphically in Figure 12.7, using arrows pointing right to show a tendency for x to increase, pointing up to show a tendency for y to increase, and so on.

We see no tendency to converge to $(0.4, 0.4)$. Indeed, if the rates at which x and y change when away from this point have just the right magnitudes, the dynamics can cycle around this point, as shown by the circular loop in Figure 12.7 with arrows indicating the direction of motion. This closed orbit could be small, wound tightly around the ESS, or it could be a bigger loop, depending on the starting point. So an ESS need not always be dynamically stable in the sense that all paths converge to it!

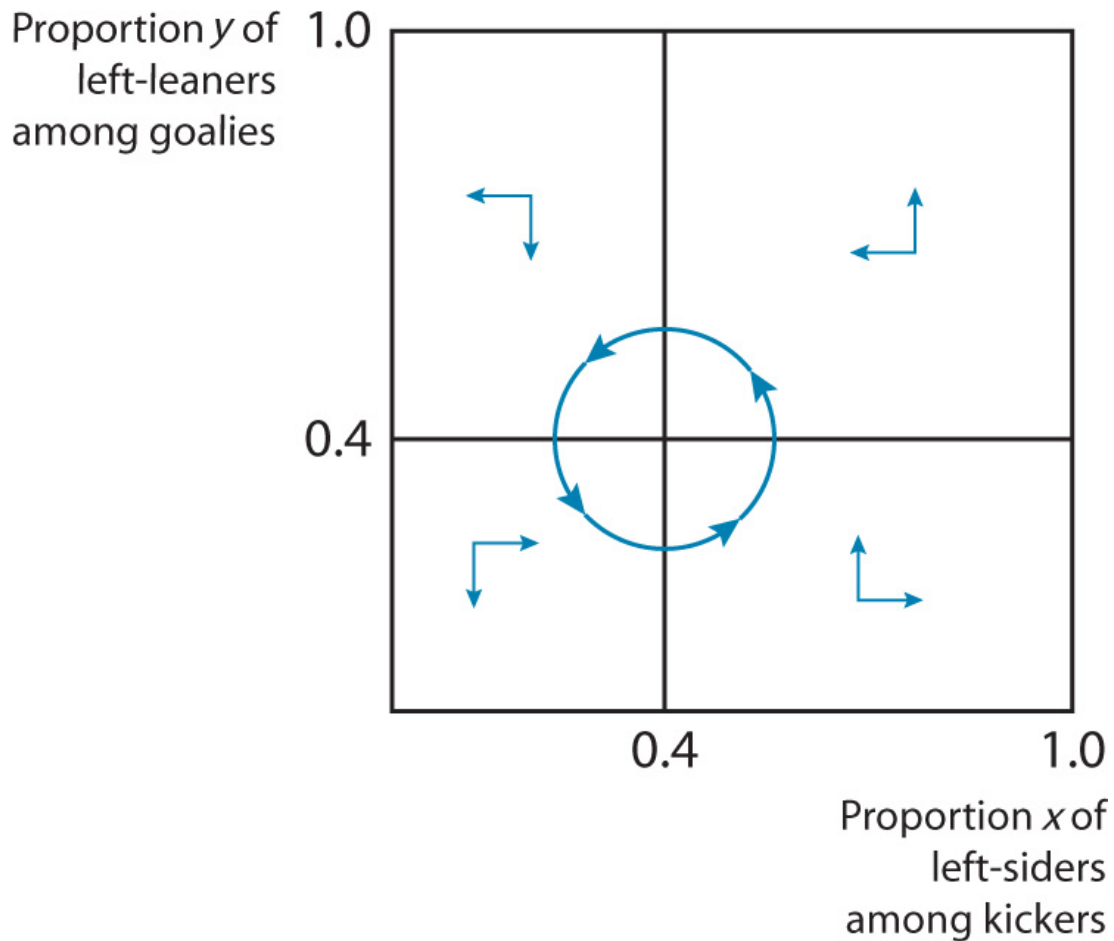


FIGURE 12.7 Dynamics for the Soccer Penalty-Kick Game

Such cyclical outcomes are well-known in some biological contexts. Consider two populations, one a predator species and the other its prey. If predators are numerous, the number of prey will decline. But when there are so few prey that the predators have little food, their numbers will also decline. Then prey can flourish again. But the resulting abundance of food enables the predators to thrive and multiply again. And so on. Such dynamics were observed and mathematically modeled by Alfred Lotka and Vito Volterra in the early twentieth century, and are named after them. [9](#)

F. Multiplayer Evolutionary Games

It is possible that evolutionary interactions may entail an entire population playing at once, rather than being matched in pairs. Such simultaneous interactions among more than two individuals are referred to as *playing the field*. In biology, all members of a flock of animals with a mixture of genetically determined behaviors may compete among themselves for some resource or territory. In such settings, the fitness of each animal depends on the strategy (phenotype) mix in the whole population. In economics or business, many firms in an industry, each following the strategy dictated by its corporate culture, may similarly compete all with all.

Such evolutionary games stand in the same relation to the rationally played collective-action games of [Chapter 11](#) as do the pair-by-pair evolutionary games of the preceding sections to the rationally played two-person games of [Chapters 4](#) through [7](#). Just as we converted the expected payoff tables of those chapters into the fitness graphs in Figures 12.3 and 12.5, we can convert the graphs for collective-action games (Figures 11.6 - 11.8) into fitness graphs for evolutionary games.

For example, consider an animal species whose members all come to a common feeding ground. There are two phenotypes: One fights for food aggressively, and the other hangs around and sneaks what it can. If the proportion of aggressive individuals is small, they will do better, but if there are too many of them, the sneakers will do better by ignoring the ongoing fights. The result will be a collective game of chicken whose fitness diagram will be exactly like Figure 11.7. Because no new principles or techniques are required, we leave it to you to pursue this idea further.

Endnotes

- Literally, the fraction of any particular type in the population can only take on values such as $1/1,000,000$, $2/1,000,000$, and so on. But if the population is sufficiently large and we show all such values as points on a straight line, as in Figure 12.3, then these points are very tightly packed together, and we can regard them as forming a continuous line. This approach amounts to letting the fractions take on any real value between 0 and 1, and it allows us to talk of the population *proportion* of a certain behavioral type, as we do throughout this chapter. By the same reasoning, if one individual member goes to jail and is removed from the population, her removal does not change the population's proportions of the various phenotypes. [Return to reference 5](#)
- *The Invasion of the Mutant Wimps* could be an interesting science-fiction comedy movie. [Return to reference 6](#)
- There can also be evolutionarily stable mixed outcomes in which each member of the population adopts a mixed strategy. We investigate this idea further in Section 5.E. [Return to reference 7](#)
- We leave this as a simple exercise, as you are by now experienced game theorists. [Return to reference 8](#)
- For readers who are good at calculus, we offer a quick proof for the cyclical behavior of the Lotka-Volterra predator-prey system. Suppose the speeds of change in each of x and y are proportional to the distance from the level of the other at which the two types are equally fit. That is, $dx/dt = a(0.4 - y)$ and $dy/dt = b(x - 0.4)$, where a and b are positive constants. Then

$$\frac{d}{dt} \left[\frac{(x-0.4)^2}{a} + \frac{(y-0.4)^2}{b} \right] = \frac{2(x-0.4)}{a} a(0.4-y) + \frac{2(y-0.4)}{b} b(x-0.4) = 0.$$

Therefore, the point (x, y) moves along a curve where $(x-0.4)^2 / a + (y-0.4)^2 / b$ is constant. This movement forms an ellipse centered on the ESS $(0.4, 0.4)$. The value of the constant, and therefore the size of the ellipse, depends on the initial positions of x and y .

[Return to reference 9](#)

3 THE REPEATED PRISONERS' DILEMMA

We saw in [Chapter 10](#) that a repetition of the prisoners' dilemma permitted consciously rational players to sustain cooperation for their mutual benefit. Let us see if a similar possibility exists in the evolutionary story. We return to the prisoners' dilemma in pricing analyzed in [Section 2.A](#). From the population of restaurateurs, we choose a pair of players to play the dilemma multiple times in succession. The overall payoff to a player (restaurateur) from such an interaction is the sum of what she gets across all rounds.

Each individual player is still programmed to play just one strategy, but that strategy has to be a complete plan of action. In a game with three rounds, for instance, a strategy can stipulate an action in the second or third play that depends on what happened in the first or second play. For example, "I will always cooperate no matter what" and "I will always defect no matter what" are valid strategies. But "I will begin by cooperating in the first round and cooperate in any later round if you cooperated in the preceding round, but defect in any later round if you defected in the preceding round" is also a valid strategy; in fact, this last strategy is tit-for-tat (TFT).

To keep the initial analysis simple, we suppose in this section that there are just two types of strategies that can possibly exist in the population of restaurateurs: always defect (A) and tit-for-tat (T). Pairs are randomly selected from the population, and each selected pair plays the game a specified number of times. The fitness of each player is simply the sum of her payoffs from all the repetitions played against her specific opponent. We examine what happens with two, three, and more generally, n such repetitions in each pair.

COLUMN	
A	T

COLUMN			
		A	T
ROW	A	576, 576	648, 504
	T	504, 648	648, 648

FIGURE 12.8 Outcomes in the Twice-Repeated Restaurant Prisoners’ Dilemma (in Hundreds of Dollars per Month)

A. Twice-Repeated Play

Figure 12.8 shows the payoff table for the game in which two members of the restaurateur population meet and play against each other exactly twice. If both players are A types, both defect both times, and we can refer back to Figure 12.1 to see that each then gets 288 each time, for a total of 576. If both are T types, defection never starts, and each gets 324 each time, for a total of 648. If one is an A type and the other a T type, then on the first play, the A type defects and the T type cooperates, so the former gets 360 and the latter 216. On the second play, both defect and get 288. So the A type's total payoff is $360 + 288 = 648$, and the T type's total is $216 + 288 = 504$.

In the twice-repeated dilemma, we see that A is weakly dominant. If the population is all A, then T-type mutants cannot invade, and A is an ESS. But if the population is all T, then A-type mutants cannot do any better than the T types. Does this mean that T must be another ESS, just as it would be a Nash equilibrium in the rational-player game-theoretic analysis of this game? The answer is no. If the population were initially all T types and a few A mutants entered, then the mutants would meet the predominant T types most of the time and would do as well as a T would do against another T. But occasionally an A mutant would meet another A mutant, and in this match she would do better than would a T against an A. Thus, the mutants would have just *slightly* higher fitness than would a member of the predominant phenotype. This advantage would lead to an increase, albeit a slow one, in the proportion of mutants in the population. Therefore, an all-T population *could* be invaded successfully by A mutants; T is not an ESS.

The static test for an ESS thus has two parts. First, we see if the mutant does better or worse than the predominant phenotype when each is matched against the predominant type. If this *primary criterion* gives a clear answer, that settles the matter. But if the primary criterion gives a tie, then we use a tie-breaking, or secondary, criterion: Does the mutant fare better or

worse than the predominant phenotype when each is matched against a mutant? Ties are exceptional, and most of the time we do not need the *secondary criterion*, but it is there in reserve for situations such as the one illustrated in Figure 12.8.[10](#)

COLUMN			
		A	T
ROW	A	864, 864	936, 792
	T	792, 936	972, 972

FIGURE 12.9 Outcomes in the Thrice-Repeated Restaurant Prisoners’ Dilemma (in Hundreds of Dollars per Month)

B. Threefold Repetition

Now suppose each matched pair from the (A, T) population plays the game three times. Figure 12.9 shows the fitness outcomes, summed over the three meetings, for each type of player when matched against a rival of each type.

To see how these fitness numbers arise, consider a couple of examples. When two T players meet each other, both cooperate the first time, and therefore both cooperate the second time and the third time as well; both get 324 each time, for a total of 972 each over 3 months. When a T player meets an A player, the latter does well the first time (360 for the A type versus 216 for the T player), but then the T player also defects the second and third times, and each gets 288 in both of those plays (for totals of 936 for A and 792 for T).

The relative fitnesses of the two types depend on the composition of the population. If the population is almost wholly A types, then A is fitter than T (because A types meeting mostly other A types earn 864 most of the time, but T types most often get 792). But if the population is almost wholly T types, then T is fitter than A (because T types earn 972 most of the time when they meet mostly other Ts, but A types earn 936 in such a situation). Each type is fitter when it already predominates in the population. Therefore, T cannot invade successfully when the population is all A, and vice versa. Now there are two possible evolutionarily stable configurations of the population: In one configuration, A is the ESS, and in the other, T is the ESS.

Next, consider the evolutionary dynamics when the initial population is made up of a mixture of the two types. How will the composition of the population evolve over time? Suppose a fraction x of the population is T and the rest, $(1 - x)$, is A. An individual A player, pitted against various opponents chosen from such a population, gets 936 when confronting a T player, which happens a fraction x of the times, and 864 against another A

player, which happens a fraction $(1 - x)$ of the times. This gives an average expected payoff of

$$936x + 864(1 - x) = 864 + 72x$$

for each A player. Similarly, an individual T player gets an average expected payoff of

$$972x + 792(1 - x) = 792 + 180x.$$

Then a T player is fitter than an A player if the former earns more on average; that is, if

$$792 + 180x > 864 + 72x$$

$$108x > 72$$

$$x > \frac{2}{3}.$$

In other words, if more than two-thirds (67%) of the population is already T, then T players will be fitter, and their proportion will grow until it reaches 100%. If the population starts with less than 67% T, then A players will be fitter, and the proportion of T players will go on declining until there are 0% of them, or 100% A players. The evolutionary dynamics move the population toward one of the two extremes, each of which is a possible ESS. These dynamics lead to the same conclusion as the static test of mutant invasion. Many, but not all, evolutionary games share this feature, that the static and dynamic tests lead to the same conclusions about ESS.

Thus, we have identified two evolutionarily stable configurations of the population. In each one, the population is all of one type (monomorphic). For example, if the population is initially 100% T, then even after a small number of mutant A types arise, the population mix will still be more than 66.66 . . . % T; T will remain the fitter type, and the mutant A strain will die out. Similarly, if the population is initially 100% A, then a small number of T-type mutants will leave the population mix with less than 66.66 . . . % T, so the A types will be fitter and the mutant T strain will die out.

If the initial population has exactly 66.66 . . . % T players (and 33.33 . . . % A players), then the two types are equally fit. However, such a polymorphism is not evolutionarily stable. The population can sustain this delicately balanced outcome only until a mutant of either type surfaces. By chance, such a mutant must arise sooner or later. The mutant's arrival will tip the fitness calculation in favor of the mutant type, and the advantage will accumulate until the ESS with 100% of that type is reached. Thus, this configuration does not meet the secondary criterion for evolutionary stability. We will sometimes loosely speak of such a configuration as an *unstable equilibrium*, so as to maintain the parallel with ordinary game theory, where mutations are not a consideration and a delicately balanced equilibrium can persist. But in the strict logic of the biological process, it is not an equilibrium at all.

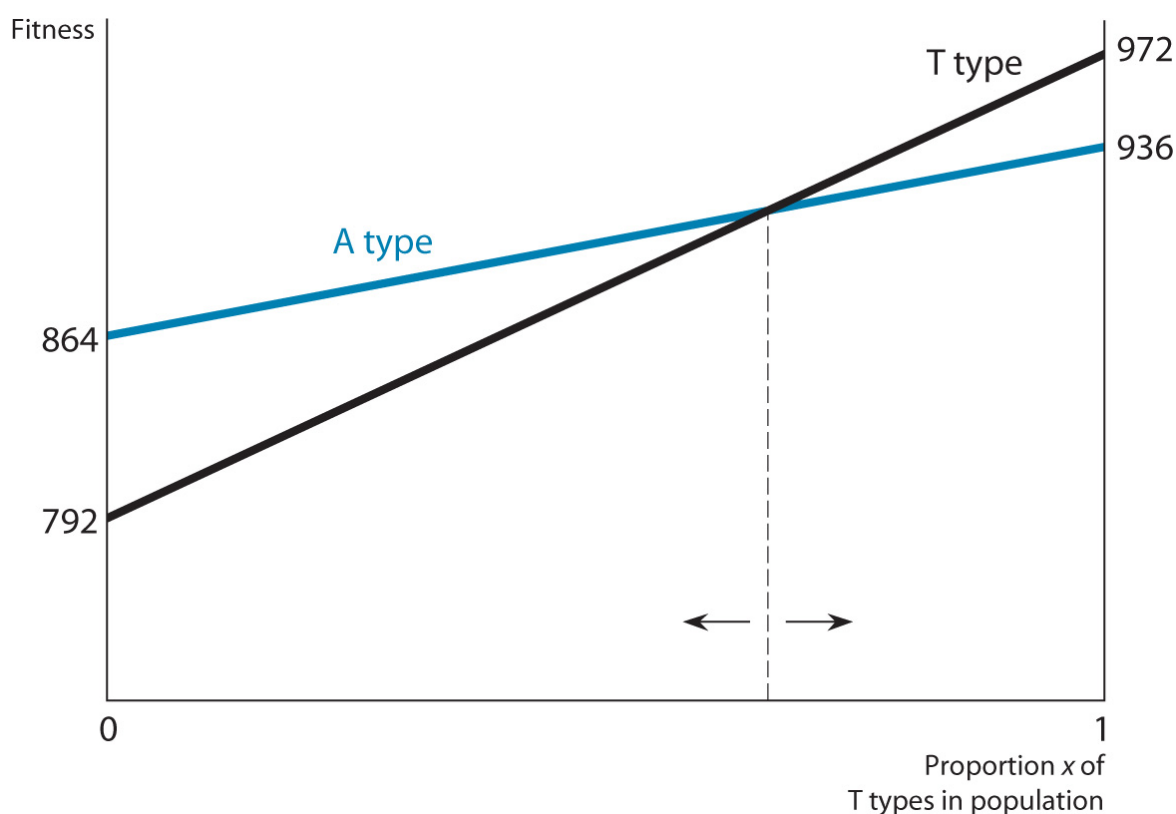


FIGURE 12.10 Fitness Graphs and Equilibria for the Thrice-Repeated Prisoners' Dilemma

This reasoning can be shown in a simple graph that closely resembles the graphs that we drew when calculating the proportions in a mixed-strategy equilibrium with consciously rational players. The only difference is that in the evolutionary context, the proportion in which the separate strategies are played is not a matter of choice by any individual player, but a property of the whole population, as shown in Figure 12.10. Along the horizontal axis, we measure the proportion x of T players in the population, which can range from 0 to 1. We measure fitness along the vertical axis. Each line shows the fitness of one type. The line for the T type starts lower (at 792, compared with 864 for the A-type line) and ends higher (972 against 936). The two lines cross where $x = 0.66$ To the right of this point, the T type is fitter, so its population proportion increases over time, and x *increases* toward 1. Similarly, to the left of this point, the A type is fitter, so its population proportion increases over time, and x *decreases* toward 0. [11](#)

C. Multiple Repetitions

What if each pair plays some unspecified number of repetitions of the game? Let us focus on a population consisting of A and T types in which interactions between random pairs occur n times (where $n > 2$). Figure 12.11 shows the total payoffs from playing n repetitions. When two A types meet, they defect, and each earns 288, every time, so each gets $288n$ in n plays. When two T types meet, they begin by cooperating, and neither is the first to defect, so each earns 324 every time, for a total of $324n$. When an A type meets a T type, on the first play the T type cooperates and the A type defects, so the A type gets 360 and the T type gets 216; thereafter, the T type retaliates against the preceding defection of the A type for all of the remaining plays, and each gets 288 in each of the remaining $(n - 1)$ plays. Thus, the A type earns a total of $360 + 288(n - 1) = 288n + 72$ in n plays against a T type, whereas the T type gets $216 + 288(n - 1) = 288n - 72$ in n plays against an A type.

		COLUMN	
		A	T
ROW	A	$288n, 288n$	$288n + 72, 288n - 72$
	T	$288n - 72, 288n + 72$	$324n, 324n$

FIGURE 12.11 Outcomes in the N -fold - Repeated Dilemma

If the proportion of T types in the population is x , then a typical A type gets $x(288n + 72) + (1 - x)288n$ on average, and a typical T type gets $x(324n) + (1 - x)(288n - 72)$ on average. Therefore, the T type is fitter if

$$x(324n) + (1 - x)(288n - 72) > x(288n + 72) + (1 - x)288n$$

$$36xn > 72$$

$$x > \frac{72}{36n} = \frac{2}{n}.$$

Once again we have two monomorphic ESSs, one all T (or $x = 1$, to which the process converges starting from any $x > 2/n$) and the other all A (or $x = 0$, to which the process converges starting from any $x < 2/n$). As in Figure 12.10, there is also an unstable polymorphic equilibrium at the balancing point $x = 2/n$.

Notice that the proportion of T at the balancing point depends on n ; it is smaller when n is larger. When $n = 10$, it is $2/10$, or 0.2. So if the population is initially 20% T players, and each pair plays 10 repetitions, the proportion of T types will grow until they reach 100% of the population. Recall that when pairs played three repetitions ($n = 3$), the T players needed an initial strength of 67% or more to achieve the same outcome, and only two repetitions meant that T types needed to begin with 100% of the population to survive. (We see the reason for this outcome in our expression for the balancing point for x , which shows that when $n = 2$, x must be above 1 before the T types are fitter.) Remember, too, that a population consisting of all T players achieves cooperation. Thus, cooperation emerges from a larger range of initial conditions when the game is repeated more times. In this sense, with more repetition, cooperation becomes more likely. What we are seeing is the result of the fact that the value of establishing cooperation increases as the length of the interaction increases.

Endnotes

- This game is just one example of a twice-repeated dilemma. With other payoffs in the basic game, twofold repetition may not produce ties. That is so in the Husband - Wife prisoners' dilemma of Chapter 4. If both the primary and secondary tests yield ties, neither phenotype satisfies our definition of ESS, and additional conceptual tools are needed that, unfortunately, lie beyond the scope of our introductory treatment of the subject. [Return to reference 10](#)
- You should now be able to draw a similar graph for the twice-repeated case. You will see that the A line is above the T line for all values of $x < 1$, but the two meet on the right-hand edge of the graph where $x = 1$. [Return to reference 11](#)

4 THE HAWK – DOVE GAME

The [hawk – dove game](#) was the first example biologists studied in their development of the theory of evolutionary games. It has instructive parallels with our analyses so far of the prisoners' dilemma and chicken, so we describe it here to reinforce and improve your understanding of the concepts involved.

The game is played not by birds of these two species, but by two animals of the same species, and Hawk and Dove are merely the names for their strategies. The context is competition for a resource. The Hawk strategy is aggressive and fights to try to get the whole resource of value V . The Dove strategy is to offer to share and avoid a fight. When two Hawk types meet each other, they fight. Each animal is equally likely (probability one-half) to win and get V or to lose, be injured, and get $-C$. Thus, the expected payoff for each is $(V - C)/2$. When two Dove types meet, they share without a fight, so each gets $V/2$. When a Hawk type meets a Dove type, the latter retreats and gets 0, whereas the former gets V . Figure 12.12 shows the payoff table for this game.

The analysis of the hawk – dove game is similar to that for the prisoners' dilemma and chicken, except that the numerical payoffs have been replaced with algebraic symbols. Here we compare the equilibria of this game when each player has a specified strategy, and success is rewarded with faster reproduction of that strategy in the population.

A. Rational Strategic Choice and Equilibrium

If $V > C$, then the game is a prisoners' dilemma in which the Hawk strategy corresponds to Defect and the Dove strategy corresponds to Cooperate. Hawk is the dominant strategy for each player, but (Dove, Dove) is the jointly better outcome.

If $V < C$, then it's a game of chicken. Now $(V - C)/2 < 0$, and so Hawk is no longer a dominant strategy. Rather, there are two pure-strategy Nash equilibria: (Hawk, Dove) and (Dove, Hawk). There is also a mixed-strategy equilibrium, where B's probability p of choosing Hawk is such as to keep A indifferent, as defined by

$$\frac{p(V - C)}{2} + (1 - p)V = p \times 0 + \frac{(1 - p)V}{2}$$
$$p = \frac{V}{C}.$$

B. Evolutionary Stability for $V > C$

We start with a population in which Hawks predominate and test whether it can be invaded by mutant Doves. Following the convention used in analyzing such games, we could write the population proportion of the mutant phenotype as m , for mutant, but for clarity in our case we will use d for mutant Dove. The population proportion of Hawks is then $(1 - d)$. Then, in a match against a randomly drawn opponent, a Hawk will meet a Dove a proportion d of the time and get V on each of those occasions, and will meet another Hawk a proportion $(1 - d)$ of the time and get $(V - C)/2$ on each of those occasions. Therefore, the fitness of a Hawk is $dV + (1 - d)(V - C)/2$. Similarly, the fitness of a mutant dove is $d(V/2) + (1 - d) \times 0$. Because $V > C$, it follows that $(V - C)/2 > 0$. Also, $V > 0$ implies that $V > V/2$. Then, for any value of d between 0 and 1, we have

$$dV + \frac{(1 - d)(V - C)}{2} > d\frac{V}{2} + (1 - d) \times 0,$$

and so the Hawk type is fitter. The Dove mutants cannot successfully invade. The Hawk strategy is evolutionarily stable, and the population is monomorphic (all Hawk).

		B	
		Hawk	Dove
A	Hawk	$(V - C)/2, (V - C)/2$	$V, 0$
	Dove	$0, V$	$V/2, V/2$

FIGURE 12.12 Payoff Table for the Hawk - Dove Game

The same holds true for any population proportion of Doves—that is, for all values of d . Therefore, from any initial mix, the proportion of Hawks will grow, and they will predominate. In addition, if the population is initially all Doves, mutant Hawks

can invade and take over. Thus, the dynamics confirm that the Hawk strategy is the only ESS. This algebraic analysis affirms and generalizes our earlier finding for the numerical example of the prisoners' dilemma of restaurant pricing (see Figure 12.1).

C. Evolutionary Stability for $V < C$

If the initial population is again predominantly Hawks, with a small proportion d of Dove mutants, then each has the same fitness function derived in [Section 4.B](#). When $V < C$, however, $(V - C)/2 < 0$. We still have $V > 0$, and so $V > V/2$. But because d is very small, the comparison of the terms with $(1 - d)$ is much more important than that of the terms with d , so

$$d\frac{V}{2} + (1 - d) \times 0 > dV + \frac{(1 - d)(V - C)}{2}.$$

Thus, the Dove mutants are fitter than the predominant Hawks and can invade successfully.

But if the initial population is almost all Doves, then we must consider whether a small proportion h of Hawk mutants can invade. (Note that because the mutant is now a Hawk, we have used h for the proportion of the mutant invaders.) The Hawk mutants have a fitness of $h(V - C)/2 + (1 - h)V$ compared with $h \times 0 + (1 - h)(V/2)$ for the Doves. Again, $V < C$ implies that $(V - C)/2 < 0$, and $V > 0$ implies that $V > V/2$. But when h is small, we get

$$\frac{h(V - C)}{2} + (1 - h)V > h \times 0 + (1 - h)\frac{V}{2}.$$

This inequality shows that Hawks are fitter and will successfully invade a Dove population. Thus, mutants of each type can invade populations of the other type. The population cannot be monomorphic, and neither pure phenotype can be an ESS. The algebra again confirms our earlier finding for the numerical example of chicken (see Figures 12.2 and 12.3).

What happens in the population, then, when $V < C$? There are two possibilities. In one, every player follows a pure strategy, but the population has a stable mix of players following different strategies. This is the polymorphic equilibrium developed for

chicken in [Section 2.C](#). The other possibility is that every player uses a mixed strategy. We begin with the polymorphic case.

D. $V < C$: Stable Polymorphic Population

When the population proportion of Hawks is h , the fitness of a Hawk is $h(V - C)/2 + (1 - h)V$, and the fitness of a Dove is $h \times 0 + (1 - h)(V/2)$. The Hawk type is fitter if

$$\frac{h(V - C)}{2} + (1 - h)V > (1 - h)\frac{V}{2},$$

which simplifies to

$$\frac{h(V - C)}{2} + (1 - h)\frac{V}{2} > 0$$

$$V - hC > 0$$

$$h < \frac{V}{C}.$$

The Dove type is then fitter when $h > V/C$, or when $(1 - h) < 1 - (V/C) = (C - V)/C$. Thus, each type is fitter when it is rarer. Therefore, we have a stable polymorphic equilibrium at the balancing point, where the proportion of Hawks in the population is $h = V/C$. This is exactly the probability with which each individual member plays the Hawk strategy in the mixed-strategy Nash equilibrium of the game under the assumption of rational behavior, as calculated in [Section 4.A](#). Again, we have an

evolutionary “justification” for the mixed-strategy outcome in chicken.

