# 3 Behavioral Economics

## 3.1 Introduction

Behavioral economics involves both the introduction and the development within economic theory of insights about behavior drawn from different domains of psychology. It has been widely recognized for a long time (since at least the early contributions of Allais and Ellsberg) that some of the central assumptions of standard economic analysis may reflect an unrealistic representation of human behavior. To be sure, the standard approach can produce wrong predictions on some occasions. Empirical and experimental research has cataloged a wide range of anomalies: observed choices that do not match the predictions of standard analysis. Behavioral economics is an attempt to improve the predictive power of economics by building in more realistic descriptions of individual behavior. It is useful for behavioral analysis to view people as departing from standard economic behavior in three distinct ways: limited rationality, limited will power, and limited self-interest.

In this chapter we consider public intervention motivated by the fact that people can make invalid choices or mistakes. This form of public intervention is aimed at protecting people against the consequences of their own decisions. There is a long history of opposition to any form of paternalistic public intervention where the state claims to know better what is "truly" good for people than they know for themselves. The nature of paternalism is frequently illustrated by the analogy to the interference of parents in their children's choices and the limitation of those choices. We may readily accept that parents can limit the freedom of choices of their children and make some choices on their behalf, but the same argument is rarely accepted with adults and state intervention. Libertarians are the most opposed to that kind of public intervention by asking "are you too incompetent to know what's best for yourself?" Traditional (public) economics assumes that people know what is best for themselves and they can get themselves to act according to their own best interest. In that context there is no need for public intervention beyond redistribution and the correction of market failure.

To provide the motivation for the approach we take, it is necessary to clarify some basic issues. We can distinguish three forms of public policies: paternalism, welfarism, and behavioralism. *Paternalism* refers to policies aimed at benefiting individuals who cannot be relied upon to pursue self-interest. It is most appropriate for children and others with behavioral disorders deemed unable to make rational choice. *Welfarism*, as

we have noted many times already, denotes government policy aimed at resolving the trade-off between efficiency and equity in order to maximize a social objective function. Improving efficiency is uncontentious, since it implies a Pareto improvement and as such will be unanimously supported. In contrast, attaining aims of equity implies redistribution, so there will be some consumers that lose from this (but others that gain). Welfarism represents policies intended to benefit individuals when self-interest (broadly defined) cannot be relied upon because of the presence of market failures or the need for redistribution. *Behavioralism* represents policies intended to benefit individuals when self-interest (again broadly defined) cannot be relied upon because of the presence of internal conflicts (*internalities*). The central idea is that if people make systematic mistakes or biased choices they may regret later on, public policy can manipulate those mistakes and biases that hurt people (1) to help them to protect themselves and at the same time (2) to respect their autonomy of choice.

There are two types of mistakes people can make and that can motivate public intervention. First, people do not know what is best for themselves through a mere *lack of knowledge* of the needed information. Suppose that people cannot easily distinguish healthy and nonhealthy food. Then a policy enforcing effective "healthy" food labeling is beneficial to individuals (if not to the food industry). Second, people do know what is best for themselves, but they cannot get themselves to act in the correct way because of a *lack of self-discipline*. For example, I may plan to stop smoking tomorrow or plan to start jogging tomorrow, but when tomorrow arrives I cannot get myself to do it. Public policy aimed at helping me to do it will be beneficial.

The purpose of this chapter is to describe how limited rationality, which is the possibility that individuals can make biased choices or invalid choices, opens up the scope for a wide range of policies intended to benefit individuals who cannot be relied upon because of internal conflicts. The policies must help people help themselves, and at the same time the policies must respect their autonomy of choice. In that sense the policies are not paternalistic. We will start by describing several examples of invalid choices. We will then describe policies that can effectively help people help themselves to do what is best for them. We will subsequently discuss the welfare evaluation of those policies, and more important, what welfare criteria to adopt when revealed preferences through individual choices cannot be relied upon to make welfare judgments, since people can make invalid choices that do not represent their "true" preferences. The final section of the chapter discusses another important deviation from the standard model, namely when people are not selfish optimizing agents but display concern over the material consumption of other agents. This section reflects the general idea that people do care about giving and receiving fair treatment in a wide set of circumstances.