We leave it to you to draw a graph similar to that in Figure 12.3 for this case. Doing so will require you to determine the dynamics by which the population proportions of each type converge to the stable equilibrium mix.

E. $V < C$: Each Player Mixes Strategies

Recall the equilibrium mixed strategy of the rational-play game calculated in [Section 4.A](#), in which $p = V/C$ was the probability of choosing to be a Hawk, while $(1 - p)$ was the probability of choosing to be a Dove. Is there a parallel in the evolutionary version, with a phenotype playing a mixed strategy? Let us examine this possibility. We still have types who play the pure Hawk strategy, which we call H, and types who play the pure Dove strategy, called D. But now a third phenotype, called M, can exist; such a type plays a mixed strategy in which it is a Hawk with probability $p = V/C$ and a Dove with probability $1 - p = 1 - (V/C) = (C - V)/C$.

When an H or a D meets an M, their expected payoffs depend on p , the probability that M is playing H, and on $(1 - p)$, the probability that M is playing D. Then each player gets p times her payoff against an H, plus $(1 - p)$ times her payoff against a D. So when an H type meets an M type, she gets the expected payoff

$$\begin{aligned} p \frac{V - C}{2} + (1 - p)V &= \frac{V}{C} \frac{V - C}{2} - \frac{C - V}{C} V \\ &= -\frac{1}{2} \frac{V}{C} (C - V) + \frac{V}{C} (C - V) \\ &= V \frac{C - V}{2C}. \end{aligned}$$

And when a D type meets an M type, she gets

$$p \times 0 + (1 - p) \frac{V}{2} = \frac{C - V}{V} \frac{V}{2} = \frac{V(C - V)}{2V}.$$

The two fitnesses are equal. This should not be a surprise; the proportions of the mixed strategy are determined to achieve exactly this equality. So an M type meeting another M type also gets the same expected payoff. For brevity of future reference, we call this common payoff K , where $K = V(C - V)/2C$.

But these equalities create a problem when we test M for evolutionary stability. Suppose the population consists entirely of M types and that a few mutants of the Hawk type, constituting a very small proportion h of the total population, invade. Then the typical mutant gets the expected payoff $h(V - C)/2 + (1 - h)K$. To calculate the expected payoff of an M type, note that she faces another M type in a fraction $(1 - h)$ of the interactions and gets K in each instance. She then faces an H type in a fraction h of the interactions; in these interactions, she plays H a fraction p of the time and gets $(V - C)/2$, and she plays D a fraction $(1 - p)$ of the time and gets 0. Thus, the M type's total expected payoff (fitness) is

$$\frac{hp(V - C)}{2} + (1 - h)K.$$

Because h is very small, the fitnesses of the M type and the mutant H type are almost equal. The point is that when there are very few mutants, both the H type and the M type meet only M types most of the time, and in this interaction the two have equal fitness, as we just saw.

Evolutionary stability hinges on whether the M type is fitter than the mutant H type when each is matched against one of the few mutants. Algebraically, M is fitter than H against other mutant H types when $pV(C - V)/2C = pK > (V - C)/2$. In our example here, this condition holds because $V < C$ [so $(V - C)$ is negative] and because K is positive. Intuitively, this condition tells us that an H-type mutant will always do badly against another H-type mutant because of the high cost of fighting, but the M type fights only part of the time and therefore suffers

this cost only a fraction p of the time. Overall, the M type does better than the H type when matched against a mutant.

Similarly, the success of a Dove invasion of the M population depends on the comparison between a mutant Dove's fitness and the fitness of an M type. As before, the mutant faces another D in a fraction d of the interactions and faces an M in a fraction $(1 - d)$ of the interactions. An M type faces another M type in a fraction $(1 - d)$ of the interactions, but in a fraction d of the interactions, the M faces a D; she plays H a fraction p of those times, thereby gaining pV , and plays D a fraction $(1 - p)$ of those times, thereby gaining $(1 - p)V/2$. The Dove's fitness is then $dV/2 + (1 - d)K$, while the fitness of the M type is $d[pV + (1 - p)V/2] + (1 - d)K$. The final term in each fitness expression is the same, so a Dove invasion is successful only if $V/2$ is greater than $pV + (1 - p)V/2$. This condition does not hold; the latter expression includes a weighted average of V and $V/2$ that must exceed $V/2$ whenever $V > 0$. Thus, the Dove invasion cannot succeed either.

This analysis tells us that M is an ESS. Thus, if $V < C$, the population can exhibit either of two evolutionarily stable outcomes. One entails a mixture of types (a stable polymorphism), and the other entails a single type that mixes its strategies in the same proportions that define the polymorphism.

F. Some General Theory

We can now generalize the ideas illustrated in this section to get a theoretical framework and a set of tools that can then be applied further. This generalization unavoidably requires some slightly abstract notation and a bit of algebra. Therefore, we cover only monomorphic equilibria in a single species. Readers who are adept at this level of mathematics can readily develop the polymorphism cases with two species by analogy. Readers who are either unprepared for this material or uninterested in it can omit this section without loss of continuity.[¹²](#)

We consider random pairings from a single species whose population has available strategies I, J, K Some of them may be pure strategies; some of them may be mixed. Each individual member of the population is hardwired to play just one of these strategies. We let $E(I, J)$ denote the payoff to an I player in a single encounter with a J player. The payoff of an I player meeting another of her own type is $E(I, I)$ in the same notation. We write $W(I)$ for the fitness of an I player. This is simply her expected payoff in encounters with randomly picked opponents, when the probability of her meeting a type is simply the proportion of that type in the population.

Suppose the population is all I type. We consider whether this can be an evolutionarily stable configuration. To do so, we imagine that the population is invaded by a few J-type mutants, so the proportion of mutants in the population is a very small number, m . Now the fitness of an I type is

$$W(I) = mE(I, J) + (1 - m)E(I, I),$$

and the fitness of a mutant is

$$W(J) = mE(J, J) + (1 - m)E(J, I).$$

Therefore, the difference in fitness between the population's predominant type and its mutant type is

$$W(I) - W(J) = m[E(I, J) - E(J, J)] + (1 - m)[E(I, I) - E(J, I)].$$

Because m is very small, the predominant type's fitness will be higher than the mutant's if the second half of the preceding expression is positive; that is,

$$W(I) > W(J) \text{ if } E(I, I) > E(J, I).$$

Then the population meets the [primary criterion](#) for evolutionary stability: It cannot be invaded, because the predominant type is fitter than the mutant type when each is matched against a member of the predominant type. Conversely, if $W(I) < W(J)$, owing to $E(I, I) < E(J, I)$, the J-type mutants can invade successfully, and an all-I population cannot be evolutionarily stable.

However, it is possible that $E(I, I) = E(J, I)$, as indeed happens if the population initially consists of a single phenotype that plays a strategy of mixing between the pure strategies I and J (a monomorphic equilibrium with a mixed strategy), as was the case in our final variant of the hawk-dove game (see [Section 4.E](#)). Then the difference between $W(I)$ and $W(J)$ is governed by how each type fares against the mutants.¹³ When $E(I, I) = E(J, I)$, we get $W(I) > W(J)$ if $E(I, J) > E(J, J)$, which indicates that the population meets the [secondary criterion](#) for the evolutionary stability of I. This criterion is invoked only if the primary criterion is inconclusive—that is, only if $E(I, I) = E(J, I)$.

If the secondary criterion is invoked—because $E(I, I) = E(J, I)$ —there is the additional possibility that it may also be inconclusive. That is, it may also be the case that $E(I, J) = E(J, J)$. If both the primary and secondary criteria for the evolutionary stability of I are inconclusive, then I is considered a *neutral ESS*. Thorough analysis of this special case is beyond the scope of our introductory text, so we leave it to those with a deeper interest in the topic to explore more advanced treatments of the theory.

The primary criterion says that if the strategy I is evolutionarily stable, then for all other strategies J that a mutant might try, $E(I, I) \geq E(J, I)$. This means that I is the

best response to itself. In other words, if the members of this population suddenly started playing as rational calculators, all members playing I would be a Nash equilibrium. We explained this in the context of earlier examples; here we see the result in its general theoretical definition.

This is a remarkable result. If you were dissatisfied with the rational-play assumption underlying the theory of Nash equilibria given in earlier chapters and you came to the theory of evolutionary games looking for a better explanation, you will find that it yields the same results. The very appealing biological description—fixed nonmaximizing behavior, but selection in response to resulting fitness—does not yield any new outcomes. If anything, it provides a backdoor justification for Nash equilibrium. When a game has several Nash equilibria, the evolutionary dynamics may even provide a good argument for choosing among them.

However, your reinforced confidence in Nash equilibrium should be cautious. Our definition of evolutionary stability is static rather than dynamic. It requires only that the configuration of the population (monomorphic, or polymorphic in just the right proportions) that we are testing for equilibrium cannot be successfully invaded by a small proportion of mutants. It does not test whether, starting from an arbitrary initial population mix, all the types other than the fittest will die out and the equilibrium configuration will be reached. And the test is carried out for those particular classes of mutants that are deemed logically possible; if the theorist has not specified this classification correctly and some type of mutant that she overlooked could actually arise, that mutant might invade successfully and destroy the supposed equilibrium. Finally, as we showed in the soccer penalty kick example of [Section 2.E](#), evolutionary dynamics can fail to converge at all.

Endnotes

- Conversely, readers who want more details can find them in Maynard Smith, *Evolution and the Theory of Games*, especially pp. 14–15. John Maynard Smith is a pioneer in the theory of evolutionary games. [Return to reference 12](#)
- If the initial population is polymorphic and m is the proportion of J types, then m may not be “very small” any more. The size of m is no longer crucial, however, because the second term in $W(I) - W(J)$ is now assumed to be zero. [Return to reference 13](#)

Glossary

[hawk - dove game](#)

An evolutionary game where members of the same species or population can breed to follow one of two strategies, Hawk and Dove, and depending on the payoffs, the game between a pair of randomly chosen members can be either a prisoners' dilemma or chicken.

[primary criterion](#)

Comparison of the fitness of a mutant with that of a member of the dominant population, when each plays against a member of the dominant population.

[secondary criterion](#)

Comparison of the fitness of a mutant with that of a member of the dominant population, when each plays against a mutant.

5 EVOLUTION OF COOPERATION AND ALTRUISM

Evolutionary game theory rests on two fundamental ideas: first, that individual organisms are engaged in games with others of their own species or with members of other species, and second, that the genotypes that lead to higher-payoff (fitter) strategies proliferate and increase in their proportions of the population while the rest decline. These ideas suggest a vicious struggle for survival like that depicted by some interpreters of Darwin, who understood “survival of the fittest” in a literal sense and who conjured up images of a “nature red in tooth and claw.” In fact, nature shows many instances of cooperation (in which individual animals behave in a way that yields the greatest benefit to everyone in a group) and even altruism (in which individual animals incur significant costs in order to benefit others). Beehives and ant colonies are only the most obvious examples. Can such behavior be reconciled with the perspective of evolutionary games?

Biologists use a fourfold classification of the ways in which cooperation and altruism can emerge among selfish animals (or phenotypes or genes): (1) family dynamics, (2) reciprocal altruism, (3) selfish teamwork, and (4) group altruism.¹⁴ The behavior of ants and bees is probably the easiest to understand as an example of family dynamics. All the individual members of an ant colony or a beehive are closely related and have genes in common to a substantial extent. Most worker ants in a colony are full sisters and therefore have half their genes in common; the survival and proliferation of one ant’s genes is helped just as much by the survival of two of her sisters as by her own survival. Most worker bees in a hive are half-sisters and therefore

have a quarter of their genes in common. An individual ant or bee does not make a fine calculation of whether it is worthwhile to risk her own life for the sake of two or four sisters, but the underlying genes of those groups whose members exhibit such behavior (phenotype) will proliferate. The idea that evolution ultimately operates at the level of the gene has had enormous implications for biology, although it has been misapplied by many people, just as Darwin's original idea of natural selection has been misapplied.¹⁵ The interesting idea is that a "selfish gene" may prosper by behaving unselfishly in a larger organization of genes, such as a cell. Similarly, a cell and its genes may prosper by participating cooperatively and accepting their allotted tasks in a body.

Reciprocal altruism can arise among unrelated individual members of the same or different species. This behavior is essentially an example of the resolution of prisoners' dilemmas through repetition, in which the players use strategies that are remarkably like tit-for-tat. For example, some small fish and shrimp thrive on parasites that collect in the mouths and gills of some large fish; the large fish let the small ones swim unharmed through their mouths in return for this "cleaning service." Another fascinating, although more gruesome, example is that of vampire bats, who share blood with those who have been unsuccessful in their own hunting. In an experiment in which bats from different home sites were brought together and selectively starved, "only bats that were on the verge of starving (that is, would die within 24 hours without a meal) were given blood by any other bat in the experiment. But, more to the point, individuals were given a blood meal only from bats they already knew from their site. . . . Furthermore, vampires were much more likely to regurgitate blood to the specific individual(s) from their site that had come to their aid when they needed a bit of blood."¹⁶ Once again, it is not to be supposed that each animal consciously calculates whether it

is in its individual interest to continue the cooperation or to defect. Instead, its behavior is instinctive.

Selfish teamwork arises when it is in the interest of each individual organism to choose cooperation when all others are doing so. In other words, this concept of cooperative behavior applies to the selection of the good outcome in assurance games; for example, populations are more likely to engage in selfish teamwork in harsh environments than in mild ones. When conditions are bad, shirking by any one animal in a group could bring disaster to the whole group, including the shirker. In such conditions, each animal is crucial for the group's survival, and none shirk so long as others are also pulling their weight. In milder environments, each may hope to become a free rider on the others' efforts without thereby threatening the survival of the whole group, including itself.

The next step goes beyond biology and into sociology: A body (and its cells and, ultimately, its genes) may benefit by behaving cooperatively in a collection of bodies—namely, a society. This idea suggests that cooperation can arise even among individual members of a group who are not close relatives. We do indeed find instances of such behavior, which falls into the final category, group altruism. Groups of predators such as wolves are a case in point, and groups of apes often behave like extended families. Even among species of prey, such cooperation arises, as when individual fish in a school take turns looking out for predators. And cooperation can also extend across species. The general idea is that a group whose members behave cooperatively is more likely to succeed in its interactions with other groups than one whose members seek the benefit of free riding within the group. If, in a particular context of evolutionary dynamics, between-group selection is a stronger force than within-group selection, then we will see group altruism. [17](#)

An instinct is hardwired into an individual organism's brain by genetics, but reciprocity and cooperation can arise from more purposive thinking or experimentation within the group and can spread by socialization—through explicit instruction or observation of the behavior of elders—instead of genetics. The relative importance of the two channels—nature and nurture—will differ from one species to another and from one situation to another. One would expect socialization to be relatively more important among humans, but there are instances of its role among other animals. We cite a remarkable one. The expedition that Robert F. Scott led to the South Pole in 1911–1912 used teams of Siberian sled dogs. This group of dogs, brought together and trained for this specific purpose, developed within a few months a remarkable system of cooperation and sustained it by using punishment schemes. “They combined readily and with immense effect against any companion who did not pull his weight, or against one who pulled too much . . . their methods of punishment always being the same and ending, if unchecked, in what they probably called justice, and we called murder.” ¹⁸

This is an encouraging account of how cooperative behavior can be compatible with evolutionary game theory, and one that suggests that dilemmas of selfish action can be overcome. Indeed, scientists investigating altruistic behavior have recently reported experimental support for the existence of *altruistic punishment*, or *strong reciprocity* (as distinguished from reciprocal altruism), in humans. Their evidence suggests that people are willing to punish those who don't pull their own weight in a group setting, even when it is costly to do so and when there is no expectation of future gain. This tendency toward strong reciprocity may even help to explain the rise of human civilization if groups with this trait were better able to survive the traumas of war and other catastrophic events. ¹⁹ Despite these findings, strong reciprocity may not be widespread in the animal world.

“Compared to nepotism, which accounts for the cooperation of

ants and every creature that cares for its young, reciprocity has proved to be scarce. This, presumably, is due to the fact that reciprocity requires not only repetitive interactions, but also the ability to recognize other individuals and keep score.” ²⁰ In other words, precisely those conditions that our theoretical analysis in [Section 2.D](#) of [Chapter 10](#) identified as being necessary for a successful resolution of the repeated prisoners’ dilemma are seen to be relevant in the context of evolutionary games.

Endnotes

- For an excellent exposition, see Lee Dugatkin' s *Cheating Monkeys and Citizen Bees: The Nature of Cooperation in Animals and Humans* (Cambridge, Mass.: Harvard University Press, 2000). [Return to reference 15](#)
- In this very brief account, we cannot begin to do justice to these debates or the issues involved. An excellent popular account, and the source of many examples cited in this section, is Matt Ridley, *The Origins of Virtue* (New York: Penguin, 1996). We should also point out that we do not examine the connection between genotypes and phenotypes, or the role of sex in evolution, in any detail. Another book by Ridley, *The Red Queen* (New York: Penguin, 1995), gives a fascinating treatment of this subject. [Return to reference 15](#)
- Dugatkin, *Cheating Monkeys*, p. 99. [Return to reference 16](#)
- Group altruism used to be thought impossible according to the strict theory of evolution that emphasizes selection at the level of the gene, but the concept is being revived in more sophisticated formulations. See Dugatkin, *Cheating Monkeys*, pp. 141 – 145, for a fuller discussion. [Return to reference 17](#)
- Apsley Cherry-Garrard, *The Worst Journey in the World* (London: Constable, 1922; reprint, New York: Carroll and Graf, 1989), pp. 485 – 86. [Return to reference 18](#)
- For the evidence on altruistic punishment, see Ernst Fehr and Simon Gächter, “Altruistic Punishment in Humans,” *Nature*, vol. 415 (January 10, 2002), pp. 137 – 40. [Return to reference 19](#)
- Ridley, *Origins of Virtue*, p. 83. [Return to reference 20](#)

SUMMARY

The biological theory of evolution parallels the theory of games used by social scientists. Evolutionary games are played by behavioral *phenotypes* with genetically predetermined, rather than rationally chosen, strategies. In an evolutionary game, phenotypes with higher *fitness* survive repeated interactions with others to reproduce and thus increase their representation in the population. A population containing one or more phenotypes in certain proportions is called *evolutionarily stable* if it cannot be *invaded* successfully by other, *mutant* phenotypes or if it is the endpoint of the dynamics of proliferation of fitter phenotypes. If one phenotype maintains its dominance in the population when faced with an invading mutant type, that phenotype is said to be an *evolutionarily stable strategy* (*ESS*), and a population consisting of that phenotype alone is said to be *monomorphic*. If two or more phenotypes coexist in an evolutionarily stable population, it is said to be *polymorphic*.

When the theory of evolutionary games is applied more generally to nonbiological games, the strategies followed by individual players are understood to be standard operating procedures or rules of thumb, instead of being genetically fixed. The process of reproduction stands for more general methods of transmission, including socialization, education, and imitation; mutations represent experimentation with new strategies.

Evolutionary games may have payoff structures similar to those analyzed in [Chapters 4](#) and [7](#), including the prisoners' dilemma, games of chicken, and assurance games. In each case, the evolutionarily stable strategy mirrors either the pure-strategy Nash equilibrium of a game with the same structure

played by rational players or the proportions of the equilibrium mixture in such a game. In one-time play of the prisoners' dilemma, Defect is the evolutionarily stable strategy; in chicken; types are fitter when rare, so there is a polymorphic equilibrium. In the assurance game, types are less fit when rare, so the polymorphic configuration is unstable and the equilibria are at the extremes. When play is between two different types of members of each of two different species, a more complex, but similarly structured, analysis is used to determine equilibria. And in repeated play of the prisoners' dilemma, Always defect is evolutionarily stable.

The *hawk - dove game* is the classic biological example of an evolutionary game. Analysis of this game parallels that of the prisoners' dilemma and chicken; evolutionarily stable strategies depend on the specifics of the payoff structure.

KEY TERMS

[evolutionary stability](#) ([471](#))

[evolutionarily stable strategy \(ESS\)](#) ([473](#))

[fitness](#) ([470](#))

[genotype](#) ([470](#))

[hawk - dove game](#) ([490](#))

[invasion](#) ([470](#))

[monomorphism](#) ([473](#))

[mutation](#) ([470](#))

[phenotype](#) ([470](#))

[playing the field](#) ([471](#))

[polymorphism](#) ([473](#))

[primary criterion](#) ([495](#))

[secondary criterion](#) ([495](#))

[selection](#) ([470](#))

Glossary

[evolutionarily stable](#)

A population is evolutionarily stable if it cannot be successfully invaded by a new mutant phenotype.

[evolutionarily stable strategy \(ESS\)](#)

A phenotype or strategy that can persist in a population, in the sense that all the members of a population or species are of that type; the population is evolutionarily stable (static criterion). Or, starting from an arbitrary distribution of phenotypes in the population, the process of selection will converge to this strategy (dynamic criterion).

[fitness](#)

The expected payoff of a phenotype in its games against randomly chosen opponents from the population.

[genotype](#)

A gene or a complex of genes, which give rise to a phenotype and which can breed true from one generation to another. (In social or economic games, the process of breeding can be interpreted in the more general sense of teaching or learning.)

[phenotype](#)

A specific behavior or strategy, determined by one or more genes. (In social or economic games, this can be interpreted more generally as a customary strategy or a rule of thumb.)

[selection](#)

The dynamic process by which the proportion of fitter phenotypes in a population increases from one generation to the next.

[invasion](#)

The appearance of a small proportion of mutants in the population.

[mutation](#)

Emergence of a new genotype.

playing the field

A many-player evolutionary game where all animals in the group are playing simultaneously, instead of being matched in pairs for two-player games.

monomorphism

All members of a given species or population exhibit the same behavior pattern.

polymorphism

An evolutionarily stable equilibrium in which different behavior forms or phenotypes are exhibited by subsets of members of an otherwise identical population.

hawk - dove game

An evolutionary game where members of the same species or population can breed to follow one of two strategies, Hawk and Dove, and depending on the payoffs, the game between a pair of randomly chosen members can be either a prisoners' dilemma or chicken.

primary criterion

Comparison of the fitness of a mutant with that of a member of the dominant population, when each plays against a member of the dominant population.

secondary criterion

Comparison of the fitness of a mutant with that of a member of the dominant population, when each plays against a mutant.

SOLVED EXERCISES

1. Two travelers buy identical handcrafted souvenirs and pack them in their respective suitcases for their return flight. Unfortunately, the airline manages to lose both suitcases. Because the airline doesn't know the value of the lost souvenirs, it asks each traveler to report a value independently. The airline agrees to pay each traveler an amount equal to the minimum of the two reports. If one report is higher than the other, the airline takes a penalty of \$20 away from the traveler with the higher report and gives \$20 to the traveler with the lower report. If the two reports are equal, there is no reward or penalty. Neither traveler remembers exactly how much the souvenir cost, so that value is irrelevant; each traveler simply reports the value that her type determines she should report.

There are two types of travelers. The High type always reports \$100, and the Low type always reports \$50. Let h represent the proportion of High types in the population.

1. Draw the payoff table for the game played between two travelers selected at random from the population.
 2. Graph the fitness of the High type, with h on the horizontal axis. On the same figure, graph the fitness of the Low type.
 3. Describe all the equilibria of this game. For each equilibrium, state whether it is monomorphic or polymorphic and whether it is stable.
2. Consider a population in which there are two phenotypes: natural-born cooperators (who do not confess under questioning) and natural-born defectors (who confess readily). If two members of this population are drawn at random to play a prisoners' dilemma game, their payoffs in a single play are the same as those in the Husband - Wife prisoners' dilemma game of [Chapter 4](#), as shown below.

1.

		COLUMN	
		Confess (Defect)	Not (Cooperate)
ROW	Confess (Defect)	10 yr, 10 yr	1 yr, 25 yr
	Not (Cooperate)	25 yr, 1 yr	3 yr, 3 yr
You may need to scroll left and right to see the full figure.			

-
1. Suppose that a pair of players plays this dilemma twice in succession. Draw the payoff table for the twice-repeated dilemma

- as in Figure 12.8, with A (always defect) and T (play tit-for-tat, starting with not defecting) as the two available strategies.
- Find all of the ESS in the two-type twice-repeated dilemma in part (a).
 - Now, suppose that there is a third phenotype, N, which never defects. Draw an expanded three-by-three payoff table for the twice-repeated dilemma, with A, T, and N as the three available strategies.
 - Suppose that the population consists of an equal mix of tit-for-tat (T) and never-defect (N) types, but there are no always-defect (A) types. Can such a population be successfully invaded by A mutants?
 - Building on part (d), suppose that fraction q_T of the population are T types and the rest are N types. Verify that there is a balancing point, q^* , $0 < q^* < 1$, such that the population can be invaded by A mutants if $q_T < q^*$, but not if $q_T > q^*$. What is q^* ?
3. Consider the thrice-repeated restaurant pricing prisoners' dilemma we studied in [Section 3.B](#), but now suppose that there are three phenotypes in the population: type A, which always defects; type T, which plays tit-for-tat; and type N, which never defects. For your convenience, the expanded payoff table for this thrice-repeated three-phenotype evolutionary game is provided below. For each of parts (a) – (e), below, be specific and explicit in your answers; make sure you use the payoff numbers in the table provided.

		COLUMN		
		A	T	N
ROW	A	864, 864	936, 792	1080, 648
	T	792, 936	972, 972	972, 972
	N	648, 1080	972, 972	972, 972
You may need to scroll left and right to see the full figure.				

- Explain why a population that is 100% type A cannot be invaded by either type N or type T mutants.
- Explain why a population that is 100% type N can be invaded by type A mutants.
- Explain why a population that is 100% type T cannot be invaded by type A mutants.
- Suppose that the population initially consists of an equal mix of the three phenotypes. Which of the three phenotypes will have the highest fitness initially and grow in number the fastest? Which

will have the lowest fitness and hence fall in number the fastest?

5. (Optional) Given the initial evolutionary dynamics described in part (d), how will population dynamics continue over time? In particular, will the population eventually consist of only one phenotype? If so, which one?
4. In the assurance (meeting coordination) game described in [Section 2.D](#), the payoffs could be thought of as describing the value of something material that the players gained in the various outcomes; they could be prizes given for a successful meeting, for example. Then other individuals in the population might observe the expected payoffs (fitnesses) of the two types, see which was higher, and gradually imitate the fitter strategy. Thus, the proportions of the two types in the population would change. But we can make a more biological interpretation. Suppose the column players are always female and the row players always male. When two players of the same type meet successfully, they pair off, and their children are of the same type as the parents. Therefore, the types would proliferate or die off as a result of successful or unsuccessful meetings. The formal mathematics of this new version of the game makes it a “two-species game” (although the biology of it does not). Thus, the proportion of S-type females in the population—call this proportion x —need not equal the proportion of S-type males—call this proportion y .
 1. Examine the dynamics of x and y by using methods similar to those used in the chapter for the penalty kick game.
 2. Find the stable outcome or outcomes of this dynamic process.
5. Recall from Exercise S1 the travelers reporting the value of their lost souvenirs. Assume that a third traveler phenotype exists in the population. The third traveler type is a mixer; she plays a mixed strategy, sometimes reporting a value of \$100 and sometimes reporting a value of \$50.
 1. Use your knowledge of mixed strategies in rationally played games to posit a reasonable mixture for the mixer phenotype to use in this game.
 2. Draw the three-by-three payoff table for this game when the mixer type uses the mixed strategy that you found in part (a).
 3. Determine whether the mixer strategy is an ESS of this game. (Hint: Test whether a mixer population can be invaded successfully by either the High type or the Low type.)
6. Consider a simplified model in which everyone gets electricity either from solar power or from fossil fuels, which are both in fixed supply. (In the case of solar power, think of the required equipment as being in fixed supply.) The up-front costs of using solar power are high, so when the price of fossil fuels is low (that is, when few

people are using fossil fuels and there is a high demand for solar equipment), the cost of solar power can be prohibitive. In contrast, when many individuals are using fossil fuels, the demand for them (and thus the price) is high, whereas the demand (and thus the price) for solar power is relatively lower. Assume the payoff table for the two types of energy consumers to be as follows:

		COLUMN	
		Solar	Fossil fuels
ROW	Solar	2, 2	3, 4
	Fossil fuels	4, 3	2, 2

-
- Does this evolutionary game have a monomorphic ESS?
 - Verify that this evolutionary game has a polymorphic ESS. What is the share s^* of the population that adopts solar in this ESS? [Hint: Construct a fitness graph like Figure 12.3 or Figure 12.5 for the two types. The balancing point (corresponding to the mixed-strategy Nash equilibrium of the game with rational players) is an ESS if, away from the balancing point, the more numerous type is less fit (as in Figure 12.3) but is not an ESS if the more numerous type is more fit (as in Figure 12.5).]
 - Suppose there are important economies of scale in producing solar equipment, such that the cost savings increase the payoffs in the (Solar, Solar) cell of the table to (y, y) , where $y > 2$. How large would y need to be for the polymorphic ESS to have $s^* = 0.75$?
 - There are two types of racers—tortoises and hares—who race against each other in randomly drawn pairs. In this world, hares beat tortoises every time without fail. If two hares race, they tie, and they are completely exhausted by the race. When two tortoises race, they also tie, but they enjoy a pleasant conversation along the way. The payoff table is as follows (where $c > 0$):

		COLUMN	
		Tortoise	Hare
ROW	Tortoise	c, c	-1, 1
	Hare	1, -1	0, 0

-
- Assume that the proportion of tortoises in the population, t , is 0.5. For what values of c will tortoises have greater fitness than hares?