## 3.2   Behavioral Individuals

In chapter 2 we described an economy in which people had only economic motives and were also fully rational. Those restrictions in the standard approach to public economics are useful as they impose a research structure and discipline in all analysis of the role and benefit of public intervention. But the immediate question is: How does the economy behave if people have noneconomic motives and rational responses (e.g., the other-regarding-preference model)? Economic motives but irrational responses (e.g., the behavioral model)? Noneconomic motives and irrational responses (psychological model)? It is obvious that the answers to important questions about how the economy behaves, and what public interventions are needed when it misbehaves, vary significantly from one model of the economy to another. Each model corresponds to a different vision of behavior in the economy. The role of the government is to set the conditions under which irrational people can be harnessed creatively (without harming the rational people) to serve the greater good. Some forms of irrational behavior may give the government new opportunities to step in with policy interventions.

Behavioral economics is interested in three forms of imperfections in decision-making due to imperfect rationality, imperfect will power, and imperfect self-interest. We now explore these imperfections and their consequences.

### 3.2.1   Simple Example: How Much to Save?

The United States Consumer Expenditure Survey suggests that households spend more out of dividend income than out of income from capital gains (Baker, Nagel, and Wurgler 2006). Rational economic theory implies that the same proportion of the extra money (when the extra is the same) should be spent in both cases. To address this issue, Shefrin and Thaler (1988) conducted an experiment that asked subjects how much they would spend out of an unexpected extra money gain of $2,400 in three possible situations (or framings).

• *First framing*   The extra money is a bonus at work paid out at constant rate of $200 a month over a year (so that the total amount is $2,400). The median saving out of the extra money was $100 monthly for a total saving of $1,200.

• *Second framing*   The extra money is a lump-sum payment of $2,400 this month. The median response was that $400 would be spent immediately, and $35 a month for the rest of the year. That makes total spending of $785, and so a total saving of $1,615.

• *Third framing*   The extra money is invested in an interest-bearing account for five years and the subjects receive at the end of five years $2,400 plus interest (so that the present value of the payment is $2,400). The median response is that none of the (future) capital would be spent during the first year.

Rational people should spend the same portion of the extra money for all the frames. Shefrin and Thaler interpret the deviation from the rational choice outcome as suggesting that people place different kinds of income in different *mental accounts*: the current income account (in frame 1), the asset income account (in frame 2), and the future income account (in frame 3), and that they spend differently from each of these mental accounts (spending much less out of the future income account). This is one of the many facets of the so-called framing effect to which we return later. A framing effect is usually said to occur when equivalent descriptions of a decision problem lead to systematically different decisions. Framing effects are commonly taken as evidence for incoherence in human decision-making, and for the empirical inapplicability of the rational choice models used by economists.

### 3.2.2   Present Bias

Present bias is an explanation of the self-control problem. Strotz (1957) analyzed how preferences that change over time would affect the saving decision, but present bias is a much more general problem where people face a choice that is liable to change their own preferences in the future. This is comparable to Odysseus when passing the Sirens and their enchantingly seductive voices. Odysseus faced an interesting decision on how to resist the temptation (either sailing close to the island to hear the Sirens but then running the risk of directing his ship onto the rocks, or sailing past the island and missing the chance to hear the Sirens).

A similar situation is the decision to initially engage in an addictive activity (gambling, smoking, drinking, etc.). There is first a decision to start the addictive activity (option $Smoke$) or not (option $No$). An initial decision to start the addictive activity, like smoking, leads to a future decision of whether to continue the addictive activity or to stop (option $Quit$). Because of the addictive nature of the activity, future preferences are such that option $Smoke$ is preferred to option $Quit$: you are hooked to this addictive activity in some sense. But initial preferences (before addiction kicks in) are such that people may want to give it a try, so that option $Quit$ seems better than option $No$, which in turn is better than option $Smoke$, which is the awful option. So the initial preferences at time $t$ are

$$Quit \succ_t No \succ_t Smoke,$$

but future preferences at time $t + 1$ when the addiction has set in are

$Smoke \succ_{t+1} Quit.$

Facing this preference reversal, a naive individual ignores the issue entirely, and plans to reach what is initially the best option $Quit$ by starting the addictive activity. But then in the next period, when the addiction has set in, the awful option $Smoke$ is the final result. On the contrary, a sophisticated individual foresees that $Smoke$ will be preferred to $Quit$ in the future so that $Quit$ is not really a feasible option. The choices boil down to either $Smoke$ or $No$, and the sophisticated agent chooses $No$, avoiding any risk of addiction. In this self-control problem there is, in effect, one individual self today and an entirely different individual self in the future. The today self and the future self have different preferences. The naive individual who ignores this duality seems obviously irrational. The sophisticated individual is rational, in a sense, but achieves rationality only by realizing the truth that there is a preference reversal between today self and future self.