2. For what values of c will tortoises be fitter than hares if $t = 0.1$?
3. If $c = 1$, can a single hare successfully invade a population of tortoises? Explain why or why not.
4. In terms of t , how large must c be for tortoises to have greater fitness than hares?
5. In terms of c , what is the level of t in a polymorphic equilibrium? For what values of c will such an equilibrium exist? Explain.
8. Consider a population with two phenotypes, X and Y, with a payoff table as follows:

		COLUMN	
		X	Y
ROW	X	2, 2	5, 3
	Y	3, 5	1, 1

1. Find the fitness for X as a function of x , the proportion of X in the population, and the fitness for Y as a function of x .

Assume that the population dynamics from generation to generation conform to the following model:

$$x_{t+1} = \frac{x_t \times F_{Xt}}{x_t \times F_{Xt} + (1 - x_t) \times F_{Yt}},$$

where x_t is the proportion of X in the population in period t , x_{t+1} is the proportion of X in the population in period $t + 1$, F_{Xt} is the fitness of X in period t , and F_{Yt} is the fitness of Y in period t .

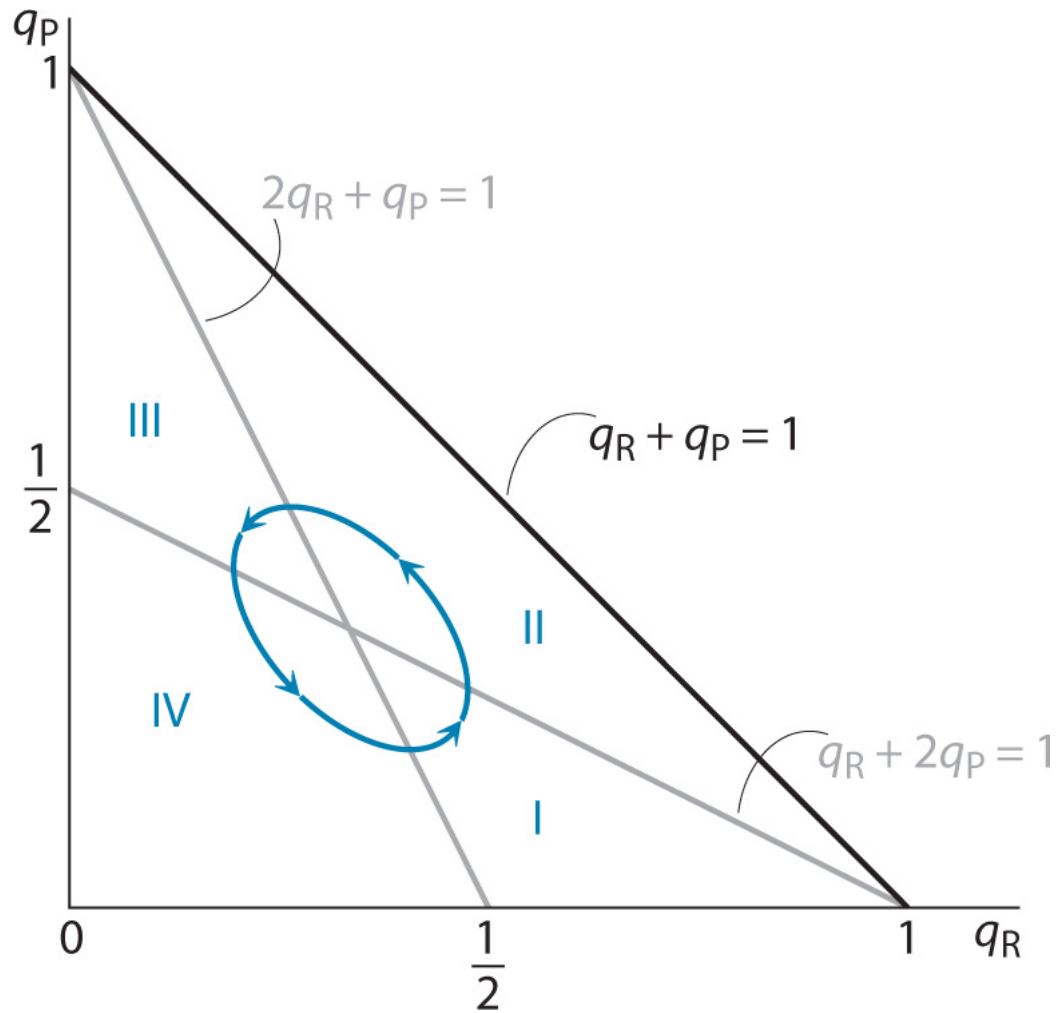
2. Assume that x_0 , the proportion of X in the population in period 0, is 0.2. What are F_{X0} and F_{Y0} ?
3. Find x_1 , using x_0 , F_{X0} , F_{Y0} , and the model given above.
4. What are F_{X1} and F_{Y1} ?
5. Find x_2 (rounded to five decimal places).
6. What are F_{X2} and F_{Y2} (rounded to five decimal places)?
9. Consider an evolutionary game between a green type and a purple type with a payoff table as follows:

		COLUMN	
		Green	Purple
ROW	Green	a , a	4, 3
	Purple	3, 4	2, 2

Let g be the proportion of greens in the population.

1. In terms of g , what is the fitness of the purple type?
 2. In terms of g and a , what is the fitness of the green type?
 3. Graph the fitness of the purple types against the fraction g of greens in the population. On the same diagram, show three lines for the fitnesses of the green type when $a = 2, 3$, and 4 . What can you conclude from this graph about the range of values of a that guarantees a stable polymorphic equilibrium?
 4. Assume that a is in the range found in part (c). In terms of a , what is the proportion of greens, g , in the stable polymorphic equilibrium?
10. At the World RPS Society, [21](#) people take the children's game of rock-paper-scissors (RPS) very seriously. You might think it's just a matter of chance who wins—Rock beating Scissors, Scissors beating Paper, Paper beating Rock, and a tie if both players make the same choice—but, for professional RPS players, “you can't just rely on luck” to win. For instance, men appear hardwired to start more often with Rock, while women are more likely to start with Paper. To maximize your chances of winning your next RPS throwdown, you should therefore start with Paper against a man or with Scissors against a woman. Of course, once enough people start following that advice, players will benefit by adapting again, and yet again. Here and in Exercise U10, you will examine the dynamics of that evolutionary dance in more detail.
1. Consider a population with three phenotypes: One always plays Rock (R type), one always plays Paper (P type), and one always plays Scissors (S type). Draw a three-by-three payoff table for a once-played RPS game between two friends, with R, P, and S as the three available phenotypes. Assign payoff 1 to any winning player, payoff -1 to any losing player, and payoff 0 to both players if there is a tie.
 2. Can there be a stable polymorphic configuration of the population with only R and P types? Use your answer to argue that, in any stable configuration, the population must include some of all three types.

Because $q_R + q_P + q_S = 1$, one can visualize all possible (q_R, q_P, q_S) configurations using a triangle as shown in the following figure. In the figure, each point in the triangle is a pair (q_R, q_P) with $q_R + q_P \leq 1$; the value for q_S is simply the difference between $q_R + q_P$ and 1 at each point. In this figure, we have highlighted (i) the point $(\frac{1}{3}, \frac{1}{3})$ corresponding to the mixed-strategy Nash equilibrium (found in Exercise S10 in [Chapter 7](#)), (ii) the lines consisting of all pairs (q_R, q_P) such that $2q_R + q_P = 1$ (line 1) and such that $q_R + 2q_P = 1$ (line 2), and (iii) the four regions of the triangle (labeled I, II, III, IV) created by these lines.



3. Verify that a player of type R gets a positive fitness whenever the population configuration corresponds to a point below line 1, and that a player of type P gets a positive fitness whenever the population configuration corresponds to a point to the right of line 2.

4. Suppose that the initial population configuration is in region I, and suppose that each type grows more numerous if it is getting a positive fitness, or grows less numerous if it is getting a negative fitness. Given these population dynamics, which of the following statements is true?
1. The population will remain in region I forever.
 2. The population will converge to the $(\frac{1}{3}, \frac{1}{3})$ configuration on a path that remains entirely inside region I.
 3. The population will transition to region II (moving counterclockwise).
 4. The population will transition to region IV (moving clockwise).

Explain your answer.

5. Prove the following statement: “If a strategy is strictly dominated in the payoff table of a game played by rational players, then in the evolutionary version of the same game it will die out, no matter what the initial population mix. If a strategy is weakly dominated, it may coexist with some other types, but not in a mixture of all types.”

UNSOLVED EXERCISES

- Consider a survival game in which members of a large population of animals meet in pairs and either fight over or share a food source. There are two phenotypes in the population: One always fights, and the other always shares. For the purposes of this question, assume that no other mutant types can arise in the population. Suppose that the value of the food source is 200 calories, and that caloric intake determines each player's fitness. If two sharing types meet one another, they each get half the food, but if a sharer meets a fighter, the sharer concedes immediately, and the fighter gets all the food.
 - Suppose that the cost of a fight is 50 calories (for each fighter) and that when two fighters meet, each is equally likely to win the fight and the food or to lose and get no food. Draw the payoff table for the game played between two random players from this population. Find all of the ESSs in the population. What type of game is being played in this case?
 - Now suppose that the cost of a fight is 150 calories for each fighter. Draw the new payoff table and find all of the ESSs for the population in this case. What type of game is being played here?
 - Using the notation of the hawk - dove game of [Section 12.4](#), indicate the values of V and C in parts (a) and (b), and confirm that your answers to those parts match the analysis presented in the chapter.
- Suppose that a single play of a prisoners' dilemma has the following payoffs:

		PLAYER 2	
		Cooperate	Defect
PLAYER 1	Cooperate	3, 3	1, 4
	Defect	4, 1	2, 2
You may need to scroll left and right to see the full figure.			

In a large population in which each member's behavior is genetically determined, each player will be either a defector (that is, always defects in any play of a prisoners' dilemma game) or a tit-for-tat player (in multiple rounds of a prisoners' dilemma, she cooperates on the first play, and on any subsequent play she does whatever her

opponent did on the preceding play). Pairs of randomly chosen players from this population will play rounds consisting of n single plays of this dilemma (where $n \geq 2$). The payoff to each player in one whole round (of n plays) is the sum of her payoffs in the n plays.

Let the population proportion of defectors be p and the proportion of tit-for-tat players be $(1 - p)$. Each member of the population plays rounds of the dilemma repeatedly, matched against a new, randomly chosen opponent for each new round. A tit-for-tat player always begins each new round by cooperating on the first play.

1. Show in a two-by-two table the payoffs to a player of each type when, in one round of plays, each player meets an opponent of each of the two types.
 2. Find the fitness (average payoff in one round against a randomly chosen opponent) for a defector.
 3. Find the fitness for a tit-for-tat player.
 4. Use the answers to parts (b) and (c) to show that, when $p > (n - 2)/(n - 1)$, the defector type has greater fitness and that, when $p < (n - 2)/(n - 1)$, the tit-for-tat type has greater fitness.
 5. If evolution leads to a gradual increase in the proportion of the fitter type in the population, what are the possible eventual equilibrium outcomes of this process for the population described in this exercise? (That is, what are the possible equilibria, and which are evolutionarily stable?) Use a diagram with the fitness graphs to illustrate your answer.
 6. In what sense does more repetition (larger values of n) facilitate the evolution of cooperation?
3. Consider the thrice-repeated restaurant pricing prisoners' dilemma we studied in [Section 3.B](#), but now suppose that there are three phenotypes in the population: type A, which always defects; type T, which plays tit-for-tat; and type S, which never defects on the first play but always defects on the second play of each round of two successive plays against the same opponent.
1. Draw the three-by-three fitness table for the game, similar to the fitness table that we provided in Exercise S3.

For each of parts (b) - (d), be specific and explicit in your answers; make sure you use payoff numbers from the fitness table you draw to answer part (a).

2. Can a population that is 100% type A be successfully invaded by type S mutants?
3. Can a population that is 100% type T be successfully invaded by type S mutants?

4. Is the newly conceived type S strategy an ESS of this game?
4. Following the pattern of Exercise S4, analyze an evolutionary version of the tennis point game (see Figure 4.17). Regard servers and receivers as separate species, and construct a figure like Figure 12.7. What can you say about the ESS and its dynamics?
5. Recall from Exercise U1 the population of animals fighting over a food source worth 200 calories. Assume that, as in part (b) of that exercise, the cost of a fight is 150 calories per fighter. Assume also that a third phenotype exists in the population. That phenotype is a mixer; it plays a mixed strategy, sometimes fighting and sometimes sharing.
 1. Use your knowledge of mixed strategies in rationally played games to posit a reasonable mixture for the mixer phenotype to use in this game.
 2. Draw the three-by-three payoff table for this game when the mixer phenotype uses the mixed strategy that you found in part (a).
 3. Determine whether the mixer strategy is an ESS of this game. (Hint: Test whether a mixer population can be invaded successfully by either the fighting type or the sharing type.)
6. Consider an evolutionary version of the game between Baker and Cutler, from Exercise U1 of [Chapter 10](#). This time, Baker and Cutler are not two individuals, but two separate species. Each time a Baker meets a Cutler, they play the following game. The Baker chooses the total prize to be either \$10 or \$100. The Cutler chooses how to divide the prize chosen by the Baker: The Cutler can choose either a 50:50 split or a 90:10 split in the Cutler's own favor. The Cutler moves first, and the Baker moves second.

There are two types of Cutlers in the population: Type F chooses a fair (50:50) split, whereas type G chooses a greedy (90:10) split. There are also two types of Bakers: Type S simply chooses the large prize (\$100) no matter what the Cutler has done, whereas type T chooses the large prize (\$100) if the Cutler chooses a 50:50 split, but the small prize (\$10) if the Cutler chooses a 90:10 split.

Let f be the proportion of type F in the Cutler population, so that $(1 - f)$ represents the proportion of type G. Let s be the proportion of type S in the Baker population, so that $(1 - s)$ represents the proportion of type T.

1. Find the fitnesses of the Cutler types F and G in terms of s .
2. Find the fitnesses of the Baker types S and T in terms of f .
3. For what value of s are types F and G equally fit?
4. For what value of f are types S and T equally fit?

5. Use the answers above to sketch a graph displaying the population dynamics. Assign f as the horizontal axis and s as the vertical axis.
6. Describe all equilibria of this evolutionary game and indicate which ones are stable.
7. Recall Exercise S7. Hares, it turns out, are very impolite winners. Whenever hares race tortoises, they mercilessly mock their slow-footed (and easily defeated) rivals. The poor tortoises leave the race not only in defeat, but with their tender feelings crushed by the oblivious hares. The payoff table is thus

		COLUMN	
		Tortoise	Hare
ROW	Tortoise	c, c	$-2, 1$
	Hare	$1, -2$	$0, 0$

-
1. For what values of c are tortoises fitter than hares if t , the proportion of tortoises in the population, is 0.5? How does this compare with your answer in Exercise S7, part (a)?
 2. For what values of c are tortoises fitter than hares if $t = 0.1$? How does this compare with your answer in Exercise S7, part (b)?
 3. If $c = 1$, will a single hare successfully invade a population of tortoises? Explain why or why not.
 4. In terms of t , how large must c be for tortoises to be fitter than hares?
 5. In terms of c , what is the value of t in a polymorphic equilibrium? For what values of c will such an equilibrium exist? Explain.
 6. Will the polymorphic equilibria found to exist in part (e) be stable? Why or why not?
 8. (Use of spreadsheet software recommended) This problem explores more thoroughly the generation-by-generation population dynamics seen in Exercise S8. Since the math can quickly become very complicated and tedious, it is much easier to do this analysis with the aid of a spreadsheet.

Again, consider a population with two types, X and Y, with a payoff table as follows:

		COLUMN	
		X	Y
ROW	X	$2, 2$	$5, 3$
	Y		

	COLUMN		
	X		Y
	Y	3, 5	1, 1

Recall that the population dynamics from generation to generation are given by

$$x_{t+1} = \frac{x_t \times F_{Xt}}{x_t \times F_{Xt} + (1 - x_t) \times F_{Yt}},$$

where x_t is the proportion of X in the population in period t , x_{t+1} is the proportion of X in the population in period $t + 1$, F_{Xt} is the fitness of X in period t , and F_{Yt} is the fitness of Y in period t .

Use a spreadsheet to extend these calculations to many generations. [Hint: Assign three horizontally adjacent cells to hold the values of x_t , F_{Xt} , and F_{Yt} , and have each successive row represent a different period ($t = 0, 1, 2, 3, \dots$). Use spreadsheet formulas to relate F_{Xt} and F_{Yt} to x_t and x_{t+1} to x_t , F_{Xt} , and F_{Yt} according to the population model given above.]

1. If there are initially equal proportions of X and Y in the population in period 1 (that is, if $x_0 = 0.5$), what is the proportion of X in the next generation, x_1 ? What are F_{X1} and F_{Y1} ?
2. Use a spreadsheet to extend these calculations to the next generation, and the next, and so on. To four decimal places, what is the value of x_{20} ? What are F_{X20} and F_{Y20} ?
3. What is x^* , the equilibrium level of x ? How many generations does it take for the population to be within 1% of x^* ?
4. Answer the questions in part (b), but with a starting value of $x_0 = 0.1$.
5. Repeat part (b), but with $x_0 = 1$.
6. Repeat part (b), but with $x_0 = 0.99$.
7. Are monomorphic equilibria possible in this model? If so, are they stable? Explain.
9. Consider an evolutionary game between a green type and a purple type, with a payoff table as follows:

COLUMN

		Green	COLUMN Purple
		Green	Purple
ROW	Green	a, a	b, c
	Purple	c, b	d, d

In terms of the parameters a , b , c , and d , find the conditions that will guarantee a stable polymorphic equilibrium.

10. (Optional, for mathematically trained students) Common side-blotched lizards (*Uta stansburiana*), residents of the California desert, play an unusual evolutionary game. During mating season, males display colorful throat patches in three shades—orange, blue, and yellow—with each color corresponding to a different mating strategy. Orange-throats (type O) tend to be somewhat larger, are aggressive, and can defend large breeding territories that host two or more females. Yellow-throats (type Y) are smaller, generally unable to defend much territory at all, and mate mainly by sneaking into other males' territories. Finally, blue-throats (type B) are monogamous, sticking close to a single mate. These differences in color and in behavior, all seemingly determined by a single gene (!), result in orange-throats beating blue-throats in competition over females, but often being bested by “sneaky” yellow-throats, who in turn can be defeated by blue-throats, who guard their mates more carefully. [22](#)

The game played among these lizard phenotypes thus has essentially the same payoff structure as rock-paper-scissors (RPS) in Exercise S10, with type O as Rock, type Y as Paper, and type B as Scissors. The actual reproductive payoffs of these phenotypes are not exactly the same as in RPS, but we suppose for simplicity that they are—with payoffs of +1, 0, or −1, depending on which phenotypes are matched to play the game—and that each phenotype grows in number over time if its fitness is positive, or declines in number if its fitness is negative.

- Let q_Y , q_B denote the proportion of lizards in the population that are yellow-throats and blue-throats, respectively, and let $q_0 = 1 - q_Y - q_B$ be the remaining proportion of orange-throats. Express the fitnesses of phenotypes Y, B, and O as functions of q_Y , q_B , and q_0 , respectively. Verify that (i) type Y grows in number if and only if $q_0 > q_B$, (ii) type B grows in number if and only if $q_Y > q_0$, and (iii) type O grows in number if and only if $q_B > q_Y$.

2. Consider the dynamics of this evolutionary system more explicitly. Let the speed of change in a variable x at time t be denoted by the derivative dx/dt . Now consider the following expressions:

$$\frac{dq_Y}{dt} = q_O - q_B, \quad \frac{dq_B}{dt} = q_Y - q_O, \quad \text{and} \quad \frac{dq_O}{dt} = q_B - q_Y.$$

Verify that these derivatives conform to the findings of part (a).

3. Define $X = (q_Y)^2 + (q_B)^2 + (q_O)^2$. Using the chain rule of differentiation, show that $dX/dt = 0$; that is, show that X remains constant over time.
4. From the definitions of these variables, we know that $q_Y + q_B + q_O = 1$. Combining this fact with the result from part (c), show that over time, in three-dimensional space, the point (q_Y, q_B, q_O) moves along a circle.
5. What does the answer to part (d) indicate regarding the stability of the evolutionary dynamics in the side-blotched lizard population? Will the population ever converge to a stable mixture of the three phenotypes?

Endnotes

- The World RPS Society hosts the website www.wrpsa.com. The “facts” presented in this exercise are discussed in that site’s “Rock Paper Scissors Beginner Strategies,” available at <https://www.wrpsa.com/rock-paper-scissors-beginner-strategies> (accessed May 1, 2019). [Return to reference 21](#)
- For more information about the side-blotched lizards, see Kelly Zamudio and Barry Sinervo, “Polygyny, Mate-Guarding, and Posthumous Fertilizations as Alternative Mating Strategies,” *Proceedings of the National Academy of Sciences*, vol. 97, no. 26 (December 19, 2000), pp. 14427 – 32. [Return to reference 22](#)

PART FOUR



Applications to Specific Strategic Situations

13 ■ Brinkmanship: The Cuban Missile Crisis

IN [CHAPTER 1](#), we explained that our basic approach to strategic games was neither pure theory nor pure case study, but a combination in which theoretical ideas would be developed by using features of particular cases or examples. Thus, we ignored those aspects of each case that were incidental to the concept being developed. However, after you have learned the theoretical ideas, a richer mode of analysis becomes available to you in which factual details of a particular case are more closely integrated with game-theoretic analysis to achieve a fuller understanding of what has happened and why. Such *theory-based case studies* have begun to appear in diverse fields—business, political science, and economic history.¹

In this chapter, we offer an case from political and military history—namely, the Cuban missile crisis of 1962. Our choice is motivated by the sheer drama of the episode, the wealth of factual information that has become available, and the applicability of the important concept of brinkmanship from game theory.

The crisis, when the world came as close to a nuclear war as it ever has, is often offered as the classic example of brinkmanship. You may think that the risk of nuclear war ended with the dissolution of the Soviet Union, making this case a historical curiosity. But nuclear arms races continue in many parts of the world, and such rivals as India and Pakistan, or Iran and Israel, may find use for the lessons taken from the Cuban missile crisis. Even major-power rivalries, including those between the United States and China, or between the United States and Russia, may heat up. More importantly for many of you, brinkmanship must be

practiced in many more common situations, from political negotiations to business-labor relations to marital disputes. Although the stakes in such games are lower than those in a nuclear confrontation between superpowers, the same principles of strategy apply.

In [Chapter 8](#), we introduced the concept of brinkmanship as a strategic move; here is a quick reminder of that analysis. A *threat* is a response rule, and the threatened action inflicts a cost on both the player making the threat and the player whose action the threat is intended to influence. However, if the threat succeeds in its purpose, this action is not actually carried out. Therefore, there is no apparent upper limit to the cost of the threatened action. But the risk of *errors*—that is, the risk that the threat may fail to achieve its purpose, or that the threatened action may occur by accident—forces the strategist to use the minimal threat that achieves its purpose. If a smaller threat is not naturally available, you can scale down a large threat by making its fulfillment probabilistic. You can do something in advance that creates a probability, but not a certainty, that the mutually harmful outcome will happen if the opponent defies you. If the need actually arose, you would not take that bad action if you had the full freedom to choose. Therefore, you must arrange in advance to let things get out of your control to some extent. [Brinkmanship](#) is the creation and deployment of such a probabilistic threat; it consists of a deliberate loss of control.

In our extended case study of the Cuban missile crisis, we will explain the concept of brinkmanship in detail. In the process, we will find that many popular interpretations and analyses of the crisis are simplistic. A deeper analysis reveals brinkmanship to be a subtle and dangerous strategy. It also shows that many detrimental outcomes in business and personal interactions—such as strikes and breakups of relationships—are examples of brinkmanship gone wrong.

Therefore, a clear understanding of the strategy, as well as its limitations and risks, is very important to all game players, which includes just about everyone.

Endnotes

- Two excellent examples of theory-based case studies are Pankaj Ghemawat, *Games Businesses Play: Cases and Models* (Cambridge, Mass.: MIT Press, 1997), and Robert H. Bates, Avner Greif, Margaret Levi, Jean-Laurent Rosenthal, and Barry Weingast, *Analytic Narratives* (Princeton, N.J.: Princeton University Press, 1998). A broader analysis of the approach can be found in Alexander L. George and Andrew Bennett, *Case Studies and Theory Development in the Social Sciences* (Cambridge, Mass.: MIT Press, 2005).
[Return to reference 1](#)

Glossary

[brinkmanship](#)

A threat that creates a risk but not certainty of a mutually bad outcome if the other player defies your specified wish as to how he should act, and then gradually increases this risk until one player gives in or the bad outcome happens.

1 A BRIEF NARRATIVE OF EVENTS

We begin with a brief story of the unfolding of the crisis. Our account draws on several books, including some that were written with the benefit of documents and statements released since the collapse of the Soviet Union.² We cannot hope to do justice to the detail, let alone the drama, of the events. We urge you to read the books that tell the story in vivid detail and to talk to any relatives who lived through it to get their firsthand memories.³

In late summer and early fall of 1962, at the height of the Cold War, the Soviet Union (USSR) started to place medium- and intermediate-range ballistic missiles (MRBMs and IRBMs) armed with nuclear weapons in Cuba. The MRBMs had a range of 1,100 miles and could hit Washington, D.C.; the IRBMs, with a range of 2,200 miles, could hit most major U.S. cities and military installations. The missile sites were guarded by the latest Soviet SA-2 surface-to-air missiles (SAMs), which could shoot down U.S. high-altitude U-2 reconnaissance planes. They were also defended by IL-28 bombers and tactical nuclear weapons called Luna by the Soviets and FROGs (free rockets over ground) by the United States, which could be used against invading troops.

This was the first time that the Soviets had ever attempted to place their missiles and nuclear weapons outside Soviet territory. Had they been successful, it would have increased their offensive capability against the United States manyfold. It is now believed that the Soviets had very few (U.S. aerial reconnaissance showed only four) operational intercontinental ballistic missiles (ICBMs) in their own country capable of reaching the United States (*War*, 464, 509-10; *Doomsday*, 158). Their initial installation in Cuba had about 40 MRBMs and IRBMs, which was a substantial

increase. But the United States would still have retained vast superiority in the nuclear balance between the superpowers. Also, as the Soviets built up their fleet of nuclear-armed submarines, the relative importance of land-based missiles near the United States would have decreased. But the missiles had more than direct military value to the Soviets. Successful placement of missiles so close to the United States would have been an immense boost to Soviet prestige throughout the world, especially in Asia and Africa, where the superpowers were competing for political and military influence. Finally, the Soviets had come to think of Cuba as a poster child for socialism. The opportunity to deter a feared U.S. invasion of Cuba, and to counter Chinese influence in Cuba, weighed importantly in the calculations of the Soviet leader and premier, Nikita Khrushchev. (See *Gamble*, 182-83, for an analysis of Soviet motives.)

The whole operation was attempted in utmost secrecy, and the Soviets hoped to conceal the missiles under palm trees (*Gamble*, Chapter 10)! But this attempt did not work. U.S. surveillance of Cuba and of shipping lanes during the late summer and early fall of 1962 had indicated some suspicious activity. When questioned about it by U.S. diplomats, the Soviets denied any intentions to place missiles in Cuba. Later, faced with irrefutable evidence, they said that their intention was defensive, to deter the United States from invading Cuba. It is hard to believe this, although we know that an offensive weapon *can* serve as a defensive deterrent threat.

On Sunday and Monday, October 14 and 15, an American U-2 spy plane took photographs over western Cuba that showed unmistakable signs of construction on MRBM launching sites. (Evidence of IRBMs was found later, on October 17.) They were shown to U.S. President John F. Kennedy the following day (October 16). He immediately convened an ad hoc group of top-level advisers, later called the Executive Committee of the

National Security Council (ExComm), to discuss the available options.⁴ He decided to keep the matter totally secret until he was ready to act, mainly because if the Soviets knew that the Americans knew about the missiles, they might speed up their installation and deployment before the Americans were ready to act, but also because spreading the news without announcing a clear response would create panic in the United States. During the rest of that week (October 16-21), ExComm met numerous times. To preserve secrecy, the president continued his normal schedule, including travel to speak for Democratic candidates in the upcoming midterm congressional elections. He kept in constant touch with ExComm. He dodged press questions about Cuba and persuaded one or two trusted media owners or editors to preserve the facade of business as usual.

Different members of ExComm had widely differing assessments of the situation and supported different actions. The military Joint Chiefs of Staff thought that the missile placement changed the balance of military power substantially; Defense Secretary Robert McNamara thought it had changed "not at all," but regarded the problem as politically important nonetheless (*Tapes*, 89). Kennedy pointed out that the first placement, if ignored by the United States, could grow into something much bigger, and that the Soviets could use the threat of missiles so close to the United States to try to force the withdrawal of U.S., British, and French forces from West Berlin. Kennedy was also aware that the placement of the missiles was a part of the *geopolitical* struggle between the United States and the Soviet Union (*Tapes*, 92).

It now appears that Kennedy was very much on the mark in this assessment. The Soviets planned to expand their presence in Cuba into a major military base (*Tapes*, 677). They expected to complete the missile placement by mid-November. Khrushchev had planned to sign a treaty with Cuba's prime minister,

Fidel Castro, in late November, then travel to New York to address the United Nations and issue an ultimatum for a settlement of the Berlin issue (*Tapes*, 679; *Gamble*, 182), using the missiles in Cuba as a threat for this purpose. Khrushchev thought Kennedy would accept the missile placement as a *fait accompli*. Khrushchev appears to have made these plans on his own. Some of his top advisers privately thought them too adventurous, but the top governmental decision-making body of the Soviet Union, the Presidium, supported him, although it acted largely as a rubber stamp for the premier's decisions (*Gamble*, 180). Castro was at first reluctant to accept the missiles, fearing that they would trigger a U.S. invasion (*Tapes*, 676 - 78), but in the end he, too, accepted them. The prospect gave him great confidence and lent some swagger to his statements about the United States (*Gamble*, 186 - 87, 229 - 30).

In all ExComm meetings up to and including the one on the morning of Thursday, October 18, all the participants appear to have assumed that the U.S. response would be purely military. The only options that they discussed seriously during this time were (1) an air strike directed exclusively at the missile sites and (probably) the SAM sites nearby, (2) a wider air strike including Soviet and Cuban aircraft parked at airfields, and (3) a full-scale invasion of Cuba. If anything, attitudes hardened when the evidence of the presence of the longer-range IRBMs arrived. In fact, at the Thursday meeting, Kennedy discussed a timetable for air strikes to commence that weekend (*Tapes*, 148).

McNamara had first mentioned a blockade of Cuba toward the end of the meeting on Tuesday, October 16, and he developed the idea (in a form uncannily close to the course of action actually taken) in a small group after the formal meeting had ended (*Tapes*, 86, 113). George Ball argued that an air strike without warning would be a "Pearl Harbor" and that the United States should not do it (*Tapes*, 115); he got important

support for this view from Robert Kennedy (*Tapes*, 149). The civilian members of ExComm further shifted toward the blockade option when they found that what the military Joint Chiefs of Staff wanted was a massive air strike; the military regarded a limited strike aimed only at the missile sites as so dangerous and ineffective that “they would prefer taking no military action than to take that limited strike” (*Tapes*, 97).

Between October 18 and Saturday, October 20, the majority opinion within ExComm gradually coalesced around the idea of starting with a blockade, simultaneously issuing an ultimatum with a short deadline (periods from 48 to 72 hours were mentioned), and proceeding to military action if necessary after this deadline expired. International law required a declaration of war to set up a blockade, but this problem was ingeniously resolved by calling it a “naval quarantine” of Cuba (*Tapes*, 190 – 96).

Some people held the same positions throughout these ExComm discussions (October 16 – 21)—for example, the military Joint Chiefs of Staff constantly favored a major air strike—but others shifted their views, at times dramatically. National Security Adviser McGeorge Bundy initially favored doing nothing (*Tapes*, 172) and then switched toward a preemptive surprise air attack (*Tapes*, 189). President Kennedy’s own position also shifted away from an air strike toward a blockade. He wanted the U.S. response to be firm. Although his reasons undoubtedly were mainly military and geopolitical, as a good politician he was also fully aware that a weak response would hurt the Democratic Party in the imminent congressional elections. On the other hand, the responsibility of starting an action that might lead to nuclear war weighed very heavily on him. He was impressed by the CIA’s assessment that some of the missiles were already operational, which increased the risk that any air strike or invasion could lead the Soviets to fire those missiles and

produce large U.S. civilian casualties (*Gamble*, 235). In the second week of the crisis (October 22–28), his decisions seemed constantly to favor the lowest-key options discussed by ExComm.

By the end of the first week's discussions, the choice lay between a blockade and an air strike. In a straw vote on October 20, the blockade won 11 to 6 (*War*, 516). Kennedy made the decision to impose a blockade and announced it in a television address to the nation on Monday, October 22. He demanded a halt to the shipment of Soviet missiles to Cuba and a prompt withdrawal of those already there.

Kennedy's speech brought all of the drama and tension of the crisis into the public arena. The United Nations held several dramatic but unproductive debates. Other world leaders and the usual policy wonks of international affairs offered advice and mediation.

Between October 23 and October 25, the Soviets at first tried bluster and denial; Khrushchev called the blockade "banditry, a folly of international imperialism," and said that his ships would ignore it. The Soviets, in the United Nations and elsewhere, claimed that their intentions were purely defensive and issued statements of defiance. In secret, they explored ways to end the crisis. This exploration included some direct messages from Khrushchev to Kennedy. It also included some very indirect and lower-level approaches by the Soviets. In fact, as early as Monday, October 22—before Kennedy's TV address—the Soviet Presidium had decided not to let this crisis lead to war. By Thursday, October 25, its members had decided that they were willing to withdraw from Cuba in exchange for a promise by the United States not to invade Cuba, but they had also agreed to "look around" for better deals (*Gamble*, 241, 259). The United States was not aware of any of this Soviet thinking.

In public as well as in private communications, the USSR suggested a swap: withdrawal of U.S. missiles from Turkey and of Soviet ones from Cuba. This possibility had already been discussed by ExComm. The missiles in Turkey were obsolete; the United States wanted to remove them anyway and replace them with a Polaris submarine stationed in the Mediterranean Sea. But it was thought that Turkey would regard the presence of U.S. missiles as a matter of prestige and that it might be difficult to persuade it to accept the change. (Turkey might also correctly regard missiles, fixed on Turkish soil, as a firmer signal of the U.S. commitment to its defense than an offshore submarine, which could move away on short notice; see *Tapes*, 568.)

The blockade went into effect on Wednesday, October 24. Despite their public bluster, the Soviets were cautious in testing it. Apparently, they were surprised that the United States had discovered the missiles in Cuba before the whole installation program was completed; Soviet personnel in Cuba had observed the U-2 overflights but had not reported them to Moscow (*Tapes*, 681). The Soviet Presidium ordered the ships carrying the most sensitive materials (actually, the IRBM missiles) to stop or turn around. But it also ordered General Issa Pliyev, the commander of the Soviet troops in Cuba, to get his troops combat-ready and to use all means except nuclear weapons to meet any attack (*Tapes*, 682). In fact, the Presidium twice prepared (then canceled without sending) orders authorizing him to use tactical nuclear weapons in the event of a U.S. invasion (*Gamble*, 242 - 43, 272, 276). The U.S. side saw only that several Soviet ships (which were actually carrying oil and other nonmilitary cargo) continued to sail toward the blockade zone. The U.S. Navy showed some moderation in its enforcement of the blockade; a tanker was allowed to pass without being boarded, and the tramp steamer *Marucla*, carrying industrial cargo, was boarded but allowed to proceed after only a cursory inspection. But tension was

mounting, and neither side's actions were as cautious as the top-level politicians on both sides would have liked.

On the morning of Friday, October 26, Khrushchev sent Kennedy a conciliatory private letter offering to withdraw the missiles in exchange for a U.S. promise not to invade Cuba. But later that day he toughened his stance. It seems that he was emboldened by two items of evidence. First, he saw that the U.S. Navy was not being excessively aggressive in enforcing the blockade. Second, some dovish statements had appeared in U.S. newspapers. Most notable among them was an article by the influential and well-connected syndicated columnist Walter Lippman, who suggested the swap whereby the United States would withdraw its missiles in Turkey in exchange for the USSR's withdrawing its missiles in Cuba (*Gamble*, 275). Khrushchev sent another letter to Kennedy on Saturday, October 27, offering this swap, and this time he made the letter public. The new letter was presumably a part of the Presidium's strategy of "looking around" for the best deal. Members of ExComm concluded that the first letter expressed Khrushchev's own thoughts, but that the second was written under pressure from hard-liners in the Presidium—or was even evidence that Khrushchev was no longer in control (*Tapes*, 498, 512–13). In fact, both of Khrushchev's letters were discussed and approved by the Presidium (*Gamble*, 263, 275).

ExComm continued to meet, and opinions within it hardened. There was a growing feeling that the blockade by itself would not work. Kennedy's TV speech had imposed no firm deadline; in the absence of a deadline, a compellent threat is vulnerable to the opponent's procrastination.⁵ Kennedy had seen this quite clearly, and as early as Monday, October 22, he commented, "I don't think we're gonna be better off if they're just sitting there" (*Tapes*, 216). But a hard, short deadline was presumably thought to be too rigid. By Thursday, others in ExComm were realizing the problem; Bundy, for

example, said, "A plateau here is the most dangerous thing" (*Tapes*, 423). The hardening of the Soviet position, as shown by the public "Saturday letter" that followed the conciliatory private "Friday letter," was another concern. More ominously, that Friday, U.S. surveillance had discovered the presence of tactical nuclear weapons (FROGs) in Cuba (*Tapes*, 475). This discovery showed the Soviet presence there to be vastly greater than thought before, but it also made an invasion more dangerous to U.S. troops. Also on Saturday, a U.S. U-2 plane was shot down over Cuba, and Cuban anti-aircraft defenses fired at lower-flying U.S. reconnaissance planes. The grim mood in ExComm throughout that Saturday was well encapsulated by Douglas Dillon: "We haven't got but one more day" (*Tapes*, 534).

On that Saturday, U.S. plans leading to escalation were being put in place. An air strike was planned for the following Monday, or Tuesday at the latest, and Air Force reserves were called up (*Tapes*, 612-13). Invasion was seen as the inevitable culmination of events (*Tapes*, 537-38). A tough private letter to Khrushchev from President Kennedy was drafted and was handed over by Robert Kennedy to the Soviet ambassador in Washington, Anatoly Dobrynin. In it, Kennedy made the following offer: (1) The Soviet Union withdraws its missiles and IL-28 bombers from Cuba with adequate verification (and ships no new ones). (2) The United States promises not to invade Cuba. (3) The U.S. missiles in Turkey will be removed after a few months, but this offer is void if the Soviets mention it in public or link it to the Cuban deal. An answer was required within 12 to 24 hours; otherwise "there would be drastic consequences" (*Tapes*, 605-7).

On the morning of Sunday, October 28, just as prayers and sermons for peace were being offered in many churches in the United States, Soviet radio broadcast the text of a letter that Khrushchev was sending to Kennedy, in which he announced that construction of the Cuban missile sites was being halted

immediately and that the missiles already installed would be dismantled and shipped back to the Soviet Union. Kennedy immediately sent a reply welcoming this decision, which was broadcast to Moscow by Voice of America radio. It now appears that Khrushchev's decision to back down was made before he received Kennedy's letter through Dobrynin, but that the letter only reinforced it (*Tapes*, 689).

That did not quite end the crisis. The U.S. Joint Chiefs of Staff remained skeptical of the Soviets and wanted to go ahead with their air strike (*Tapes*, 635). In fact, construction activity at the Cuban missile sites continued for a few days. Verification of the missiles' withdrawal by the United Nations proved problematic. The Soviets tried to make the Turkey part of the deal semipublic. They also tried to keep the IL-28 bombers in Cuba out of the withdrawal. Not until November 20 was the deal finally clinched and the withdrawal begun (*Tapes*, 663 - 65; *Gamble*, 298 - 310).

Endnotes

- Our sources (which we cite in the text using the word underlined here for each book, followed by the relevant page numbers) include Robert Smith Thompson, *The Missiles of October* (New York: Simon & Schuster, 1992); James G. Blight and David A. Welch, *On the Brink : Americans and Soviets Reexamine the Cuban Missile Crisis* (New York: Hill and Wang, 1989); Richard Reeves, *President Kennedy: Profile of Power* (New York: Simon & Schuster, 1993); Donald Kagan, *On the Origins of War and the Preservation of Peace* (New York: Doubleday, 1995); Aleksandr Fursenko and Timothy Naftali, *One Hell of a Gamble : The Secret History of the Cuban Missile Crisis* (New York: W. W. Norton, 1997); *The Kennedy Tapes : Inside the White House during the Cuban Missile Crisis*, ed. Ernest R. May and Philip D. Zelikow (Cambridge, Mass.: Harvard University Press, 1997); Michael Dobbs, *One Minute to Midnight : Kennedy, Khrushchev and Castro on the Brink of Nuclear War* (New York: Knopf, 2008); and Daniel Ellsberg, *The Doomsday Machine: Confessions of a Nuclear War Planner* (New York: Bloomsbury Publishing, 2017). Graham T. Allison's *Essence of Decision: Explaining the Cuban Missile Crisis* (Boston: Little Brown, 1971) remains important not only for its narrative, but also for its analysis and interpretation. Our view differs from his in some important respects, but we remain in debt to his insights. We follow and extend the ideas in Avinash Dixit and Barry Nalebuff, *Thinking Strategically* (New York: W. W. Norton, 1991), [Chapter 8](#). [Return to reference 2](#)
- For those of you with no access to firsthand information, or those who seek a beginner's introduction to both the details and the drama of the missile crisis, we recommend the film *Thirteen Days* (2000, New Line Cinema). A relatively short book by Sheldon Stern uses the evidence

from the Kennedy administration tapes to present as accurate a view of the crisis and its later analysis as possible. His book is perhaps the best short read for interested parties. See Sheldon Stern, *The Cuban Missile Crisis in American Memory: Myths versus Reality* (Stanford, Calif.: Stanford University Press, 2012).

[Return to reference 3](#)

- Members of ExComm who figured most prominently in the discussions were Secretary of Defense Robert McNamara; National Security Adviser McGeorge Bundy; the chairman of the Joint Chiefs of Staff, General Maxwell Taylor; Secretary of State Dean Rusk and Undersecretary George Ball; Attorney General Robert Kennedy (who was also the president's brother); Secretary of the Treasury Douglas Dillon (also the only Republican in the cabinet); and Llewellyn Thompson, who had recently returned from being U.S. ambassador in Moscow. During the two weeks that followed, they would be joined by or would consult with several others, including the U.S. ambassador to the United Nations, Adlai Stevenson; Dean Acheson, former secretary of state and a senior statesman of U.S. foreign policy; and the chief of the U.S. Air Force, General Curtis LeMay. [Return to reference 4](#)
- See the discussion of salami tactics in Chapter 8, section 6.E. [Return to reference 5](#)

2 A SIMPLE GAME-THEORETIC EXPLANATION

At first sight, the game-theoretic aspect of the Cuban missile crisis looks very simple. The United States wanted the Soviet Union to withdraw its missiles from Cuba; thus the U.S. objective was to achieve compellence. For this purpose, the United States deployed a threat: Soviet failure to comply would lead to a nuclear war. This was sufficiently frightening to Khrushchev that he complied. The prospect of nuclear annihilation was equally frightening to Kennedy, but that is in the nature of a threat. All that is needed is that the threat be sufficiently costly to the other side to induce it to act in accordance with our wishes; then we don't have to carry out the bad action anyway.

A somewhat more formal statement of this argument proceeds by drawing a game tree like that shown in Figure 13.1. The Soviets have installed the missiles, and now the United States has the first move. It chooses between doing nothing and issuing a threat. If the United States does nothing, this is a military and political achievement for the Soviets; so we score the payoffs as -1 for the United States and 1 for the Soviets. If the United States issues its threat, the Soviets get to move, and they can either withdraw or defy. Withdrawal is a humiliation for the Soviets and a reaffirmation of U.S. military superiority, so we score it 1 for the United States and -1 for the Soviets. If the Soviets defy the U.S. threat, there will be a nuclear war. This outcome is terrible for both, so we score this -10 for each. The numbers are (intentionally) chosen to be same as those in the game of chicken (see Figure 4.16) except that the disaster payoff is much worse. The conclusions do not depend on the precise numbers that we have chosen, however. If you disagree with our choice, you can substitute other numbers you think to be a more accurate representation; as long as the *relative* ranking of the outcomes is the same, you will get the same subgame-perfect equilibrium.

Now we can easily find that equilibrium. If faced with the U.S. threat, the Soviets get -1 from withdrawal and -10 by defiance; so they prefer to withdraw. Looking ahead to this outcome, the United States reckons on getting 1 if it issues the threat and -1 if it does not; therefore it is optimal for the United States to make the threat. The outcome gives payoffs of 1 to the United States and -1 to the Soviets.

But a moment's further thought shows this interpretation to be unsatisfactory. One might start by asking why the Soviets would deploy the missiles in Cuba at all, when they could look ahead to this unfolding of the subsequent game in which they would come out the losers. But even more importantly, several facts about the situation and several events in the course of its unfolding do not fit into this picture of a simple threat.

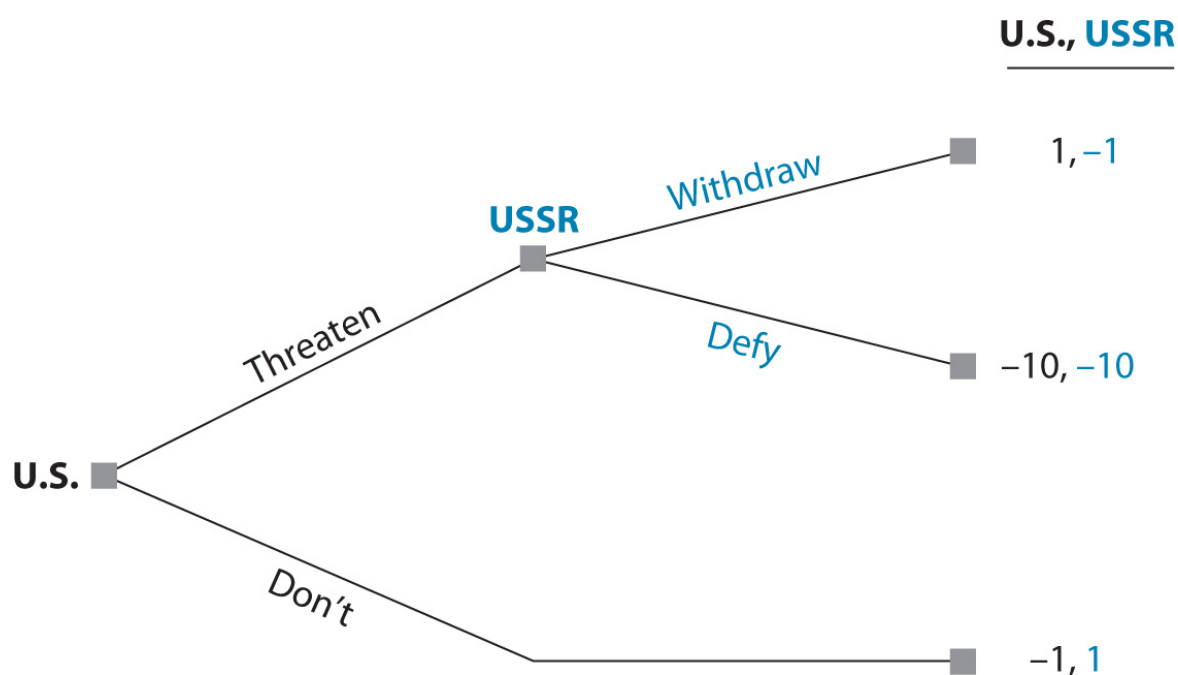


FIGURE 13.1 The Simple-Threat Model of the Crisis

So, let us develop a more satisfactory game-theoretic argument. As we pointed out before, the idea that a threat has only a lower limit on its size—namely, that it be large enough to frighten the opponent—is correct only if the threatener can be absolutely sure that everything will go as planned. But almost all games have some element of uncertainty. You cannot be certain about

your opponent's value system, and you cannot be completely sure that any player's intended actions will be accurately implemented. Therefore, a threat carries a twofold risk. Your opponent may defy it, requiring you to carry out the costly threatened action; or your opponent may comply, but the threatened action may occur by mistake anyway. When such risks exist, the cost of a threatened action to oneself becomes an important consideration.

The Cuban missile crisis was replete with such uncertainties. Neither side could be sure of the other's payoffs—that is, of its relative valuations of winning or losing the Cold War nor the costs of a hot war. Also, the choices of “blockade” and “air strike” were much more complex than those simple phrases suggest, and there were many weak links and random effects between an order given in Washington or Moscow and its implementation in the Atlantic Ocean or in Cuba.

Graham Allison's excellent book *Essence of Decision* brings out all of these complexities and uncertainties. They led him to conclude that the Cuban missile crisis cannot be explained in game-theoretic terms. He considers two alternatives: one explanation based on the fact that bureaucracies have their set rules and procedures, and another based on the internal politics of U.S. and Soviet governance and military apparatuses. He concludes that the political explanation is best.

We broadly agree, but interpret the Cuban missile crisis differently. It is not the case that game theory is inadequate for understanding and explaining the crisis; rather, the crisis was *not a two-person game*—United States versus USSR, or Kennedy versus Khrushchev. Each of the two “sides” was itself a complex coalition of players with differing objectives, information, actions, and means of communication. The players within each side were engaged in other games, and some of those players were also directly interacting with their counterparts on the other side. In other words, the crisis can be seen as a complex multiplayer game with players aligned into two broad coalitions. Kennedy and Khrushchev can be regarded as the top-level players in this game, but each was constrained by others in his own coalition with

divergent views and information, and neither had full control over the actions of those others. We argue that this more subtle game-theoretic perspective is not only a good way to look at the crisis, but also essential in understanding how to practice brinkmanship. We begin with some items of evidence that Allison emphasizes, as well as others that emerge from other writings.

First, there are several indications of divisions of opinion on each side. On the U.S. side, as already noted, there were wide differences within ExComm. In addition, Kennedy found it necessary to consult others, such as former president Eisenhower and leading members of Congress. Some of them had very different views; for example, Senator William Fulbright said in a private meeting that the blockade “seems to me the worst alternative” (*Tapes*, 271). The media and the political opposition would not give the president unquestioning support for too long either. Kennedy could not have continued on a moderate course if the opinion among his advisers and the public became decisively hawkish.

Individuals also *shifted* positions in the course of the two weeks. For example, McNamara was at first quite dovish, arguing that the missiles in Cuba were not a significant increase in the Soviet threat (*Tapes*, 89) and favoring a blockade and negotiations (*Tapes*, 191), but ended up more hawkish, claiming that Khrushchev’s conciliatory letter of Friday, October 26, was “full of holes” (*Tapes*, 495, 585) and urging an invasion (*Tapes*, 537). Most importantly, the U.S. military chiefs always advocated a far more aggressive response. Even after the crisis was over and nearly everyone thought the United States had won a major round in the Cold War, Air Force General Curtis LeMay remained dissatisfied and wanted action: “We lost! We ought to just go in there today and knock ’em off,” he said (*Essence*, 206; *Profile*, 425).

Even though Khrushchev was the dictator of the Soviet Union, he was not in full control of the situation. Differences of opinion on the Soviet side are less well documented, but for what it is worth, later memoirists have claimed that Khrushchev made the decision to install the missiles in Cuba almost unilaterally, and

that when he informed the members of the Presidium, they thought it a reckless gamble (*Tapes*, 674; *Gamble*, 180). There were limits to how far he could count on the Presidium to rubber-stamp his decisions. Indeed, two years after the crisis, the disastrous Cuban adventure was one of the main charges leveled against Khrushchev when the Presidium dismissed him from office (*Gamble*, 353 - 55). It has also been claimed that Khrushchev wanted to defy the U.S. blockade, and that only the insistence of First Deputy Premier Anastas Mikoyan led to the Soviets' cautious response (*War*, 521).

Various parties on the U.S. side had very different information and a very different understanding of the situation, and at times these differences led to actions that were inconsistent with the intentions of the leadership, or even against their explicit orders. The concept of an "air strike" to destroy the missiles is a good example. The nonmilitary people in ExComm thought this would be a very narrowly targeted attack that would not cause significant Cuban or Soviet casualties, but the Air Force intended a much broader attack. Luckily, this difference came out in the open early, leading ExComm to decide against an air strike and the president to turn down an appeal by the Air Force (*Essence*, 123, 209). As for the blockade, the U.S. Navy had set procedures for such an action. The political leadership wanted a different and softer process: form the ring closer to Cuba to give the Soviets more time to reconsider, allow obviously nonmilitary cargo ships to pass unchallenged, cripple but do not sink the ships that defy challenge. Despite McNamara's explicit instructions, the Navy mostly followed its standard procedures (*Essence*, 130 - 32).

There was a similar lack of information and communication, as well as weakness in the chain of command and control, on the Soviet side. For example, the construction of the missile sites was done according to standard bureaucratic procedures. The Soviets, used to constructing ICBM sites in their own country, where they did not face significant risk of air attack, laid out the sites in Cuba, where they would have been much more vulnerable, in a similar way.

All these factors made the outcome of any decision by the top-level commander on each side somewhat *unpredictable*. This gave rise to a substantial risk of the “threat going wrong.” And this risk was rising as the crisis continued. On the day the blockade went into effect, Kennedy thought that the chances of war were 20% (*Midnight*, 107); others attribute to him the higher estimate of “between one out of three and even” (*Essence*, 1).

Brinkmanship can use such uncertainty to strategic advantage to make one’s threat probabilistic and credible. In effect, Kennedy was saying to Khrushchev, “Neither of us wants nuclear war. But I can’t accept the missiles as a *fait accompli*. The quarantine I have set in motion creates a risk of war. You can end the game, and the risk, by withdrawing the missiles.” To achieve the advantage, Kennedy’s brinkmanship had to generate a risk high enough that Khrushchev would comply with his wishes, but low enough that he could tolerate it himself.

Neither the precise calculation of these limits nor the implementation of brinkmanship in practice is easy. We will develop two distinct models to bring out different aspects of the use of brinkmanship. In [Section 3](#), we will consider how the United States could calculate the limits of players’ tolerance for well-controlled risk. In [Section 4](#), we will examine some of the more practical difficulties of controlling risk. Then we will build a model in which risk rises gradually over time, and each player has to decide how long it is willing to accept the rising risk of disaster rather than conceding the game.