The self-control problem is a dynamic consistency problem in which the naive individual fails to take into account the future preference reversal at the initial stage. There is an inconsistency between what he would like to do tomorrow, and what he would do in effect tomorrow. Odysseus solved the present-bias problem in a different way: by precommitment. He created an extra option by being tied to the mast of the ship that allowed him to hear the Sirens without being able to direct his ship onto the rocks. The introduction of methods of precommitment is always a useful public policy intervention.

### 3.2.3   The $(\beta, \delta)$ Model of Self-Control

Rational choice implies time consistency, namely that decisions are not sensitive to timing. Time consistency means that an initial consumption plan can be constructed, and that this plan will not need to be revised as times passes. In many domains early decisions are not carried out because consumption plans change over time. There are many illustrations such as: next month I will quit smoking, next week I will study and catch up on my homework, tomorrow morning I will wake up early and exercise; after the Christmas vacation I will start eating better; next weekend I will send in this form, next month I will start saving, and so forth. Early plans tend to have gratification up front and the "good" behavior to follow later (lie in today, but get up at 6 am tomorrow to finish the problem set), but when tomorrow comes instant gratification is again chosen and the "good" part of the plan is delayed. It is as if people are inhabited by multiple selves that disagree (internal conflict). Early selves make plans and choices that later

selves will not want to follow. Plans made at a distance tend to be more patient than choices made in the present. Dynamic inconsistency implies a conflict between early and late selves (i.e., a preference reversal). It is in essence a self-control problem (like procrastination, laziness, addiction, or compulsive consumption) where people cannot act according to plans.

There is a simple way to model the self-control problem. As it involves intertemporal choice, we need to use discounting factors for future utility levels. The self-control problem can be represented by the $(\beta, \delta)$ model (or *quasi-hyperbolic utility*)

$$U = u_0 + \beta\delta u_1 + \beta\delta^2 u_2 + \ldots + \beta\delta^T u_T, \tag{3.1}$$

where $\delta < 1$ is the standard discount rate and $\beta \leq 1$ is self-control discounting. For $\beta = 1$, there is no self-control problem. For $\beta < 1$, the immediate future is more heavily discounted than the more-distant future (this is the "present bias"): in the long-run we are relatively more patient than we are in the short run. The discount rate between any two periods in the future is $\delta$ whereas the discount rate between the present (time 0) and the immediate future (time 1) is $\beta\delta$.

To see how the model works, consider a consumption decision involving utility $u_1$ at $t = 1$ and delayed utility $u_2$ at $t = 2$. If it is an *investment* good (like exercising, studying, training, or savings), it has the feature that people must trade off the cost $u_1 < 0$ against a future benefit $u_2 > 0$. If it is a *temptation* good (like compulsive consumption, unhealthy food, surfing on the web, or credit card usage), it has the feature that people must trade off the reward $u_1 > 0$ against a future cost $u_2 < 0$. What is the consumption decision from an ex ante perspective?

If the consumer could commit to a choice in advance, say at time $t = 0$, she would consume if $\beta\delta u_1 + \beta\delta^2 u_2 \geq 0$ or, equivalently,

$$u_1 + \delta u_2 \geq 0. \tag{3.2}$$

Note that the parameter $\beta$ cancels out in the desired (future) consumption choice. However, the consumer actually consumes at time $t = 1$ if

$$u_1 + \beta\delta u_2 \geq 0. \tag{3.3}$$

Compared to the desired level of consumption set in advance, the naive individual (with parameter $\beta < 1$) underconsumes investment goods (with delayed benefit $u_2 > 0$) and overconsumes temptation goods (with delayed cost $u_2 < 0$), since $\beta\delta < \delta$. This is the self-control problem. Compared to the actual consumption, the naive individual overestimates the consumption of investment good ($u_2 > 0$) and underestimates the consumption of temptation good ($u_2 < 0$). Conversely, a sophisticated agent (without

self-control problem, so $\beta = 1$) will actually consume according to the plan since $\beta\delta = \delta$.

### 3.2.4   Reference-Dependence Bias

One explanation for some of the observed anomalies is that people assess alternative options by comparing them to a reference point. The reference point might be a previous level of consumption or a target level of consumption. Whatever the explanation, the key assumption is that utility is measured relative to the reference point.