3 BRINKMANSHIP WITH WELL-CONTROLLED RISK

In this section, we present a model that considers one particular form of the uncertainty inherent in the Cuban missile crisis—namely, the United States' lack of knowledge of the Soviets' true motives in installing their missiles on the island. We analyze the effect of this type of uncertainty formally, and we draw conclusions about when and how President Kennedy could have hoped to use brinkmanship successfully. Similar analyses and conclusions can be drawn for other forms of uncertainty involved in this case.

A. When Is a Simple Threat Too Large?

Reconsider the game tree shown in Figure 13.1. Suppose the Soviet payoffs from withdrawal and defiance are the opposite of what they were before: -10 for withdrawal and -1 for defiance. In this alternative scenario, the Soviets are hard-liners. They prefer nuclear annihilation to the prospect of a humiliating withdrawal and the prospect of living in a world dominated by the capitalist United States; their slogan is “Better dead than red-white-and-blue.” We show the game tree for this scenario in Figure 13.2. Now, if the United States makes the threat, the Soviets defy it. So the United States stands to get -10 from the threat, but only -1 if it makes no threat and accepts the presence of the missiles in Cuba. It takes the lesser of the two evils. In the subgame-perfect equilibrium of this version of the game, the Soviets “win,” and the U.S. threat does not work.

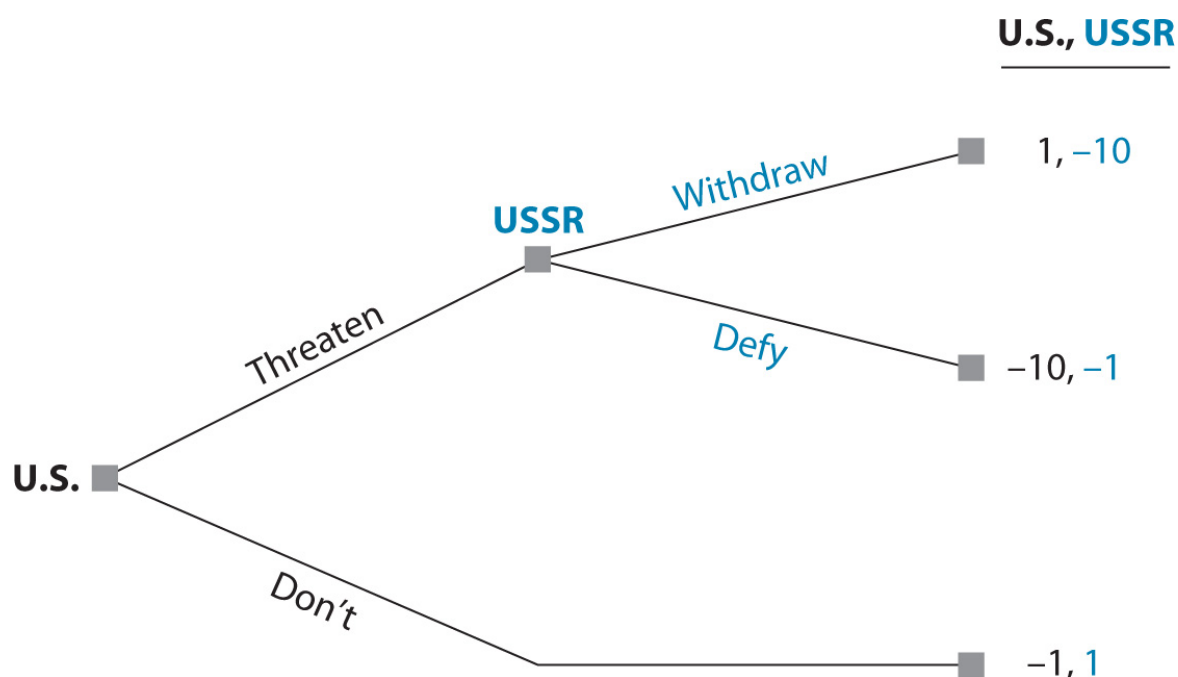


FIGURE 13.2 The Game with Hard-Line Soviets

In reality, when the United States makes its move, it does not know whether the Soviets' stance is hard-line, as in Figure

13.2, or softer, as in Figure 13.1. The United States can try to estimate the probabilities of the two scenarios—for example, by studying past Soviet actions and reactions in different situations. We can regard Kennedy's statement that the probability of the blockade leading to war was between one-third and one-half as his estimate of the probability that the Soviets were hard-line. Because the estimate is imprecise over a range, we work with a general symbol, p , for the probability, and we examine the consequences of different values of p .

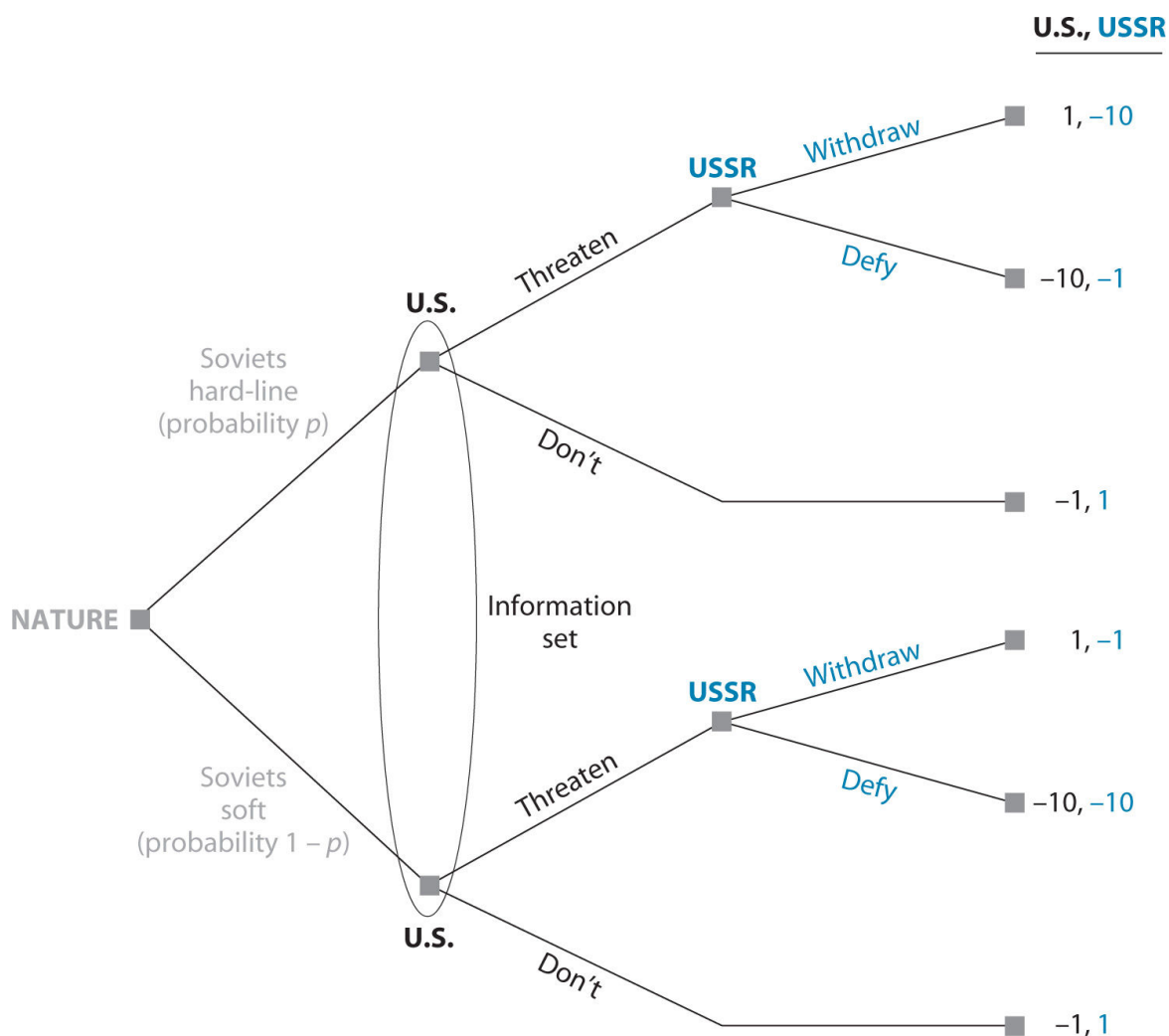


FIGURE 13.3 The Game with Unknown Soviet Payoffs

The tree for this more complex game is shown in Figure 13.3. The game starts with an outside force (here labeled Nature) determining the Soviets' type. Along the upper branch leading

from Nature's choice, the Soviets are hard-line. This branch leads to the upper node where the United States makes its decision whether to issue its threat, and the rest of the tree is exactly like that for the game in Figure 13.2. Along the lower branch leading from Nature's choice, the Soviets are soft. This branch leads to the lower node where the United States makes its decision whether to issue its threat, and the rest of the tree is exactly like that for the game in Figure 13.1. But the United States does not know at which node it is making its choice. Therefore, the two U.S. nodes constitute an *information set*, as indicated by the oval that encloses them. Its significance is that the United States cannot take different actions at the two nodes within the set, such as issuing the threat only if the Soviets are soft. It must take the same action at both nodes, either threatening at both nodes or not threatening at both. It must make this decision in light of the probabilities that the game is located at one node or the other—that is, by calculating the *expected* payoffs of the two actions.

The Soviets themselves know what type they are. So we can do some rollback near the end of the game. Along the upper path of play, the hard-line Soviets will defy a U.S. threat, and along the lower path, the soft Soviets will withdraw in the face of the threat. Therefore, the United States can look ahead and calculate that a threat will yield a -10 if the game is actually moving along the upper path (a probability of p) and a 1 if it is moving along the lower path (a probability of $1 - p$). The expected U.S. payoff from making the threat is therefore $-10p + (1 - p) = 1 - 11p$.

If the United States does not make the threat, it gets a -1 along either path; so its expected payoff is -1 . Comparing the expected payoffs of the two actions, we see that the United States should make the threat if $1 - 11p > -1$, or $11p < 2$, or $p < 2/11 = 0.18$.

If the threat were sure to work, the United States would not care how bad its payoff might be if the Soviets defied it, whether -10 or even far more negative. But the risk that the Soviets might be hard-liners and might thus defy a threat makes the -10 relevant in the U.S. calculations. Only if the probability, p , of the

Soviets being hard-line is small enough will the United States find it acceptable to make the threat. Thus, the upper limit of $2/11$ on p is also the upper limit of this U.S. tolerance, given the specific numbers that we have chosen. If we choose different numbers, we will get a different upper limit; for example, if we rate a nuclear war as -100 for the United States, then the upper limit on p will be only $2/101$. But the idea of a large threat being “too large to make” if the probability of its going wrong is above a critical limit holds in general.

In this instance, Kennedy’s estimate was that p lay somewhere in the range from $1/3$ to $1/2$. The lower end of this range, 0.33 , is unfortunately above our upper limit 0.18 for the risk that the United States is willing to tolerate. The simple bald threat “If you defy us, there will be nuclear war” is too large, too risky, and too costly for the United States to make.

B. The Probabilistic Threat

Suppose Kennedy makes a different kind of threat, one that reduces the large, but simple, threat described above by creating merely a probability, rather than a certainty, that the Soviets will incur dire consequences if they do not comply. With a [probabilistic threat](#) of this type, Kennedy is effectively declaring, “If you defy us, there is a probability q (< 1) that nuclear war will result.” It is important that the random mechanism that generates this outcome is out of Kennedy’s control after the fact, but he can set the probability q in advance. This game is like Russian roulette (a perfect metaphor in this context?). In that potentially lethal game of chance, you load one of the six chambers in a handgun, creating a one-sixth probability that a bullet will actually be fired. But when you pull the trigger after having spun the cylinder containing the chambers, you have no control over whether the chamber that gets fired is actually loaded.

Brinkmanship, the making and deployment of a probabilistic threat, requires creating and controlling a suitable risk of this kind. It requires two apparently inconsistent things. On the one hand, you must let matters get enough out of your control that you will not have full freedom after the fact to refrain from taking the dire action, so that your threat will remain credible. On the other hand, you must retain sufficient control to keep the risk of the action from becoming too large and your threat too costly. Such “controlled lack of control” looks difficult to achieve, and it is. We will consider in [Section 5](#) how to attempt the trick when you need to. Just one hint: All the complex differences of judgment, asymmetries of information, and difficulties of enforcing orders that made a simple threat too risky are exactly the forces that make it possible to create a risk of war and therefore make brinkmanship credible. The real difficulty is not how to lose control, but how to do so in a controlled way.

In the case of the U.S. - Soviet tensions over the Cuban missile sites, we need to ask what levels of q (the probability of nuclear war in Kennedy's brinkmanship threat) will be both *effective* in compelling Khrushchev to withdraw and tolerable (*acceptable*) to Kennedy given his estimated range of p (the probability that he faces hard-line Soviet opponents). To answer this question, we slightly alter the game of Figure 13.3 to get Figure 13.4. Here, if the Soviets defy the United States, war will occur with probability q . With the remaining probability, $(1 - q)$, the United States will give up and accept the presence of Soviet missiles in Cuba. Remember that if the game gets to the point where the Soviets defy the United States, the latter does not have a choice in the matter. The Russian-roulette revolver has been set for the probability q , and chance determines whether the firing pin hits a loaded chamber (that is, whether nuclear war actually happens).

Thus, nobody knows the precise outcome or the payoffs that will result if the Soviets defy this brinkmanship threat, but they know the probability, q , and can calculate expected values. For the United States, the outcome is -10 with the probability q and -1 with the probability $(1 - q)$, so the expected value is $-10q - (1 - q) = -1 - 9q$. For the Soviets, the expected payoff depends on whether they are hard-line or soft (and only they know their own type). If hard-line, they get -1 from war, which happens with probability q , and 1 if the United States gives up, which happens with probability $(1 - q)$. The hard-line Soviets' expected payoff is $-q + (1 - q) = 1 - 2q$. If they were to withdraw, they would get a -10 , which is clearly worse no matter what value q takes. Thus, the hard-line Soviets will defy the brinkmanship threat.

The calculation is different if the Soviets are soft. Reasoning as before, we see that they get the expected payoff $-10q + (1 - q) = 1 - 11q$ from defiance and the sure payoff -1 if they withdraw. For them, withdrawal is better if $-1 > 1 - 11q$, or $11q > 2$, or $q > 0.18$. Thus, U.S. brinkmanship must contain at least an 18% probability of war; otherwise, it will not deter the Soviets, even if they are the soft type. We call this lower limit on the probability q the [effectiveness condition](#).

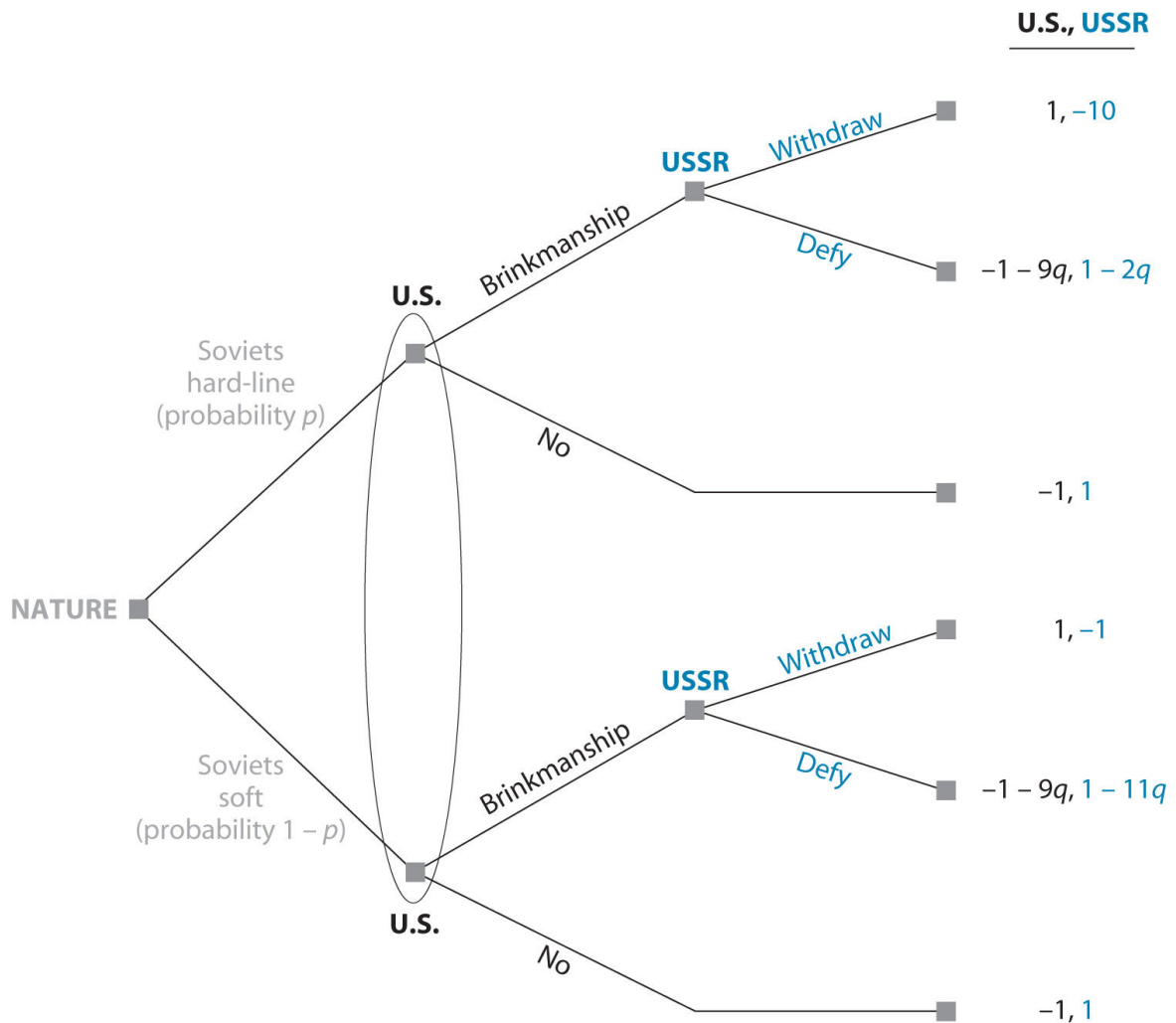


FIGURE 13.4 The Brinkmanship Model of the Crisis

Observe how the expected payoffs for U.S. brinkmanship and Soviet defiance shown in Figure 13.4 relate to the simple-threat model of Figure 13.3; the latter can now be thought of as a special case of the general brinkmanship-threat model of Figure 13.4, corresponding to the extreme value $q = 1$.

We can solve the game shown in Figure 13.4 in the usual way. We have already seen that along the upper path, the Soviets, being hard-line, will defy the United States, and that along the lower path, the soft Soviets will comply with U.S. demands if the effectiveness condition is satisfied. If the effectiveness condition is not satisfied, then both types of Soviets will defy the United States—in which case the United States would do

better never to make this threat at all. So let us proceed by assuming that the soft Soviets will comply; we now look at the U.S. choices. Basically, how risky can the threat be and still remain tolerable to the United States?

If the United States makes the threat, it runs the risk p that it will encounter the hard-line Soviets, who will defy the threat. Then the expected U.S. payoff will be $-1 - 9q$, as calculated before. The probability is $(1 - p)$ that the United States will encounter the soft Soviets. We are assuming that they comply; then the United States gets a 1. Therefore, the expected payoff to the United States from the probabilistic threat, assuming that it is effective against the soft Soviets, is

$$(-1 - 9q) \times p + 1 \times (1 - p) = -9pq - 2p + 1.$$

If the United States refrains from making a threat, it gets a -1 . Therefore, the condition for the United States to make the threat is

$$-9pq - 2p + 1 > -1, \text{ or}$$

$$q < \frac{2}{9} \frac{1 - p}{p} = \frac{0.22(1 - p)}{p}.$$

That is, the probability of war must be small enough to satisfy this expression, or the United States will not make the threat at all. We call this upper limit on q the [acceptability condition](#). Note that p enters the formula for the maximum value of q that will be acceptable to the United States; the larger the chance that the Soviets will not give in, the smaller the risk of mutual disaster that the United States finds acceptable.

For the probabilistic threat to benefit the United States, it should satisfy both the effectiveness condition and the acceptability condition. We can determine the resulting

probability of war by using Figure 13.5. The horizontal axis is the probability, p , that the Soviets are hard-line, and the vertical axis is the probability, q , that war will occur if they defy the U.S. threat. The horizontal line $q = 0.18$ gives the effectiveness condition. The probabilistic threat will be made only if its associated (p, q) combination is above this line, since otherwise, it will not deter even the soft-type Soviets. The curve $q = 0.22(1 - p)/p$ gives the acceptability condition. The probabilistic threat will be made only if (p, q) is below this curve, since otherwise, the resulting risk of war will be intolerable to the United States, even if soft-type Soviets are successfully deterred. Therefore, an effective and acceptable threat should fall somewhere between these two lines, above and to the left of their point of intersection at $p = 0.55$ and $q = 0.18$ (shown as a gray wedge in Figure 13.5).

The curve reaches $q = 1$ when $p = 0.18$. For values of p less than 0.18, the dire threat (certainty of war) is acceptable to the United States and is effective against the soft-type Soviets. This simply confirms our analysis in [Section 3.A](#).

For values of p in the range from 0.18 to 0.55, the dire threat with $q = 1$ puts (p, q) above the acceptability condition and is too large to be tolerable to the United States. But a scaled-down threat can be found. For this range of values of p , some values of q are low enough to be acceptable to the United States and yet high enough to compel the soft-type Soviets. Brinkmanship (using a probabilistic threat) can do the job in this situation, whereas a simple dire threat would be too risky.

If p exceeds 0.55, then no value of q satisfies both conditions. If the probability that the Soviets will never give in is greater than 0.55, then any threat large enough to work against the soft-type Soviets ($q \geq 0.18$) creates a risk of war too large to be acceptable to the United States. If $p \geq 0.55$, therefore, the United States cannot use the brinkmanship strategy to its advantage.

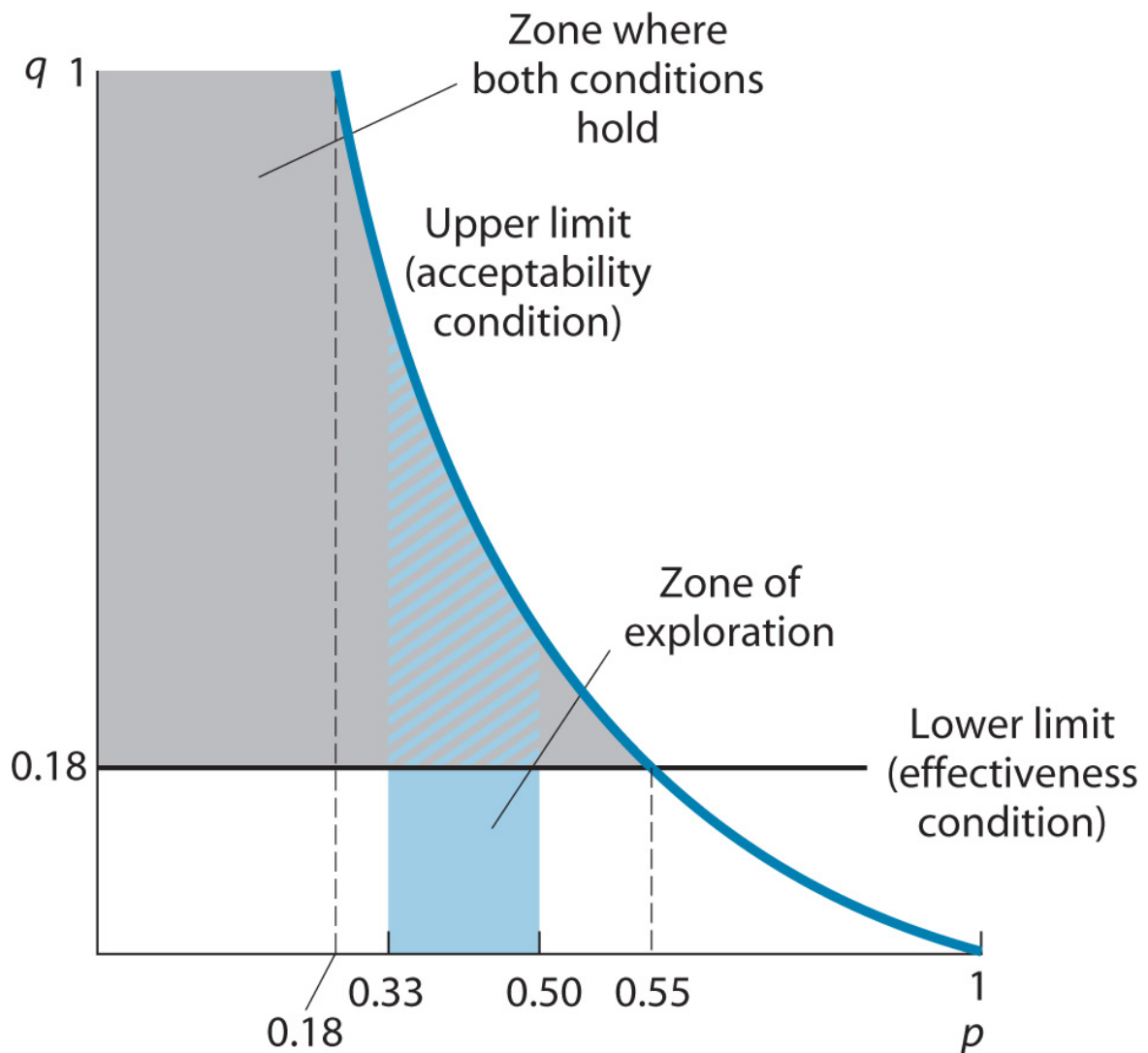


FIGURE 13.5 Conditions of Successful Brinkmanship

If we regard Kennedy's estimate of the probability of war (originally 20% and later increased to "between one-third and one-half") as his estimate that the Soviets were hard-line, it falls squarely within the range of p where brinkmanship can be successfully practiced (p between 0.18 and 0.55). So, should we regard the episode as a result of Kennedy's brilliant assessment of the situation and exercise of brinkmanship?⁶ His calculation would have required assumptions about the payoffs for the two types of Soviets, as well as clarity about his own payoffs. In Exercise U1 we will ask you to calculate the consequences of one change in these assumptions. Now we turn to the practical difficulties of controlling risk.

Endnotes

- Or as vindication of our brilliant choice of numbers? 😊
[Return to reference 6](#)

Glossary

[probabilistic threat](#)

A strategic move in the nature of a threat, but with the added qualification that if the event triggering the threat (the opponent's action in the case of deterrence or inaction in the case of compellence) comes about, a chance mechanism is set in motion, and if its outcome so dictates, the threatened action is carried out. The nature of this mechanism and the probability with which it will call for the threatened action must both constitute prior commitments.

[effectiveness condition](#)

A lower bound on the probability of fulfillment in a brinkmanship threat, expressed as a function of the probability of error, showing the lower limit of risk that will induce the threatened player to comply with the wishes of the threatener.

[acceptability condition](#)

An upper bound on the probability of fulfillment in a brinkmanship threat, expressed as a function of the probability of error, showing the upper limit of risk that the player making the threat is willing to tolerate.

4 BRINKMANSHIP WITH UNCONTROLLED RISK : A GAME OF DYNAMIC CHICKEN

Our account of brinkmanship, while a substantial improvement on the pure threat model of [Section 2](#), remains inadequate. First, it still does not explain why the Soviets installed the missiles in Cuba in the first place. Why didn't they look ahead to the subsequent game of brinkmanship outlined above and conclude that their best strategy was not to play at all? Presumably the Soviets thought it sufficiently likely that the United States would prove to be soft itself and accept the missiles in Cuba, just as the Soviets had learned to live with U.S. missiles in Turkey. In other words, the Soviets themselves were unsure whether the United States was a soft or hard-line type; there was *two-sided* uncertainty about payoffs.

Even more importantly, our discussion in [Section 3](#) assumed the ability to fix q —that is, to achieve a *controlled* loss of control. However, particularly in the last couple of days of the crisis, events were rapidly spinning out of the control of the principals, Kennedy and Khrushchev, into a realm of uncontrolled loss of control!

Daniel Ellsberg, who later became famous for leaking the Pentagon Papers, was a young researcher for the Department of Defense during the crisis. At that time, he and his immediate superiors estimated the probability of war to be very low, between one in a thousand and one in a hundred (*Doomsday*, 189, 199). The United States had overwhelming superiority in nuclear weapons and missiles, so they thought that, given any degree of rationality, Khrushchev simply had to back down, even without a significant concession from the United States. Ellsberg was astounded to hear the higher numbers offered by more senior people in ExComm, but later realized that they were correct. He and other relatively junior participants had not known the extent to which the top people were losing control of the situation. In fact, those leaders themselves came to realize the extent of their loss of control very late in the game (*Doomsday*, 201–22).

At numerous points—for example, when the U.S. Navy was trying to stop and board the freighter *Marucla*—the people involved might have set off an incident with alarming consequences by taking some action in fear of the immediate situation. Most dramatically, a Soviet submarine crew, warned to surface when approaching the quarantine line on October 27, considered firing a nuclear-tipped torpedo that it carried onboard (unknown to the U.S. Navy). The firing-authorization rule required the approval of three officers, only two of whom agreed; the third officer by himself may have prevented all-out nuclear war.⁷

The U.S. Air Force created even greater dangers. A U-2 plane that drifted “accidentally” into Soviet air space almost caused a serious setback. General Curtis LeMay, acting without the president's knowledge or authorization, ordered the Strategic Air Command's nuclear bombers to fly past their “turnaround” points and some distance toward Soviet air space to positions where they would be detected by Soviet radar. Fortunately, the Soviets responded calmly; Khrushchev merely protested to Kennedy.⁸

On Saturday, October 27, Castro ordered his anti-aircraft forces to fire on all U.S. planes overflying Cuba and refused the Soviet ambassador's request to rescind the order (*War*, 544). On the same day, an overflying U.S. U-2 plane was shot down by a Soviet surface-to-air missile; a lower-level local commander had interpreted his orders more broadly than Moscow had intended [*War*, 537; *Tapes*, 682].

With each day that the crisis continued, the probability of nuclear war was increasing, and the top leaders were losing control over it. Combine this rising risk of disaster with each side's uncertainty about the other's payoffs, and we have a game of dynamic chicken, which

we introduced in [Chapter 9](#). Instead of two teenagers driving their cars toward each other as the risk of a collision increases, testing each other's bravery (or foolhardiness) in deciding when to swerve, we have the two leaders of the two major powers of the day, testing each other's resolve as the risk of Armageddon increases. Unlike the model in [Section 13.3](#), where Kennedy controlled the level of risk, this model shows the risk increasing steadily through time, and Kennedy and Khrushchev must decide each day whether to continue the confrontation at the new, higher level of risk or to concede the game (Kennedy by accepting the missiles or Khrushchev by withdrawing them).

In [Chapter 9](#), we solved a two-step version of such a game. Here, we apply the method developed in [Chapter 9](#), but we do so for a game with 13 steps, one for each day of the crisis.⁹ In our original dynamic chicken example, the value each player placed on winning, W , was that player's private information; in keeping with the context here, we cast each side's perceived cost of catastrophe, C , in that role. We still have two-sided asymmetric information, as in [Chapter 9](#), but the value unknown to the other player has changed. We fix the value of winning at $W = 1$ and the cost of losing at $L = 1$, and take each player's cost of catastrophe C to be uniformly distributed over the interval from 0 to 10.¹⁰ Because you would surely defy a threat made by the other player if your C was below L , our choice of values makes the implied probability that each side is hard-line equal to one-tenth.¹¹ A player at the soft end (with a high cost of catastrophe, C , close to 10) would concede quickly. Each day, as the probability of disaster rises, only successively tougher players remain; the threshold value of C , beyond which a player concedes, falls.

Figure 13.6 shows the results of our calculations. For easy reference, we include a brief summary of the major events of each day of the crisis in the second column.¹² The third column shows the values we have chosen (our best "guesstimates") for the rising probabilities of disaster. These probabilities start very low, rise slowly over the middle days to conform to Kennedy's estimates cited previously, and then rise sharply in the last two days of the crisis because of the loss of control detailed previously. On the last day of the crisis (Sunday, October 28), the risk rises to 1.000 because if Khrushchev had not conceded then, the U.S. invasion plans would have gone into effect.

The last four columns in Figure 13.6 show the results of the calculations for the dynamic chicken model. In the fourth column, we report the threshold levels of C ; a player with a C value above the threshold listed for a particular step of the game would concede at that step. The last three columns give probability calculations: the probability that each player concedes at a particular step, the probability of disaster at each step, and the cumulative probability of disaster at each step. Observe that on the first day of the crisis (October 16), the threshold value of C reported in the fourth column is 9.852, close to the maximum of the range of possible values of C (which is 10). Only the most peace-loving, dovish players (with C values in the top 1.5% of the full range) would concede at that stage.

Our results for the probabilities of concession clarify two features of the crisis at once. First, we see why Khrushchev may have embarked on this adventure at all. He clearly put significant value on the prestige that a victory in this confrontation would yield to the USSR (and to him as its leader) in other communist countries, in the third world, and in many countries of the Western bloc as well. A victory would also strengthen his position in other confrontations with the United States, especially over Berlin. With a C not too close to 10, it would be rational for Khrushchev to launch the game, and not to concede for quite a while. Combine this with the high likelihood that he expected Kennedy to be weak (to have a high C , in the language of our model). Kennedy's performances during their meeting in Vienna in 1961, and during the attempted invasion of Cuba by exiles at the Bay of Pigs the same year, were widely perceived to be weak; this may have led Khrushchev to believe that Kennedy would concede quickly, accepting the missiles in Cuba as a *fait accompli*. Second, we can infer that both principals, Kennedy and Khrushchev, must have been very hard-line to have lasted as long as they did without conceding. Even if we accept October 22 as the date

when the Soviet Presidium decided to concede, that date implies that the Soviets had a C value of 4.329, significantly above L ($=1$) and well below the maximum possible in our model. However, it is also possible that both sides underestimated the risk of war for quite a while, and that only the near-confrontations that occurred in the last two days of the crisis brought home to them the extent to which they had lost control of the situation.

Figure 13.6 also shows the probability of concession decreasing in the final few days, dropping to 0.16 on the final Sunday; this probability is quite small because, by then, only the toughest players are left. Our result for the conditional probability of disaster on Saturday (October 27, day 12) is high enough at 46.4% to make McNamara's fear as he left the White House that beautiful fall evening that this "might be the last sunset I saw" (*Doomsday*, 201) quite understandable. And the predicted cumulative probability of disaster over the whole duration of the crisis is high enough (at 58.3%) that we should indeed be thankful that it did not end in a nuclear disaster.

Date (October 1962)	Major events	Probability of disaster at this step if neither concedes	Concession threshold	Probability that each player concedes at this step	Probability of disaster at this step (conditional on it being reached)	Cumulative probability of disaster up to and including this step
15	U-2 photos show MRBM missile sites					
16	JFK briefed, convenes ExComm	0.010	9.852	0.015	0.010	0.010
17	New overflight, evidence of IRBMs	0.025	9.120	0.074	0.021	0.030
18	ExComm continues meetings. Differences of opinion; changes of mind; options debated.	0.040	7.971	0.126	0.031	0.055
19		0.050	6.853	0.140	0.037	0.080
20		0.060	5.919	0.136	0.045	0.109
21	ExComm majority for blockage/quarantine	0.075	5.054	0.146	0.055	0.144

You may need to scroll left and right to see the full figure.

Date (October 1962)	Major events	Probability of disaster at this step if neither concedes	Concession threshold	Probability that each player concedes at this step	Probability of disaster at this step (conditional on it being reached)	Cumulati probabil of disas up to an includin this ste
22	JFK addresses nation. Soviets deny, bluster. Soviet Presidium meets in secret and decides USSR will eventually withdraw but will first explore best possible deals in exchange.	0.100	4.329	0.143	0.073	0.186
23	UN Security Council meeting	0.150	3.533	0.184	0.100	0.240
24	Blockade goes into effect. Some Soviet ships approaching quarantine line stop or reverse, others continue.	0.200	2.889	0.182	0.134	0.297
25	U-2 overflights increase. <i>Marucla</i> pursued.	0.300	2.330	0.194	0.195	0.371

You may need to scroll left and right to see the full figure.

Date (October 1962)	Major events	Probability of disaster at this step if neither concedes	Concession threshold	Probability that each player concedes at this step	Probability of disaster at this step (conditional on it being reached)	Cumulati probabil of disas up to an includin this ste
26	NSK' s conciliatory letter (asking "no invasion" promise). <i>Marucla</i> boarded. Castro authorizes anti-aircraft fire on low-flying U.S. reconnaissance flights. Soviet junior officer in Cuba nearly fires missile at overflying U-2 while boss is away from desk. Attitudes harden in ExComm.	0.500	1.801	0.227	0.299	0.456
27	NSK' s hard-line letter (Cuba/Turkey swap); JFK independently offers this as secret deal in letter RFK hands over to Dobrynin. Overflying U-2 shot down. Soviet sub crew considers firing nuclear torpedo but fails in getting required unanimity of three officers. ExComm feeling: "We haven' t got but one more day." U.S. air strike planned for Monday, October 29.	0.750	1.416	0.214	0.464	0.531

You may need to scroll left and right to see the full figure.

Date (October 1962)	Major events	Probability of disaster at this step if neither concedes	Concession threshold	Probability that each player concedes at this step	Probability of disaster at this step (conditional on it being reached)	Cumulati probabil of disas up to an includin this ste
28	NSK speech withdrawing missiles	1.000	1.190	0.160	0.706	0.583

You may need to scroll left and right to see the full figure.

FIGURE 13.6 A Dynamic Chicken Model of the Crisis

Endnotes

- This story became public in a conference held in Havana, Cuba, in October 2002, to mark the 40th anniversary of the missile crisis. See Kevin Sullivan, “40 Years After Missile Crisis, Players Swap Stories in Cuba,” *Washington Post*, October 13, 2002, p. A28. Vadim Orlov, who was a member of the Soviet submarine crew, identified the officer who refused to fire the torpedo as Vasili Arkhipov, who died in 1999. See also *Doomsday*, pp. 216 – 17. [Return to reference 7](#)
- Richard Rhodes, *Dark Sun: The Making of the Hydrogen Bomb* (New York: Simon & Schuster, 1995), pp. 573 – 75. LeMay, renowned for his extreme views and his constant chewing of large unlit cigars, is supposed to be the original inspiration for General Jack D. Ripper, in the 1963 movie *Dr. Strangelove*, who orders his bomber wing to launch an unprovoked attack on the Soviet Union. [Return to reference 8](#)
- The mathematical details of the model are in our working paper, “We Haven’ t Got But One More Day: The Cuban Missile Crisis as a Dynamic Chicken Game,” (June 2019). Available at SSRN: <https://ssrn.com/abstract=3406265>. [Return to reference 9](#)
- In the formal algebra of the model, we should say the payoffs are $L=-1$ and C between 0 and -10 , but the positive values are easier to write in a verbal account. [Return to reference 10](#)
- Note that the probabilities in this dynamic chicken interpretation of brinkmanship are not directly comparable to those we calculated in the one-sided version in Section 3.B. [Return to reference 11](#)
- JFK is, of course, John F. Kennedy, RFK is his brother Robert, and NSK is Nikita Sergeyevich Khrushchev. [Return to reference 12](#)

5 PRACTICING BRINKMANSHIP

The very features of the Cuban missile crisis that make it inaccurate to regard it as a two-person game made it easier for the players to practice brinkmanship. The blockade of Cuba was a relatively small action, unlikely to start a nuclear war at once. But once Kennedy had set the blockade in motion, its operation, escalation, and other features were not totally under his control. So Kennedy was not saying to Khrushchev, “If you defy me (cross a sharp brink), I will coolly and deliberately launch a nuclear war that will destroy both our peoples.” Rather, he was implicitly saying, “The wheels of the blockade have started to turn and are gathering their own momentum. The longer you defy me, the more likely it is that some operating procedure will slip up, the domestic political pressure on me will rise to a point where I must give in to the hawks, or some military guy will run amok. If any of these things come to pass, I may be unable to prevent nuclear war, no matter how much I may regret it at that point. Only you can now defuse the tension by complying with my demand to withdraw the missiles.” And Khrushchev was making similar implicit statements to Kennedy up until the moment when he decided to concede.

We believe that this perspective gives a much better and deeper understanding of the crisis than can most analyses based on simple threats. It tells us why the *risk* of war played such an important role in all discussions. It even makes Allison’s compelling arguments about bureaucratic procedures and internal divisions on both sides an integral part of the picture: These features allowed the top-level players on both sides to lose some control credibly—that is, to practice brinkmanship.

One important condition remains to be discussed. In [Chapter 8](#), we saw that every threat has an associated implicit

affirmation—namely, that the bad consequence will not take place if your opponent complies with your wishes. The same is required for brinkmanship. If, as you are increasing the level of risk, your opponent does comply, you must be able to “go into reverse” —begin reducing the risk immediately and quickly remove it from the picture. Otherwise, the opponent would not gain anything by compliance. This might have been a problem in the Cuban missile crisis. If the Soviets had feared that Kennedy could not control hawks such as LeMay (“We ought to just go in there today and knock ’ em off”), they would have gained nothing by giving in.

To review and sum up, brinkmanship is the strategy of exposing your rival and yourself to a gradually increasing risk of mutual harm. The actual occurrence of the harmful outcome is not totally within the threatener’ s control. But the loss of control itself needs to be controlled; the probability of disaster must be kept within certain bounds to ensure that the risk is acceptable to the threatener. This “controlled loss of control” is difficult to achieve, and in the last few days of the Cuban missile crisis, the situation was close to becoming totally uncontrolled by the two principals.

Because the level of risk in a game of brinkmanship is difficult to control, it makes sense to start with a low level of risk and let it rise gradually in its own way. This approach allows you to find out whether the opponent’ s tolerance for the rising risk runs out before your own does. In other words, this approach allows you to stand eyeball to eyeball with your opponent and to see who blinks first.^{[13](#)}

A game of brinkmanship can end in one of three ways: with success for one side and loss for the other, or in mutual disaster. Fortunately, the Cuban missile crisis did not end disastrously; if it had, none of us would be here to analyze it as a case study.

Viewed as a strategy that entails increasing the risk of disaster for both sides of an interaction, brinkmanship is everywhere. In most confrontations—for example, between a company and a labor union, a husband and a wife, a parent and a child, or the president and Congress—one player cannot be sure of the other players' objectives and capabilities. Therefore, most threats carry a risk of error, and every threat must contain an element of brinkmanship. We hope that we have given you some understanding of this strategy and that we have impressed on you the risks that it carries. Unsuccessful brinkmanship can lead to a labor strike, the dissolution of a marriage, a shutdown of the U.S. government, or some other disaster that, while small compared with nuclear annihilation, looms large in the context of the game you are playing. You will have to face brinkmanship, or conduct it yourself, on many occasions in your personal and professional lives. Please do so carefully, with a clear understanding of its potentialities and risks.

To help you do so, we now recapitulate the important lessons learned from the handling of the Cuban missile crisis, reinterpreted here in the context of a labor union leader contemplating a strike in pursuit of the union's demand for higher wages, unsure whether this action will result in the whole firm's shutting down:

1. Start small and safe. Your first step should not be an immediate walkout; it should be to schedule a membership meeting at a date a few days or weeks hence, while negotiations continue.
2. Let the risks increase gradually. Your public and private statements, as well as the stirring up of the sentiments of the membership, should induce management to believe that acceptance of its current low-wage offer is becoming less and less likely. If possible, stage small incidents—for example, a few one-day strikes or local walkouts.

3. As this process continues, read and interpret signals in management' s actions to figure out whether the firm has enough profit potential to afford the union' s high-wage demand.
4. Try to retain enough control over the situation; that is, retain the power to induce your membership to ratify the agreement that you will reach with management; otherwise management will think that the risk of a strike will not decrease even if it concedes to your demands.
5. Remain alert for signs that the situation is getting out of your control, and be ready to reassert control and de-escalate, either when the opponent concedes or by your own concession.

Endnotes

- This image was invoked by Secretary of State Dean Rusk when Soviet ships bound for Cuba appeared to stop or reverse on October 24 (*Midnight*, 88). [Return to reference 13](#)

SUMMARY

In some game situations, the risk of error in the presence of a threat may call for the use of as small a threat as possible. When a large threat cannot be reduced in other ways, it can be scaled down by making its fulfillment probabilistic. Strategic use of a *probabilistic threat*, in which you expose your rival and yourself to an increasing risk of harm, is called *brinkmanship*.

Brinkmanship requires a player to relinquish some control over the outcome of the game without completely losing control. You must create a threat with a risk level that is both large enough to be *effective* in compelling or deterring your rival and small enough to be *acceptable* to you. To do so, you must test the limit of your opponent's risk tolerance, while going up to your own limit if necessary, through a *gradual escalation of the risk of mutual harm*.

The Cuban missile crisis of 1962 serves as a case study in the use of brinkmanship. Analyzing the crisis as an example of a simple threat, with the U.S. blockade of Cuba establishing credibility for its threat, is inadequate. A better analysis accounts for the many complexities and uncertainties inherent in the situation and the likelihood that a simple threat was too risky. Because the actual crisis included numerous political and military players, Kennedy could attempt "controlled loss of control" by ordering the blockade and gradually letting incidents and tension escalate, until Khrushchev yielded in the face of the rising risk of nuclear war. A model of dynamic chicken captures such uncertainty and gives better explanations, and even plausible numerical magnitudes, for the unfolding of the crisis.

KEY TERMS

[acceptability condition](#) ([534](#))

[brinkmanship](#) ([518](#))

[effectiveness condition](#) ([533](#))

[probabilistic threat](#) ([531](#))

Glossary

brinkmanship

A threat that creates a risk but not certainty of a mutually bad outcome if the other player defies your specified wish as to how he should act, and then gradually increases this risk until one player gives in or the bad outcome happens.

probabilistic threat

A strategic move in the nature of a threat, but with the added qualification that if the event triggering the threat (the opponent's action in the case of deterrence or inaction in the case of compellence) comes about, a chance mechanism is set in motion, and if its outcome so dictates, the threatened action is carried out. The nature of this mechanism and the probability with which it will call for the threatened action must both constitute prior commitments.

effectiveness condition

A lower bound on the probability of fulfillment in a brinkmanship threat, expressed as a function of the probability of error, showing the lower limit of risk that will induce the threatened player to comply with the wishes of the threatener.

acceptability condition

An upper bound on the probability of fulfillment in a brinkmanship threat, expressed as a function of the probability of error, showing the upper limit of risk that the player making the threat is willing to tolerate.

SOLVED EXERCISES

1. Consider a game between a union and the company that employs the union membership. The union can threaten to strike (or not) to get the company to meet its wage and benefit demands. When faced with a threatened strike, the company can choose to concede to the demands of the union or to defy its threat of a strike. The union, however, does not know the company's profit position when it decides whether to make its threat; that is, it does not know whether the company is sufficiently profitable to meet its demands—and the company's assertions in this matter cannot be believed. Nature determines whether the company is profitable; the probability that the firm is unprofitable is p .

The payoff structure is as follows: (i) When the union makes no threat, the union gets a payoff of 0 (regardless of the profitability of the company). The company gets a payoff of 100 if it is profitable, but a payoff of 10 if it is unprofitable. A passive union leaves more profit for the company if there is any profit to be made. (ii) When the union threatens to strike and the company concedes, the union gets 50 (regardless of the profitability of the company) and the company gets 50 if it is profitable but -40 if it is not. (iii) When the union threatens to strike and the company defies the union's threat, the union must strike and gets -100 (regardless of the profitability of the company). The company gets -100 if it is profitable and -10 if it is not. Defiance is very costly for a profitable company but not so costly for an unprofitable one.

1. What happens when the union uses the pure threat to strike unless the company concedes to the union's

demands?

2. Suppose that the union sets up a situation in which there is some risk, with probability $q < 1$, that it will strike after the company defies its threat. This risk may arise from the union leadership's imperfect ability to keep the membership in line. Draw a game tree similar to Figure 13.4. for this game.
 3. What happens when the union uses brinkmanship, threatening to strike with some probability q unless the company accedes to its demands?
 4. Derive the effectiveness and acceptability conditions for this game, and determine the values for p and q for which the union can use a pure threat, brinkmanship, or no threat at all.
2. The professor teaching a course has established a rule that any late homework must receive a failing grade, but ultimately a grader decides what grades to assign. Consider the resulting game between the grader and a student. The grader moves first, either threatening to follow through on the harsh official policy or announcing a more lenient approach. The student then decides whether to complete the homework on time or be late. The grader wants homework to be on time, but also does not want to fail the student, since then the student is sure to complain and try to get a grade change from the professor, dragging the grader into an unpleasant process. The grader gets payoff +1 when homework is on time, payoff 0 when homework is late and not failed, or payoff -3 when homework is late and failed. The student prefers not to fail but also prefers to do the work late because he is busy with other activities. With probability p , those other activities are so important that the student will choose to be late even if that leads to failure: This "unwilling type" gets payoff -10 from being on time, 0 from being late and not failing, and -5 from failing. The rest of the time, the

student prefers to be on time rather than fail: This “willing type” gets payoff -1 from being on time, 0 from being late and not failing, and -5 from failing.

1. What happens when the grader uses the pure threat to fail the student unless homework is on time?
 2. Suppose that the grader sets up a situation in which there is some risk, with probability $q < 1$, that she will fail the student for late homework. Such risk could arise from uncertainty regarding whether the professor is truly serious about the policy. Draw a game tree similar to Figure 13.4 for this game.
 3. What happens when the grader uses brinkmanship, threatening to fail the student with some probability q unless homework is on time?
 4. Derive the effectiveness and acceptability conditions for this game, and determine the values for p and q for which the grader can use a pure threat, brinkmanship, or no threat at all.
3. Scenes from many movies illustrate the concept of brinkmanship. Analyze the following descriptions from this perspective. What are the risks the two sides face? How do those risks increase during the course of the execution of the brinkmanship threat?
1. In the 1980 film *The Gods Must Be Crazy*, the only survivor of a rebel team that tried to assassinate the president of an African country has been captured and is being interrogated. He stands blindfolded with his back to the open door of a helicopter. Above the noise of the helicopter rotors, an officer asks him, “Who is your leader? Where is your hideout?” The man does not answer, and the officer pushes him out the door. In the next scene, we see that although its engine is running, the helicopter is actually on the ground, and the man has fallen 6 feet and landed on his back. The officer appears at the door and says, laughing, “Next time it will be a little higher.”

2. In the 1998 film *A Simple Plan*, two brothers remove some of a \$4.4 million ransom payment that they find in a crashed airplane. After many intriguing twists of fate, the remaining looter, Hank, finds himself in conference with an FBI agent. The agent, who suspects, but cannot prove, that Hank has some of the missing money, fills Hank in on the story of the money's origins and tells him that the FBI possesses the serial numbers of about 1 of every 10 of the bills in that original ransom payment. The agent's final words to Hank are, "Now it's simply a matter of waiting for the numbers to turn up. You can't go around passing \$100 bills without eventually sticking in someone's memory."
4. In this exercise, we provide two examples of the successful use of brinkmanship, where "success" consists of the two parties' reaching a mutually acceptable deal. For each example, (i) identify the interests of the parties; (ii) describe the nature of the uncertainty inherent in the situation; (iii) give the strategies the parties used to escalate the risk of disaster; (iv) discuss whether those strategies were good ones; and (v) (Optional) if you can, set up a small mathematical model of the kind presented in this chapter. In each case, we provide a few readings to get you started; you should locate more by using the resources of your library and online resources such as Lexis-Nexis.
 1. The Uruguay Round of international trade negotiations that started in 1986 and led to the formation of the World Trade Organization in 1994. *Reading*: John H. Jackson, *The World Trading System*, 2nd ed. (Cambridge, Mass.: MIT Press, 1997), pp. 44 - 49 and Chapters 12 and 13.
 2. The Camp David Accords between Israel and Egypt in 1978. *Reading*: William B. Quandt, *Camp David: Peacemaking and Politics* (Washington, D.C.: Brookings Institution, 1986).

5. The following examples illustrate the unsuccessful use of brinkmanship, where brinkmanship is considered “unsuccessful” when the mutually bad outcome (disaster) occurs. Answer the questions listed in Exercise S4 for each example.
 1. The confrontation between the Chinese communist regime and the student prodemocracy demonstrators in Beijing in June 1989. *Readings*: Donald Morrison, ed., *Massacre in Beijing: China’s Struggle for Democracy* (New York: Time Magazine Publications, 1989); Suzanne Ogden, Kathleen Hartford, L. Sullivan, and D. Zweig, eds., *China’s Search for Democracy: The Student and Mass Movement of 1989* (Armonk, N.Y.: M. E. Sharpe, 1992).
 2. The Caterpillar strike, from 1991 to 1998. *Readings*: “The Caterpillar Strike: Not Over Till It’s Over,” *Economist*, February 28, 1998; “Caterpillar’s Comeback,” *Economist*, June 20, 1998; Aaron Bernstein, “Why Workers Still Hold a Weak Hand,” *BusinessWeek*, March 2, 1998.
6. Answer the questions listed in Exercise S4 for these potential opportunities for brinkmanship in the future:
 1. A Taiwanese declaration of independence from the People’s Republic of China. *Reading*: Ian Williams, “Taiwan’s Independence,” *Foreign Policy in Focus*, December 20, 2006, available at www.fpif.org/fpiftxt/3815.
 2. The militarization of space; for example, the positioning of weapons in space or the shooting down of satellites. *Reading*: “Disharmony in the Spheres,” *Economist*, January 17, 2008, available at www.economist.com/node/10533205.

UNSOLVED EXERCISES

1. In the calculations of [Section 3](#), we assumed that the payoff to the United States is -10 when Soviets of either type (hard-line or soft) defy the U.S. threat, as illustrated in Figure 13.3. Suppose now that this payoff is -25 rather than -10 .
 1. Incorporate this change in payoff into a game tree similar to the one in Figure 13.4.
 2. Using the payoffs from your game tree in part (a), find the effectiveness condition for this version of the U.S. - USSR brinkmanship game.
 3. Using the payoffs from part (a), find the acceptability condition for this game.
 4. Draw a diagram similar to that in Figure 13.5, illustrating the effectiveness and acceptability conditions found in parts (b) and (c).
 5. For what values of p , the probability that the Soviets are hard-line, is the pure threat ($q = 1$) acceptable? For what values of p is the pure threat unacceptable but brinkmanship still possible?
 6. If Kennedy was correct in believing that p lay between one-third and one-half, does your analysis of this version of the game suggest that an effective *and* acceptable probabilistic threat existed? Use this example to explain how a game theorist's assumptions about player payoffs can affect the predictions that arise from the theoretical model.
2. A corrupt Russian policeman is interrogating a mobster who knows where \$1 million is hidden. The policeman can threaten to kill the mobster unless he reveals the money's location. The policeman moves first, deciding whether to make the threat, after which the mobster decides whether to reveal the secret location. The policeman wants the money, but also prefers not to murder

the mobster, since then there might be an investigation that would uncover his own crimes. The policeman gets \$1,000,000 (normalized to payoff +1) if the mobster reveals the money's location, payoff 0 if the mobster stays quiet and is left alive, and payoff -0.5 if the mobster stays quiet and is killed. With probability $(1 - p)$, the mobster is willing to tell the secret to save his life. This "loose-lipped type" gets payoff 0 when staying quiet and living, payoff -2 when telling the secret and living, and payoff -10 when staying quiet and dying. The rest of the time, with probability p , the mobster is a "tight-lipped type" who would rather die than spill the beans. This type gets payoff 0 when staying quiet and living, payoff -2 when telling the secret and living, and payoff -1 from staying quiet and dying.

1. What happens when the policeman uses the pure threat to kill the mobster unless he tells the secret?
2. The corrupt policeman has a gun that can hold six bullets. By loading the gun with $b = 1, 2, 3, 4$, or 5 bullets and then spinning the gun's cylinder, the policeman can create a risk $q = b/6 < 1$ that there is a bullet in the chamber that he is about to fire. This allows the policeman to make a probabilistic threat by committing to pull the trigger exactly once—unless the mobster reveals the secret. Draw a game tree similar to Figure 13.4 for this game.
3. What happens when the corrupt policeman uses brinkmanship, playing "Russian roulette" to threaten to kill the mobster with some probability q unless he reveals the money's secret location?
4. Derive the effectiveness and acceptability conditions for this game, and determine the values for p and q for which the corrupt policeman can use a pure threat, brinkmanship, or no threat at all.
5. Suppose that $p = 1/3$. Will the corrupt policeman find it optimal to play Russian roulette with the mobster?

If so, how many bullets will he optimally load into the gun before spinning the cylinder?

3. Answer the questions from Exercise S3 for the following movies:

1. In the 1941 movie classic *The Maltese Falcon*, the hero, Sam Spade (Humphrey Bogart), is the only person who knows the location of the immensely valuable gem-studded falcon figure, and the villain, Caspar Gutman (Sydney Greenstreet), is threatening to torture him for that information. Spade points out that torture is useless unless the threat of death lies behind it, and Gutman cannot afford to kill Spade, because then the information would die with him. Therefore, he may as well not bother with the threat of torture. Gutman replies, "That is an attitude, sir, that calls for the most delicate judgment on both sides, because, as you know, sir, men are likely to forget in the heat of action where their best interests lie and let their emotions carry them away."

2. The 1925 Soviet classic *The Battleship Potemkin* (set in the summer of 1905) closes with a squadron of ships from the tsar's Black Sea fleet chasing the mutinous and rebellious crew of the *Potemkin*. The tension mounts as the ships draw ever closer. Men on each side race to their battle stations, load and aim the huge guns, and wait nervously for the order to fire on their countrymen. Neither side wants to attack the other, but neither wants to back down or to die without defending itself. The tsar's ships have orders to take the *Potemkin* by any means necessary, and the crew knows it will be tried for treason if it surrenders.

4. Answer the questions listed in Exercise S4 for these examples of successful brinkmanship:

1. The negotiations between the South African apartheid regime and the African National Congress to establish a new constitution with majority rule, 1989 to 1994.

Reading: Allister Sparks, *Tomorrow Is Another Country* (New York: Hill and Wang, 1995).

2. Peace in Northern Ireland: disarmament of the IRA in July 2005, the St. Andrews Agreement of October 2006, the elections of March 2007, and the power-sharing government of Ian Paisley and Martin McGuinness.

Reading: “The Thorny Path to Peace and Power Sharing,” *CBC News*, March 26, 2007, available at www.cbc.ca/news2/background/northern-ireland/timeline.xhtml.

5. Answer the questions listed in Exercise S4 for these examples of unsuccessful brinkmanship:

1. The U.S. budget confrontation between President Clinton and the Republican-controlled Congress in 1995. *Readings:* Sheldon Wolin, “Democracy and Counterrevolution,” *Nation*, April 22, 1996; David Bowermaster, “Meet the Mavericks,” *U.S. News and World Report*, December 25, 1995 – January 1, 1996; “A Flight that Never Seems to End,” *Economist*, December 16, 1995.

2. The television writers’ strike of 2007 – 2008. *Readings:* “Writers Guild of America,” online archive of the *New York Times* on the Writers Guild and the strike, available at http://topics.nytimes.com/top/reference/timestopics/organizations/w/writers_guild_of_america/index.xhtml; “Writers Strike: A Punch from the Picket Line,” available at <http://writers-strike.blogspot.com>.

6. Answer the questions listed in Exercise S4 for these potential opportunities for brinkmanship in the future:

1. The West Coast states (California, Oregon, and Washington) attempt to secede from the United States. *Readings:* “What If California Attempted to Secede from the U.S.?” available at <http://www.bbc.com/future/story/20190221-what-if-california-seceded-from-the-us>; “Partition and Secession in California,” available at

https://en.wikipedia.org/wiki/Partition_and_secession_in_California.

2. A nuclear confrontation between India and Pakistan over Kashmir or other issues. *Readings:* “Could the Confrontation between India and Pakistan Lead to Nuclear War?” *Pacific Standard*, <https://psmag.com/news/could-the-conflict-between-pakistan-and-india-lead-to-nuclear-war>; “Factbox: India and Pakistan—Nuclear Arsenals and Strategies,” Reuters, available at <https://www.reuters.com/article/us-india-kashmir-pakistan-nuclear-factbo/factbox-india-and-pakistan-nuclear-arsenals-and-strategies-idUSKCN1QI405>.

14 ■ Design of Incentives

JAMES MIRRLEES WON THE NOBEL PRIZE in economics in 1996 for his pioneering work on optimal nonlinear income taxation and related policy issues. Many non-economists, and some economists, too, found his work difficult to understand. But the *Economist* magazine gave a brilliant characterization of the broad importance and relevance of the work. It said that Mirrlees showed us “how to deal with someone who knows more than you do.” ¹

In [Chapter 9](#), we observed some of the ways in which asymmetries of information affect the analysis of games. But the underlying problem for Mirrlees differed slightly from the problems we considered earlier. In his work, one player (the government) needed to devise a set of rules so that the other players’ (the taxpayers’) incentives were aligned with the first player’ s goals. Models with this general framework, in which a less informed player works to create motives for a more informed player to take actions beneficial to the less informed one, now abound, and they are relevant to a wide range of social and economic interactions. Generally, the less informed player is called the *principal* while the more informed one is called the *agent*; hence these models are termed *principal - agent* models.

The principal’ s goal in a principal - agent problem is to design an *incentive scheme* that will motivate the agent to take actions that benefit the principal, taking into account that the agent knows something (about the world or about herself) that the principal does not. Such incentive schemes are known as *mechanisms*, and the process that the principal uses to devise the best possible incentive scheme is referred to as [mechanism design](#) or [incentive design](#). As we will see, mechanism-design ideas apply in many contexts that at first glance might seem to have little in common, including the

pricing of goods or services ([Section 1](#)), highway procurement ([Section 3](#)), employee supervision ([Section 5.A](#)), insurance provision ([Section 5.B](#)), auctions ([Chapter 15](#)), and Mirrlees' s original application: government taxation.

In Mirrlees' s model, the government seeks a balance between efficiency and equity. It wants the more productive members of society to contribute effort to increase its total output; it can then redistribute the proceeds to benefit the poorer members. If the government knew the exact productive potential of every person and could observe the quantity and quality of every person' s effort, it could simply order everyone to contribute according to their ability, and it could distribute the fruits of their effort according to people' s needs. But such detailed information would be costly or even impossible to obtain, and such redistribution schemes can be equally difficult to enforce. Each person has a good idea of his abilities and needs, and chooses his own effort level, but stands to benefit by concealing this information from the government. Pretending to have less ability and more need will enable him to get away with paying less in taxes or getting larger checks from the government; in addition, his incentive to provide effort is reduced if the government takes part of the yield. The government must calculate its tax policy, or design its fiscal mechanism, taking these problems of information and incentives into account. Mirrlees' s contribution was to solve this complex mechanism-design problem within the principal - agent framework.

The economist William Vickrey shared the 1996 Nobel Prize in economics with Mirrlees for his own work in mechanism design in the presence of asymmetric information. Vickrey is best known for designing an auction mechanism to elicit truthful bidding, a topic we will study in greater detail in [Chapter 15](#). But his work extended to other mechanisms, such as

congestion pricing on highways, and he and Mirrlees laid the groundwork for extensive research in the field.

Indeed, in the past 30 years, the general theory of mechanism design has made great advances. The 2007 Nobel Prize in economics was awarded to Leonid Hurwicz, Roger Myerson, and Eric Maskin for their contributions to it. Their work, and that of many others, has taken the theory and applied it to numerous specific contexts, including the design of compensation schemes, insurance policies, and of course, tax schedules and auctions. In this chapter, we develop a few prominent applications, using our usual method of numerical examples followed by exercises.

Endnotes

- “Economics Focus: Secrets and the Prize,” *Economist*, October 12, 1996. [Return to reference 1](#)

Glossary

[mechanism design](#)

Same as **incentive design**.

[incentive design](#)

The process that a *principal* uses to devise the best possible incentive scheme (or mechanism) in a *principal-agent problem* to motivate the agent to take actions that benefit the principal. By design, such incentive schemes take into account that the agent knows something (about the world or about herself) that the principal does not know. Also called **mechanism design**.

1 PRICE DISCRIMINATION

A firm generally sells to diverse customers with different levels of willingness to pay for its product. Ideally, the firm would like to extract from each customer the maximum that he would be willing to pay. If the firm could charge each customer an individualized price based on that customer's willingness to pay, economists would say that it was practicing perfect (or first-degree) [price discrimination](#).

Such perfect price discrimination may not be possible for many reasons. The most general underlying reason is that even a customer who is willing to pay a lot prefers to pay less. Therefore, the customer will prefer a lower price, and the firm may have to compete with other firms or resellers who undercut its high price. But even if there are no close competitors, the firm usually does not know how much each individual customer is willing to pay, and customers will try to get away with pretending to be unwilling to pay a high price so as to secure a lower price. In some situations, even if the firm could detect a customer's willingness to pay, it would be illegal to practice blatant first-degree price discrimination based on the identity of the buyer. In such situations, the firm must devise a product line and prices so that customers' choices of what they buy (and therefore what they pay) go some way toward the firm's goal of increasing its profit by way of price discrimination.

In our terminology of asymmetric information games developed in [Chapter 9](#), the process by which a firm identifies a customer's willingness to pay from his purchase decisions involves *screening* to achieve *separation of types* (by *self-selection*). The firm does not know each customer's *type* (willingness to pay), so it tries to acquire that information

from the customer’ s actions. An example that should be familiar to most readers is that of airlines. These firms try to separate business travelers, who are willing to pay more for their tickets, from tourists, who are not willing to pay as much, by offering low prices in return for various restrictions on fares that the business flyers are not willing to accept, such as advance purchase and minimum stay requirements.² We develop this particular example in more detail here to make the underlying ideas more precise and quantifiable.

We consider the pricing decisions of a firm called Pie-In-The-Sky (PITS), an airline running a service from Podunk to South Succotash. It carries some business passengers and some tourists; the former type are willing to pay a higher price than the latter type for any particular ticketed seat. To serve the tourists profitably without having to offer the same low price to the business passengers, PITS has to develop a way of creating different versions of the same flight; it then needs to price these options in such a way that each type will choose a different version. As mentioned above, the airline could distinguish between the two types of passengers by offering restricted and unrestricted fares. The practice of offering first-class and economy tickets is another way to distinguish between the two groups; we will use that practice as our example.

Type of service	PITS’ s cost	RESERVATION PRICE		PITS’ S POTENTIAL PROFIT	
		Tourist	Business	Tourist	Business

You may need to scroll left and right to see the full figure.

Economy	100	140	225	40	125
First	150	175	300	25	150
You may need to scroll left and right to see the full figure.					

FIGURE 14.1 Airline Price Discrimination (in Dollars)

Suppose that 30% of PITS' s customers are business travelers and 70% are tourists. The table in Figure 14.1 shows the (maximum) willingness to pay, or the *reservation price* in economics jargon, for each type of customer for each class of service, along with the costs of providing the two types of service and the potential profits available under each option.

We begin by setting up a ticket-pricing scheme that is ideal from PITS' s point of view. Suppose it knows the type of each individual customer: its salespeople determine customers' types, for example, by observing their style of dress when they come to make their reservations. Also suppose that there are no legal prohibitions on price discrimination and no possibility that lower-priced tickets can be resold to other passengers. (Actual airlines prevent such resale by requiring positive ID for each ticketed passenger.) Then PITS could practice perfect (first-degree) price discrimination.

How much would PITS charge each type of customer? It could sell a first-class ticket to each business traveler at \$300 for a profit of $\$300 - \$150 = \$150$ per ticket or sell him an economy ticket at \$225, for a profit of $\$225 - \$100 = \$125$

per ticket. The former is better for PITS, so it would want to sell \$300 first-class tickets to business travelers. It could sell a first-class ticket to each tourist at \$175, for a profit of $\$175 - \$150 = \$25$, or sell him an economy ticket at \$140, for a profit of $\$140 - \$100 = \$40$. Here, the latter is better for PITS, so it would want to sell \$140 economy tickets to the tourists. Ideally, PITS would like to sell only first-class tickets to business travelers and only economy tickets to tourists, in each case at a price equal to the relevant group's maximum willingness to pay. PITS's total profit per 100 customers from this strategy would be

$$(\$140 - \$100) \times 70 + (\$300 - \$150) \times 30 = \$40 \times 70 + \$150 \times 30 = \$2,800 + \$4,500 = \$7,300.$$

Thus, PITS's best possible outcome earns it a profit of \$7,300 for every 100 customers it serves.

Now turn to the more realistic scenario in which PITS cannot identify the type of each customer or is not allowed to use such information for purposes of overt price discrimination. How can it use the different ticket versions to screen its customers?

The first thing PITS should realize is that the pricing scheme devised above will not be the most profitable in the absence of identifying information about each customer. Most importantly, it cannot charge the business travelers the full \$300 they are willing to pay for first-class seats while charging only \$140 for an economy seat. Then the business travelers could buy economy seats, for which they are actually willing to pay \$225, for \$140 and get an extra benefit, or *consumer surplus* in the jargon of economics, of $\$225 - \$140 = \$85$. They might use this surplus, for example, for better food or accommodation on their travels. Paying the maximum \$300 that they are willing to pay for a first-class seat would leave them no consumer surplus. Therefore, they would switch to economy class in this situation, and

screening would fail. PITS' s profit per 100 customers would drop to $(140 - 100) \times 100 = \$4,000$.

The maximum that PITS will be able to charge for first-class tickets must give business travelers at least as much extra benefit as the \$85 they can get if they buy an economy ticket. Thus, the price of first-class tickets can be at most $\$300 - \$85 = \$215$. (Perhaps it should be \$214 to give business travelers a definite positive reason to choose first class, but we will ignore the trivial difference.) PITS can still charge \$140 for an economy ticket to extract as much profit as possible from the tourists, so its total profit in this case (from every 100 customers) will be

$$(140 - 100) \times 70 + (215 - 150) \times 30 = 40 \times 70 + 65 \times 30 \\ = 2,800 + 1,950 = 4,750.$$

This profit is more than the \$4,000 that PITS would get if it tried unsuccessfully to implement its perfect price discrimination scheme despite its limited information, but less than the \$7,300 it would get if it had full information and successfully practiced perfect price discrimination.

By pricing first-class seats at \$215 and economy seats at \$140, PITS can successfully screen and separate business travelers from tourists on the basis of their self-selection of its two types of services. But PITS must sacrifice some profit to achieve this indirect discrimination. PITS loses this profit because it must charge the business travelers less than their full willingness to pay. As a result, its profit per 100 passengers drops from the \$7,300 it could achieve if it had full and complete information, to the \$4,750 it achieves by indirect discrimination based on self-selection. The difference, \$2,550, is precisely 85×30 , where 85 is the drop in the first-class fare below the business travelers' full willingness to pay for this service, and 30 is the number of these business travelers per 100 passengers served.

Our analysis shows that, in order to achieve separation of types with its ticket-pricing mechanism, PITS has to keep the first-class fare sufficiently low to give the business travelers enough incentive to choose this service. Those travelers have the option of choosing economy class if it provides more benefit (or surplus) to them; PITS has to ensure that they do not “defect” to making the choice that PITS intends for the tourists. Such a requirement, or constraint, on the screener’s strategy arises in all problems of mechanism design and is called an *incentive-compatibility constraint*.

The only way PITS could charge business travelers more than \$215 without inducing their defection would be to increase the economy fare. For example, if the first-class fare were \$240 and the economy fare were \$165, then business travelers would get equal consumer surplus from either class; their surplus would be $\$300 - \240 from first class and $\$225 - \165 from economy class, or \$60 from each. At those higher prices, they would still be (only just) willing to buy first-class tickets, and PITS could enjoy higher profits from each first-class ticket sale.

But at \$140, the economy fare is already at the limit of the tourists’ willingness to pay. If PITS raised that fare to \$165, it would lose those customers altogether. In order to keep those customers willing to buy, PITS’s pricing mechanism must meet an additional requirement, namely, the tourists’ *participation constraint*.

PITS’s pricing strategy is thus squeezed between the participation constraint of the tourists and the incentive-compatibility constraint of the business travelers. If it charges X for economy and Y for first class, it must keep $X < 140$ to ensure that the tourists still buy tickets, and it must keep $225 - X < 300 - Y$, or $Y < X + 75$, to ensure that the business travelers choose first class and not economy class. Subject to these constraints, PITS wants to charge

prices that are as high as possible. Therefore, its profit-maximizing screening strategy is to make X as close to 140 and Y as close to 215 as possible. Ignoring the small differences that are needed to preserve the less-than signs, let us call the prices 140 and 215. Then charging \$215 for first-class seats and \$140 for economy seats is the solution to PITS' s mechanism-design problem.

This pricing strategy being optimal for PITS depends on the specific numbers in our example. If the proportion of business travelers were much higher—say, 50%—PITS would have to revise its optimal ticket prices. If 50% of its customers were business travelers, the sacrifice of \$85 on each business traveler' s ticket might be too high to justify keeping the few tourists. PITS might do better not to serve the tourists at all—that is, to violate the tourists' participation constraint and raise the price of first-class service. Indeed, the strategy of screening by self-selection with these percentages of passengers would yield PITS a profit, per 100 customers, of

$$(140 - 100) \times 50 + (215 - 150) \times 50 = 40 \times 50 + 65 \times 50 \\ = 2,000 + 3,250 = 5,250.$$

By contrast, the strategy of serving only business travelers in \$300 first-class seats would yield a profit (per 100 customers) of

$$(300 - 150) \times 50 = 150 \times 50 = 7,500,$$

which is higher than its profit with the screening strategy. Thus, if there are only a relatively few customers with low willingness to pay, a seller might find it better not to serve them at all than to offer sufficiently low prices to the mass of high-paying customers to prevent their switching to the lower-priced version of its product or service.

Precisely what proportion of business travelers constitutes the borderline between the two cases? We leave this question as an exercise for you. We will simply point out that an airline's decision to offer low tourist fares may be a profit-maximizing response to the existence of asymmetric information, rather than an indication of some soft spot for vacationers!

Endnotes

- Pricing policies that offer a menu of different versions of a product at different prices are referred to as *second-degree* price discrimination. Some examples include *quantity discounts*, where different quantities of a good are offered at different per-unit prices, and *upgrades* (like the example here) or *damaged goods*, where different qualities of a product are offered. See Raymond Deneckere and Preston McAfee, “Damaged Goods,” *Journal of Economics & Management Strategy*, vol. 5, no. 2 (June 1996), pp. 149 – 174. Readers with an economics background will also be familiar with *third-degree* price discrimination, also referred to as “market segmentation,” in which the firm can observe something about the buyer (such as her age, income, or location) and can charge different prices based on that information. [Return to reference 2](#)

Glossary

[price discrimination](#)

Perfect, or first-degree, price discrimination occurs when a firm charges each customer an individualized price based on willingness to pay. In general, price discrimination refers to situations in which a firm charges different prices to different customers for the same product.

2 SOME TERMINOLOGY

We have now seen one example of mechanism design in action. There are many others, of course, and we will see additional ones in later sections of this chapter. We pause briefly here, however, to set out the specifics of the terminology used in most models of this type.

Mechanism-design problems are broadly of two kinds. The airline price-discrimination case in [Section 1](#) is an example of the first kind, in which one player is better informed (in the example, the customer knows his own willingness to pay), and his information affects the payoff of the other player (in the example, the airline's pricing and therefore its profits). In the language of [Chapter 9](#), mechanisms in this category are designed to cope with potential *adverse selection*. The less informed player designs a scheme in which the more informed player must make some choice that will reveal the information, albeit at some cost to the less informed player (in the example, the airline's inability to charge the business travelers their full willingness to pay).

In the second kind of mechanism-design problem, one player takes some action that is not observable to others. For example, an employer cannot observe the quality, or sometimes even the quantity, of the effort an employee exerts, and an insurance company cannot observe all the actions that an insured driver or homeowner takes to reduce the risk of an accident or robbery. In the language of [Chapter 9](#), mechanisms for this kind of problem are designed to cope with potential *moral hazard*. The less informed player designs a scheme—for example, offering profit sharing to an employee or imposing deductibles and copayments on the insured—that aligns the other player's incentives to some extent with those of the mechanism designer.^{[3](#)}

In each case, the less informed player designs the mechanism; she is called the [principal](#) in the strategic game. The more informed player is then called the [agent](#); this term is most accurate in the case of the employee and less so in the cases of the customer or the insured, but the jargon has become established and we will adopt it. The game is then called a [principal-agent](#), or [agency, problem](#).

In both kinds of problems, the principal designs the mechanism to maximize her own payoff, subject to two types of constraints. First, the principal knows that the agent will use the mechanism to maximize his own (the agent's) payoff. In other words, the principal's mechanism has to be consistent with the agent's incentives. As we saw in [Chapter 9, Section 5.A](#), this requirement is called the *incentive-compatibility constraint*. Second, given that the agent responds to the mechanism in his own best interest, the agency relationship has to give the agent at least as much expected payoff as he would get elsewhere—for example, by working for a different employer, or by driving instead of flying. In [Chapter 9](#), we termed this requirement the *participation constraint*. We saw specific examples of both constraints in the airline price-discrimination example in the previous section; we will meet many other examples and applications in the rest of this chapter.

Endnotes

- In formal game theory, the two kinds of mechanism-design problems we describe here are often called “hidden information” and “hidden action” problems, respectively. The distinction between the two kinds of problems was emphasized by Oliver Hart and Bengt Holmstrom in their classic paper, “The Theory of Contracts,” in Truman Bewley (ed.), *Advances in Economic Theory: Fifth World Congress*, (Cambridge: Cambridge University Press, 1987), pp. 71 – 156. [Return to reference](#)
[3](#)

Glossary

principal

The principal is the less-informed player in a principal-agent game of asymmetric information. The principal in such games wants to design a mechanism that creates incentives for the more-informed player (agent) to take actions beneficial to the principal.

agent

The agent is the more-informed player in a principal-agent game of asymmetric information. The principal (less-informed) player in such games attempts to design a mechanism that aligns the agent's incentives with his own.

principal-agent (agency) problem

A situation in which the less-informed player (principal) wants to design a mechanism that creates incentives for the more-informed player (agent) to take actions beneficial to himself (the principal).

principal-agent (agency) problem

A situation in which the less-informed player (principal) wants to design a mechanism that creates incentives for the more-informed player (agent) to take actions beneficial to himself (the principal).

3 INFORMATION-REVEALING CONTRACTS

When writing procurement contracts for certain services—perhaps highway building or office-space construction—governments and firms face mechanism-design problems of the kind we have been describing. There are two common ways of writing such contracts. In a *cost-plus contract*, the buyer agrees to pay the supplier of the services a sum equal to his cost, plus a small amount of profit. In a *fixed-price contract*, a specific price for the services is agreed upon in advance; the supplier keeps any extra profit if his actual cost turns out to be less than anticipated, and he bears the loss if his actual cost is higher.

Each type of contract has its own good and bad points. The cost-plus contract appears not to give the supplier excessive profit; this characteristic is especially important for public-sector procurement contracts, where the citizens are the ones who ultimately pay for the procured services. But the supplier typically has better information about his cost than does the buyer of his services; therefore, the supplier can be tempted to overstate or pad his costs in order to extract some benefit from the wasteful excess. The fixed-price contract, in contrast, gives the supplier every incentive to keep the cost at a minimum and thus to achieve an efficient use of resources. But with this kind of public-sector contract, the citizens have to pay the set price and give away any excess profit (to the supplier). The optimal contract should balance these two considerations.

A. Highway Construction: Full Information

We first consider the example of a state government designing a procurement mechanism for a highway-construction project. Specifically, suppose that a major highway is to be built by a construction firm to be hired by the state, and that the state government has to decide how many lanes it should have.⁴ More lanes yield more social benefit in the form of faster travel and fewer accidents (at least up to a point, beyond which the harm to the countryside will be too great). To be specific, we suppose that the social value V (measured in billions of dollars) of having N lanes on the highway is given by the formula

$$V = 15N - \frac{N^2}{2}.$$

The cost of construction per lane, including an allowance for normal profit, could be either \$3 billion or \$5 billion per lane, depending on the types of soil and minerals located in the construction zone. For now, we assume that the state government can determine the construction cost as well as the firm can. So it chooses N and writes a contract to maximize the benefit to the state (V) net of the fee paid to the firm (call it F); that is, the government's objective is to maximize net benefit, G , where $G = V - F$.⁵

Suppose first that the government knows that the actual cost per lane is 3 (billion dollars per lane of highway). At this cost level, the government has to pay $3N$ to the firm for an N -lane highway. The government then chooses N to maximize G , as above, where the appropriate formula in this situation is

$$G = V - F = 15N - \frac{N^2}{2} - 3N = 12N - \frac{N^2}{2}.$$

Recall that in the appendix to [Chapter 5](#), we gave a formula for finding the correct value to maximize this type of function. Specifically, the solution to the problem of choosing X to maximize

$$Y = A + BX - CX^2$$

is $X = B/(2C)$. Here, Y is G , X is N , and $A = 0$, $B = 12$, and $C = \frac{1}{2}$. Applying our solution formula yields the government's optimal choice of $N = 12/(2 \times \frac{1}{2}) = 12$. The best highway to choose therefore has 12 lanes, and the cost of that 12-lane highway is \$36 billion. So, the government offers the contract: "Build a 12-lane highway and we will pay you \$36 billion." ⁶ This price includes normal profit, so the firm is happy to take the contract.

Similarly, if the cost is \$5 billion per lane, the optimal N will be 10. The government will offer a \$50 billion contract for the 10-lane highway. And the firm will accept the contract.

B. Highway Construction: Asymmetric Information

Now suppose that the firm knows how to assess the relevant terrain to determine the actual building cost per lane, but the state government does not. The government can only estimate what the cost will be. We assume that it thinks that there is a two-thirds probability of the cost being 3 (billion dollars per lane) and a one-third probability of the cost being 5.

What if the government tries to go ahead with the full-information optimum we found in [Section 3.A](#) and offers a pair of contracts: “Build a 12-lane highway for \$36 billion” and “Build a 10-lane highway for \$50 billion”? If the actual cost is really only \$3 billion per lane, the firm will get more profit by taking the latter contract, even though that one was designed for the situation in which the cost is \$5 billion per lane. The true cost of the 10-lane highway will be only \$30 billion, and the firm will earn \$20 billion in excess profit.⁷

This outcome is not very satisfying for the state government. The contracts offered do not give the firm sufficient incentive to choose between them on the basis of cost; it will always take the \$50 billion contract. There must be a better way for the government to design its procurement contract system.

Let us therefore allow the government the freedom to design an optimal screening mechanism to separate the two possible types of projects (low-cost and high-cost). Suppose it offers a pair of contracts: “Contract L: Build N_L lanes and get paid R_L dollars” and “Contract H: Build N_H lanes and get paid R_H dollars.” If contracts L and H are designed correctly then, when the firm determines the true cost is low (\$3 billion per lane), it will pick contract L (L stands for “low”) and, when it determines the true cost is high (\$5 billion per lane), it will pick contract H (H stands for “high”). The numbers that the symbols

N_L , R_L , N_H , and R_H represent must satisfy certain conditions for this screening mechanism to work.

First, under each contract, a firm anticipating the relevant cost (low for contract L and high for contract H) must receive enough payment to cover its cost (inclusive of normal profit). Otherwise, it will not agree to the terms; it will not participate in the contract. Thus, the contract must satisfy two *participation constraints*: $3N_L \leq R_L$ for the firm when the cost is 3, and $5N_H \leq R_H$ for the firm when the cost is 5.

Next, the government needs the two contracts to be such that the firm would not benefit by taking contract H when it knows the true cost is low, and vice versa. That is, the contracts must also satisfy two incentive-compatibility constraints. For example, if the true cost is low, contract L will yield excess profit $R_L - 3N_L$, whereas contract H will yield $R_H - 3N_H$. (Note that in the latter expression, the number of lanes and the payment are as specified in the H contract, but the firm's cost is still only 3, not 5.) To be incentive-compatible for the low-cost case, the contracts must ensure that when the cost is 3 the firm's excess profit from contract H is no larger than that from contract L, or that the latter expression no larger than the former. Thus, we need $R_L - 3N_L \geq R_H - 3N_H$. Similarly, if the true cost is 5, the firm's excess profit from the L contract must be no larger than its excess profit from the H contract. To keep the contracts incentive-compatible, we therefore need $R_H - 5N_H \geq R_L - 5N_L$.

As before, the government wants to maximize the net benefit of the payment, G . Here, there are two possible outcomes, so the government must actually maximize the net *expected* benefit using the probabilities of the two project types as weights to calculate the expected value of G . Therefore, the government's objective here is to maximize

$$E[G] = \left(\frac{2}{3}\right)\left[15N_L - \frac{(N_L)^2}{2} - R_L\right] + \left(\frac{1}{3}\right)\left[15N_H - \frac{(N_H)^2}{2} - R_H\right].$$

The problem looks formidable, with four choice variables and four inequality constraints. But it simplifies greatly, because two of the constraints are redundant, and the other two must hold as exact equalities (or, the other two constraints *bind*), allowing us to solve and substitute for two of the variables.

Note that if the participation constraint when the true cost is high, $5N_H \leq R_H$, and the incentive-compatibility constraint when the true cost is low, $R_L - 3N_L \geq R_H - 3N_H$, both hold, then we can get the following string of inequalities (using the fact that N_H cannot be negative):

$$R_L - 3N_L \geq R_H - 3N_H \geq 5N_H - 3N_H \geq 2N_H \geq 0.$$

The first and last expressions in the inequality string tell us that $R_L - 3N_L \geq 0$. Therefore, we need not consider the participation constraint when the true cost is low, $3N_L \leq R_L$, separately; it is automatically satisfied when the two other constraints are satisfied.

It is also intuitive that when the true cost is high, the firm will not want to pretend that the cost is low; it would be compensated for the lower cost project while incurring the higher true cost. However, this intuition needs to be verified by the rigorous logic of the analysis. To confirm the intuition, we ignore the second incentive-compatibility constraint, $R_H - 5N_H \geq R_L - 5N_L$, and proceed to solve the problem with just the remaining two constraints. Then we return and verify that our solution to the two-constraint problem satisfies the ignored incentive-compatibility constraint anyway. So our solution must also be the solution to the three-constraint problem. (If something better was available, it would also work better for the less constrained problem.)

Thus, we have two constraints to consider: $5N_H \leq R_H$ (the participation constraint for the firm when the true project cost is high) and $R_L - 3N_L \geq R_H - 3N_H$ (the incentive-compatibility constraint for the firm when the true project cost is low). We write these as $R_H \geq 5N_H$ and $R_L \geq R_H + 3(N_L - N_H)$. Next, observe

that R_L and R_H each figure negatively in $E[G]$, the expected benefit to the government: It wants to make them as small as possible while still satisfying the constraints. This result is achieved by satisfying each constraint with equality; thus, we can set $R_H = 5N_H$ and $R_L = R_H + 3(N_L - N_H) = 3N_L + 2N_H$. These expressions for the contract payments can now be substituted into the expected benefit function, $E[G]$. This substitution yields

$$\begin{aligned} E[G] &= \left(\frac{2}{3}\right) \left[15N_L - \frac{(N_L)^2}{2} - 3N_L - 2N_H \right] + \left(\frac{1}{3}\right) \left[15N_H - \frac{(N_H)^2}{2} - 5N_H \right] \\ &= 8N_L - \frac{(N_L)^2}{3} + 2N_H - \frac{(N_H)^2}{6}. \end{aligned}$$

The expected benefit function now splits cleanly into two parts: One (the first two terms) involves only N_L , and the other (the second two terms) involves only N_H . We can apply our maximization formula (from [Section 3.A](#)) separately to each part. In the N_L part, $A = 0$, $B = 8$, and $C = 1/3$, so the optimal $N_L = 8/(2 \times 1/3) = 24/2 = 12$. In the N_H part, $A = 0$ again, $B = 2$, and $C = 1/6$, so the optimal $N_H = 2/(2 \times 1/6) = 12/2 = 6$.

Now we can use the optimal values for N_L and N_H to derive the optimal payment (R) values, using the formulas for R_L and R_H that we derived and substituted into $E[G]$ in the previous paragraph. Substituting $N_L = 12$ and $N_H = 6$ into those formulas gives us $R_H = 5N_H = 5 \times 6 = 30$ and $R_L = R_H + 3(N_L - N_H) = 3 \times 12 + 2 \times 6 = 48$. These calculations give us the optimal values for all of the unknowns in the government's expected benefit function.

However, remember that we ignored one of the incentive-compatibility constraints. We must ensure that the ignored constraint, $R_H - 5N_H \geq R_L - 5N_L$, holds with our calculated values for the R s and the N s. In fact, it does. The left-hand side of the expression is $30 - 5 \times 6 = 0$. And the right-hand side is $48 - 5 \times 12 = -12$, so the ignored constraint is indeed satisfied.

Our solution indicates that the government should offer the following two contracts: “Contract L: Build 12 lanes and get paid \$48 billion” and “Contract H: Build 6 lanes and get paid \$30 billion.” How can we interpret this solution so as best to understand the intuition behind it? That intuition is most easily seen when we compare the solution here with the one we found in [Section 3.A](#), when the government had full information about project cost. Figure 14.2 shows these comparisons.

The optimal mechanism with asymmetric information differs in two important respects from the one we found when information was perfect. First, although the contract that the government intends the firm to choose if the project cost is low has the same number of lanes (12) as in the full-information case, its payment to the firm is larger in the asymmetric case (48 instead of 36). Second, the contract that the government intends the firm to choose if the project cost is high has a smaller number of lanes (6 instead of 10), but pays the full cost for that number ($30 = 6 \times 5$) and no more. Both of these differences separate the project types.

With asymmetric information, the firm may be tempted to pretend that the true cost of the project is high when it is in fact low. The optimal procurement mechanism therefore incorporates both a “carrot” to reward the firm for truthfully admitting the true low cost and a “stick” to dissuade it from pretending that the true cost is high. The carrot is the excess profit, $48 - 36 = 12$, that comes from the admission the firm makes implicitly by choosing contract L. The stick is the reduction in excess profit it incurs by choosing contract H, achieved by reducing the number of lanes that will be constructed in that case. The full-information high-cost contract would have the highway be 10 lanes and would pay \$50 billion; if the firm chose this contract while knowing the true cost was low would make excess profit of $50 - 3 \times 10 = \$20$ billion. In the optimal information-constrained high-cost contract, only 6 lanes are constructed, and the firm is paid \$30 billion. If the true cost is low, it makes an excess profit of $30 - 3 \times 6 = \$12$ billion by choosing the high-cost contract. Its benefit from the pretense of high cost (implicit in its choice of contract H even though the true cost is low) is reduced. In fact, the excess profit is reduced exactly to the

amount that the firm is guaranteed by the carrot part of the mechanism, thereby exactly offsetting its temptation to pretend that the true cost of the project is high.

	N_L	R_L	N_H	R_L
Perfect information	12	36	10	50
Asymmetric information	12	48	6	30

FIGURE 14.2 Highway-Building Contract Values

Endnotes

- Generally, numerous contractors would be competing for the highway-construction contract. For this example, we restrict ourselves to the case in which there is only one contractor. [Return to reference 4](#)
- In reality, the cost per lane would not have only two discrete values, but could take any value along a continuous range of possibilities. The probabilities of each value would then correspondingly form a density function on this range. Our methods will not always yield an integer solution, N , for each possible cost along this range. But we leave these matters to more advanced treatments and confine ourselves to this simple illustrative example. [Return to reference 5](#)
- In reality, the contract will contain many clauses specifying quality, timing, inspections, and so forth. We leave out these details to keep the exposition of the basic idea of mechanism design simple. [Return to reference 6](#)
- If multiple contractors are competing for the job, the ones not selected might spill the beans about its true cost. But for large highway projects (as for many other large government projects, such as defense contracts), there are often only a few potential contractors, and they do better by colluding among themselves and not revealing their private information. For simplicity, we keep the analysis confined to the case where there is just one contractor. [Return to reference 7](#)

4 EVIDENCE CONCERNING INFORMATION-REVELATION MECHANISMS

In the cases considered so far in this chapter, the agent has some private information, which we called that player's *type* in [Chapter 9](#). Further, the principal designs a mechanism that requires the agent to take some action that reveals this information. In the terminology of [Chapter 9](#), these mechanisms are examples of screening for the separation of types by self-selection.

Price-discrimination mechanisms are ubiquitous. All firms have customers who are diverse in their willingness to pay for the firms' products. Ideally, firms would like to discriminate by giving a price break to the less willing customers without giving the same break to the more willing ones. The ability of a firm to practice price discrimination may be limited for reasons other than those of information, including anti-discrimination laws, competition from other firms, or resale by initial buyers. But here we focus on information-based examples of price discrimination, keeping other reasons in the background of the discussion.

Your local coffee shop probably has a “frequent-drinker card” ; for every ten cups you buy, for example, you get one free. Why is it in the firm's interest to do this? Frequent drinkers are more likely to be locals, who have the time and incentive to search out the best deals in the neighborhood. To attract those customers away from other competing coffee shops, your coffee shop must offer a sufficiently attractive price. In contrast, infrequent customers are more likely to be strangers in the town, or in a hurry, and have less time

and incentive to search for the best deals; when they need a cup of coffee and see a coffee shop, they are willing to pay whatever the price is (within reason). So posting a higher price and giving out frequent-drinker cards enables your coffee shop to give a price break to the price-sensitive regular customers without giving the same price break to the occasional buyers. If you don't have the card, you are revealing yourself as the latter type, willing to pay a higher price.

Many restaurants offer fixed-price three-course menus or blue-plate specials as well as regular à la carte offerings. This strategy enables them to separate diverse customer types with different tastes for soups, salads, main courses, desserts, and so on. Similarly, book publishers start selling new books in a hardcover version and issue a paperback version a year or more later. The price difference between the two versions is generally far greater than the difference in the cost of production of the two kinds of books. The idea behind this pricing scheme is to separate two types of customers, those who need or want to read the book immediately and are willing to pay more for the privilege, and those who are willing to wait until they can get a better price.

The rise of the Internet and our expanding online lives make some sorts of price discrimination more difficult, while also creating new opportunities to target consumers based on extensive, extremely personalized “big data.” For products that can be easily found through an online search, like books or consumer electronics, the Internet can make it harder for sellers to use price discrimination.⁸ But for highly customized products and services, where there is little, if any, competition, sellers can use everything they have learned about you (including past purchases, browsing history, social networks, physical location, and so on)⁹ to offer experiences and set prices that are entirely unique to

you. To take an extreme example, consider virtual reality. Everything you say and do within a virtual world (including any virtual money you might earn) is visible to the company running that world, allowing them to know literally everything about your virtual self—and to charge you accordingly.

We invite you to look for other examples of price discrimination and similar screening mechanisms in your own experience. A good source of examples is Tim Harford's *Undercover Economist*.^{[10](#)}

There is a lot of research literature on the design of procurement mechanisms of the kind we sketched in [Section 3.11](#). These mechanism-design problems pertain to situations where the buyer confronts just one potential seller, whose cost is private information. This type of interaction accurately describes how contracts for major defense weapon systems or very specialized equipment are designed, as there is usually only one reliable supplier of such products or services. However, in reality, buyers often have the choice of several suppliers, and mechanisms that set the suppliers in competition with one another are beneficial to the buyer. Many such mechanisms take the form of auctions. For example, construction contracts are often awarded by inviting bids and choosing the bidder that offers to do the job for the lowest price (after adjusting for the promised quality of the work, the speed of completion, or other relevant attributes of the bid). We discuss auctions in depth in [Chapter 15](#).

Endnotes

- A 2000 study found substantial price dispersion online for commodity products such as books and CDs, but when its authors restricted attention to large online retailers, there was less variation of prices online than in traditional large brick-and-mortar stores. (This study also found that prices on the Internet were, on average, 9%–16% lower than those in stores.) See Erik Brynjolfsson and Michael D. Smith, “Frictionless Commerce? A Comparison of Internet and Conventional Retailers,” *Management Science*, vol. 46, no. 4 (April 2000), pp. 563–85. [Return to reference 8](#)
- Geofeedia, a social-media surveillance company launched in 2011, scours social media to determine people’s locations and then sells that location-based data to businesses, police departments, and others with an interest in knowing where people are at any given time. See Lee Fang, “The CIA is Investing in Firms that Mine Your Tweets and Instagram Photos,” *The Intercept*, April 14, 2016 (available at <https://theintercept.com/2016/04/14/in-undisclosed-cia-investments-social-media-mining-looms-large>, accessed May 1, 2019) and Colin Lecher and Russell Brandom, “Facebook Caught an Office Intruder Using the Controversial Surveillance Tool It Just Blocked,” *The Verge*, October 19, 2016 (available at <https://www.theverge.com/2016/10/19/13317890/facebook-geofeedia-social-media-tracking-tool-mark-zuckerberg-office-intruder>, accessed May 1, 2019). [Return to reference 9](#)
- Tim Harford, *The Undercover Economist: Exposing Why the Rich Are Rich, the Poor Are Poor—and Why You Can Never Buy a Decent Used Car!* (New York: Oxford University

Press, 2005). The first two chapters give examples of pricing mechanisms. [Return to reference 10](#)

- Jean-Jacques Laffont and Jean Tirole, *A Theory of Incentives in Procurement and Regulation* (Cambridge, Mass.: MIT Press, 1993), is the classic of this literature. [Return to reference 11](#)

5 INCENTIVES FOR EFFORT: THE SIMPLEST CASE

We now turn from the first type of mechanism-design problem, in which the principal's goal is to achieve information revelation, to the second type, which deals with moral hazard. The principal's goal in such situations is to provide an incentive that will induce the best level of effort from the agent, even though that effort level is not observable by the principal.

A. Managerial Supervision

Suppose you are the owner of a company that is undertaking a new project. You have to hire a manager to supervise it. The success of the project is uncertain, but good supervision can increase the probability of success. Managers are only human, though; they will try to get away with as little effort as they can! If your manager's effort is observable, you can write a contract that compensates the manager for his trouble sufficiently to bring forth good supervisory effort.¹² But if you cannot observe the effort, you have to try to give him incentives based on success of the project—for example, a bonus if the project is successful. Unless good effort absolutely guarantees success, however, such bonuses make the manager's income uncertain; he gets no bonus if the project fails, even if he has exerted the required effort. And the manager is likely to be averse to this risk or loss, so you have to compensate him for it. You have to design your compensation policy to maximize your own expected profit, recognizing that the manager's choice of effort depends on the nature and amount of the compensation. The solution to this incentive-design problem is intended to cope with the moral-hazard problem of the manager's shirking.

Let us consider a numerical example. Suppose that if the project succeeds, it will earn the company a profit of \$1 million over material and wage costs. If it fails, the profit will be zero. With good supervision, the probability of success is one-half, but if supervision is poor, the probability of success is only one-quarter. The manager you want to hire is currently in a steady job elsewhere that gets him \$100,000; to get him to accept a job with your firm, you must pay him at least as much, but the extra effort costs him (in terms of the extra time diverted from family, friends, or other pursuits) an equivalent of \$50,000.

In an ideal world where effort is observable, you can write a contract that states, “If you work for me, and if you exert extra effort, I will pay you \$150,000; but if you don’t, I will pay you only \$100,000.” Your expected profit (in millions of dollars) when the manager exerts the extra effort will be

$$0.5 \times 1 + 0.5 \times 0 - 0.15 = 0.35$$

Without the payment for extra effort, the manager will shirk, and your expected profit will be

$$0.25 \times 1 + 0.75 \times 0 - 0.1 = 0.15.$$

So you prefer to include the extra effort clause, and the manager is satisfied, too.

What if the manager’s effort cannot be observed? Then the incentive must be based on something that can be observed, and the only possibility in our example is success or failure of the project. Suppose you pay the manager a salary s , plus a bonus b if the project succeeds. Now, extra effort will get the manager $s + 0.5b - 0.05$; without that effort, he will get $s + 0.25b$. Thus, to elicit the extra effort, you must choose the bonus to satisfy $s + 0.5b - 0.05 \geq s + 0.25b$, or $0.25b \geq 0.05$; that is, $b \geq 0.2$. The bonus for success has to be \$200,000! That may seem a hefty sum, but observe that the extra effort will increase the probability of the manager’s receiving the bonus only from 0.25 to 0.5—that is, by 0.25—and that this extra probability will raise his expected income by just 0.25 times \$200,000—that is, by \$50,000, exactly enough to compensate him for the cost of the effort.

You don’t want to pay the manager any more than you have to. So you want his expected earnings with the extra effort—namely, $s + 0.5b$ —to equal his earnings in his current job plus the money-equivalent cost of the extra effort, or $0.10 +$

$0.05 = 0.15$. Using $b = 0.2$, we have $s = 0.15 - 0.5 \times 2 = 0.05$; that is, \$50,000. You are now offering the manager a lower wage, but a large enough bonus for success: He will get \$50,000 if the project fails, but \$250,000 if it succeeds.

However, this situation creates a risk for the manager. He probably dislikes the prospect of facing that risk. Our discussion of insurance in [Chapter 9](#) showed that people will pay to avoid risk, or have to be paid to bear it; here, to attract the manager to this job, you, as the owner, will have to compensate him for accepting the risk. Even more importantly, if the project fails, the manager will get only \$50,000, less than he was earning in his old job. Research in psychology and behavioral economics has shown that people are especially averse to losses measured in relation to the status quo. The manager, recognizing the prospect of such a loss, will demand sufficient compensation to take the job with your company.

Suppose you pay the manager $x > 0.1$ if the project succeeds, but $y < 0.1$ if it fails. The manager values y at less than its monetary value because of his loss aversion. To keep the calculations simple, we suppose that losses get twice as much weight as gains; for example, that the manager views getting \$90,000 after losing \$10,000 as equivalent to getting \$80,000 if loss were not an issue. To reflect this mathematically, suppose the manager regards y as equivalent to z where the loss $0.1 - z = 2 \times (0.1 - y)$, or $z = 2y - 0.1$. Now we can recalculate the contract you would need to offer. The manager expects $0.5 \times x + 0.5 \times z$ if the project succeeds and $0.25 \times x + 0.75 \times z$ if it fails. So, to bring forth his extra effort, you need $0.5x + 0.5z - 0.05 \geq 0.25x + 0.75z$, or $x \geq z + 0.2$, or $x \geq 2y - 0.1 + 0.2 = 2y + 0.1$, or $x - 2y \geq 0.1$. This is the incentive-compatibility constraint.

You also need to pay the manager enough, when he exerts the extra effort, to ensure that he receives compensation to

cover his previous salary plus the cost of that effort: $0.5x + 0.5z \geq 0.1 + 0.05$, or $x + z \geq 0.3$, or $x + 2y - 0.1 \geq 0.3$, or $x + 2y \geq 0.4$. This is the participation constraint.

Adding the two constraints, we have $2x \geq 0.5$, or $x \geq 0.25$. The cheapest way to satisfy the two constraints is then to set $x = 0.25$ and $y = 0.075$; that is, to pay the manager \$75,000 if the project fails and \$250,000 if it succeeds.

Note that, compared with the contract that does account for loss aversion, you must pay the manager more after failure (\$75,000 rather than \$50,000), but the same after success (\$250,000). The reason for the difference is that the manager views getting \$75,000 *after losing \$25,000* as equivalent to getting \$50,000 without any loss. In this way, the manager's loss aversion imposes an extra cost on you.

B. Insurance Provision

Suppose you own a precious pearl necklace worth \$100,000. Your parents taught you to be careful with such prized possessions, so you make a habit of storing the necklace in a bank safe-deposit vault when you are not wearing it, and you are constantly on guard against potential thieves in your vicinity whenever you wear it. Given all these precautions, your probability of losing the necklace is only 1% per year, compared with 6% if you took no precautions. So your expected monetary loss is only \$1,000, but you are averse to the risk and the loss, and you would be willing to pay much more—say, \$3,000—to insure the necklace. Insurance companies pool the independent risks of many customers like you, as we explained in [Chapter 9](#). The premium for this coverage in the market could therefore be much less than your willingness to pay—say, \$1,500—which is a 50% margin above what the company expects to pay out for lost necklaces.

Once you have the coverage, however, your incentive to be careful with the necklace is reduced. It costs you time and effort to make all those trips back and forth to the bank, and the anxiety of being constantly on guard when wearing it detracts from your enjoyment of the social occasions when you want to look your best and most relaxed. Suppose that all those inconveniences add up to the equivalent of \$500 in extra cost to you. So long as you remain uninsured, you will naturally prefer to be careful, since the extra \$500 cost is much less than the \$5,000 saved on average by reducing the probability of loss from 6% to 1%. However, once insurance has got you covered, why bother?

Now, think again about the insurance company's perspective on this situation. Once you are no longer careful, the probability of loss increases to 6%, and the insurance

company will have to pay out \$6,000, on average, giving it an expected loss of \$4,500 (before other expenses) if it charges the lower \$1,500 premium we calculated earlier. Consequently, the insurance company cannot afford to insure the necklace for less than \$6,000, but at that price, you would rather not buy insurance at all. [13](#)

The insurance contract could offer you coverage conditional on your being careful with the necklace, if your care was observable. Perhaps the insurance company could get documented evidence of your trips to the bank. But in most situations, your caution would be hard to document. Some other solution has to be found.

The usual solution to this problem is to offer partial insurance coverage, which leaves you to bear part of the loss and so gives you sufficient incentive to avoid it. The insurance contract might include a deductible (e.g., stating that you are responsible for the first \$40,000 of losses, while the company covers the rest) or co-insurance (e.g., stating that the company covers only 60% of the loss). For a single item of known value, a deductible and a co-insurance payment amount to the same thing. More generally, they can work differently, but we must leave those differences to more advanced treatments of the subject.

Consider a partial insurance contract like the one just described. The premium for 60% coverage will be proportionately smaller than that for full coverage ($0.6 \times \$1,500 = \900), and your willingness to pay for such coverage will also be proportionally smaller ($0.6 \times \$3,000 = \$2,400$). Now, you have two kinds of choices: (i) whether to buy the partial coverage at a cost of \$900 and (ii) whether to be careful. If you buy the partial insurance and choose to be careless, you will bear \$40,000 of the loss 6% of the time, for an expected monetary loss of \$2,400. Because being careful costs you only \$500, you will clearly choose to be

careful. Anticipating this, the insurance company can predict that it will only have to pay out \$60,000 with probability 1%. Overall, then, both you and the insurance company benefit from agreeing to a partial insurance contract. You get coverage that is worth \$2,400 to you at a cost of only \$900, while the insurance company gets \$900 but has to pay out only \$600 on average.

Endnotes

- Most importantly, if a dispute arises, you or the manager must be able to prove to a third party, such as an arbitrator or a court, whether the manager made the stipulated effort or shirked. This requirement, often called *verifiability*, is more stringent than mere observability by the parties to the contract (you and the manager). We intend such public observability or verifiability when we use the more common term *observability*. [Return to reference 12](#)
- Presumably, you can lower your loss probability to 1% at a cost of \$500 and you are willing to pay \$3,000 to insure yourself completely against that remaining 1% risk. So, you cannot be willing to pay more than \$3,500 for insurance. [Return to reference 13](#)

6 INCENTIVES FOR EFFORT: EVIDENCE AND EXTENSIONS

The theme of the managerial-effort incentive scheme of [Section 5.A](#) was the trade-off between giving the manager a more powerful incentive to exert extra effort and requiring him to bear more of the risk of the firm's project. This trade-off is an important consideration in practice, but it must be considered in combination with other features of the relationship between a firm and its employee. The firm has many employees, and the overall outcome for the firm depends on some combination of their actions. Most firms have multiple outputs, and each employee performs multiple tasks. Consequently, it may not be possible to describe the quality and quantity of effort simply as “good” or “bad,” or outcomes simply as “success” or “failure.”

Moreover, the firm and its employees interact over a long time, and they work together on many projects. All of these features require correspondingly more complex incentive schemes. In this section, we outline a few such schemes and refer you to a rich body of literature for further details.¹⁴ The mathematics of these schemes gets correspondingly complex, so we merely give you the intuitions behind them and leave formal rigorous analyses to more advanced courses.

A. Nonlinear Incentive Schemes

Suppose that a project has three possible outcomes—failure, modest success, and huge success—and that an agent's level of effort affects the likelihood of each outcome. What sort of bonus should an employer pay to encourage good effort, and how should that bonus depend on the extent of success? If the size of the bonus is proportional to the size of the success (as measured by sales, profit, or some other metric), the resulting wage schedule is referred to as a linear incentive scheme. For instance, suppose that a salesperson receives base salary s and commission c on each sale, and that we use x to denote the number of sales. The salesperson's total wage will be $w = s + cx$, which increases linearly with x . Of course, there are many other *nonlinear* ways to motivate an agent.

To explore the question of optimal nonlinear incentive design, we return to the managerial supervision example from [Section 5.A](#), but now with three possible outcomes: zero profit (failure), \$500,000 profit (modest success), or \$1 million profit (huge success). If the manager works hard, the probability of failure is one-sixth, the probability of modest success is one-third, and the probability of huge success is one-half. If the manager shirks, these probabilities are reversed to one-half, one-third, and one-sixth, respectively.

Suppose you pay the manager a base salary s , a bonus m for modest success, and a bonus h for huge success. His expected income is $s + m/3 + h/2$ if he exerts good supervisory effort, and $s + m/3 + h/6$ if he does not. Since the money-equivalent cost of his good effort is \$50,000, or 0.05 (million dollars), the incentive-compatibility condition to induce him to exert good effort becomes $s + m/3 + h/2 > s + m/3 + h/6 + 0.05$, or $h > 0.15$. This condition does not help us fix s and m separately, only the term $s + m/3$. The rest of the solution must come from the participation condition and the manager's loss aversion. We omit those details, but you should already be able to see that the incentive scheme focuses on the outcome of huge success; it need

not be concerned with modest success, and the bonus need not be proportional to the owner's profit. What happens here is that the manager's effort shifts some probability from failure to huge success, leaving the probability of modest success unchanged at $\frac{1}{3}$. This example is, of course, a special case, but the point is that the optimal incentive design will depend on how effort changes the probabilities of various outcomes.

Nonlinear incentive schemes like this one are ubiquitous in practice. The most common form incorporates a stipulated, fixed bonus that is paid if a certain performance standard or quota is achieved. Such a *quota-bonus scheme* can create a powerful incentive if the quota can be set at such a level that an increase in the worker's effort substantially increases the probability of meeting it. As an illustration, consider a firm that wants each salesperson to produce \$1 million in sales and is willing to pay up to \$100,000 for this level of performance. If it pays a flat 10% commission, each salesperson's incremental effort in pushing sales from \$900,000 to \$1 million will bring him \$10,000. But if the firm offers a wage of \$60,000 and a bonus of \$40,000 for meeting the quota of \$1 million, then this last bit of effort will earn him \$40,000. Thus, the quota gives the salesperson a much stronger incentive to make the incremental effort.

But the quota-bonus scheme is not without its drawbacks. The level at which the quota is set must be judged quite precisely. Suppose the firm misjudges and sets the quota at \$1.2 million, and the salesperson knows that the probability of reaching that level of sales, even with superhuman effort, is quite small. The salesperson may then give up, make very little effort, and settle for earning just the base salary. The salesperson's resulting sales may fall far short of even \$1 million. (Conversely, the pure quota-bonus scheme we outlined above gives him no incentive to go beyond the \$1 million level.) Finally, the quota must be applied over a specific period, usually a calendar year. This requirement produces even more perverse incentives. A salesperson who has bad luck in the first few months of a year may realize that he has no chance of making his quota and take things easy for the rest of the year. If, in contrast, he has very good luck

and meets the quota by July, again, he has no incentive to exert himself for the rest of the year. Or he may be able to manipulate the scheme by conspiring with his customers to shift sales from one year to another to improve his chances of making the quota in both years. A linear scheme is less open to such manipulation.

Therefore, firms usually combine a quota-bonus scheme with a more graduated linear incentive scheme. For example, a salesperson may get a base salary, a low rate of commission for sales between \$500,000 and \$1 million, a higher rate of commission for sales between \$1 million and \$2 million, and so on. Managers of mutual funds, for example, are rewarded for good performance over a calendar year. Their rewards come from their firms in the form of bonuses, but also from the public when they invest more money in those managers' specific funds.