Denote the reference point by $r$. Utility depends on a combination of an absolute (possibly stochastic) consumption utility, $m(x)$, and a penalty function $\mu(m(x) - m(r))$, which is reference dependent since it is determined by the deviation $m(x) - m(r)$ between consumption utility and reference utility. When the reference point is the same as the bundle chosen, then $x = r$, and the penalty term disappears, so the model reduces to standard consumer choice. This approach creates multiple equilibria, opening up a role for marketing, advertising, and sales prices to influence preferences by creating and changing the reference point. This approach also helps explain conformism and the effect of experience on adaptive preference. Reference effects can be seen by comparing the situation of owners and dealers in either the housing market or the car market. Dealers do not expect to hold on to goods they receive. Since their reference point does not include the goods, they do not experience a loss when selling them. In contrast, the owner of a car and the owner of a house will exhibit some endowment effect according to which their reference points include those goods, and so they feel more of a loss when selling them. This is a general difference you feel when you buy a good for resale rather than for utilization.

Another illustration of reference-dependent preferences is the fact that people value income changes as well as income levels. Standard preferences involve valuing only income and consumption levels. Reference-dependent preferences assume that the value function $v(x;\ r)$ is defined over differences from a reference point $r$ instead of over the overall income level. A simplified version involves the following reference-dependent preferences over income, $x$ :

$$v(x;r) = \begin{cases} x - r & \text{for } x \geq r, \\ \lambda(x - r) & \text{otherwise.} \end{cases} \tag{3.4}$$

The parameter $\lambda > 1$ denotes a loss aversion parameter that overweights losses: the value function is steeper for losses below the reference point ($x < r$) than for gains

over the reference point ($x > r$). This simple formulation can explain the asymmetry in the valuation of small equal-sized losses and gains exhibited in many experiments.

Loss aversion can also explain the endowment effect. Consider the housing market. Homeowners willing to sell their houses are likely to fix the sale prices with reference to the initial purchase prices. A homeowner values the sale according to the extent of deviation from the purchase price. If there is loss aversion, the homeowner will overweight a price loss compared to a price gain. The homeowner who fears selling at a loss is willing to ask a higher sale price. Obviously a higher sale price will increase the utility of a sale, but it will also reduce the probability of a sale. The homeowner will trade off these two opposing effects. The loss-averse homeowner will sell above the purchase price. We observe similar application in the stock market where the tendency is to sell "winners" and hold back "losers."

The reference-dependent model can also explain some anomalies in response to price changes. Consider the labor supply response to an hourly wage increase. For a rational worker a higher hourly wage induces longer working hours (if the substitution effect dominates the income effect). However, for a reference-dependent worker, the wage increase may well reduce her labor supply if the wage increase shifts income above the target income used as the reference point used to value labor choice. The worker achieves her target income by working less. Similar reasoning applies to the response of saving rates to interest rates: the household achieves the target, saving income, by saving less when the interest rate is higher.

### 3.2.5  The Gambler's Fallacy

The gambler's fallacy is the first of three behavioral anomalies we consider that involve the mistakes people make when forming beliefs, and these mistakes distort their decisions.

The least controversial type of mistake about objective, real world facts is non-Bayesian statistical reasoning. Such errors are likely to affect investment decisions and many other economic decisions under uncertainty. Consider the simple example of coin flips. Let $\{h, t\}$ represents a lottery that pays $\$h$ if the next flip of a coin comes up heads ($H$) and $\$t$ if the next flip comes up tails ($T$). Consider the following situation: If the person observes $HHH$, she chooses a lottery with $h < t$; that is, she chooses the lottery that pays more if the next flip is $T$. If $h = t - \Delta$, this is equivalent to a bet of the amount $\Delta$ that $T$ will come up next. If instead she observes that the previous flips are $TTT$, she chooses a lottery with $h > t$; that is, the lottery that pays more if the next flip is $H$. If $h = t + \Delta$, this is equivalent to a bet of the amount $\Delta$ that $H$ will come up

next. And finally, if she has observed no flips before, she chooses a lottery such that $h = t$, considering that heads and tails are equally likely. These choices are wrong in the sense that no matter what previous outcomes have appeared the probability of $h$ arriving next (or $t$) is always one-half.

The specific pattern of mistakes in these choices are called the *gambler's fallacy*: the person believes that if the same realization of the random process has occurred a number of times, the other realization is in some sense "due" for the next draw.

### 3.2.6   Confirmation Bias

Confirmation bias arises from inferring *less* than what is justified from the observation of a recent event. This is the tendency to perceive data as more consistent with a prior hypothesis than they truly are. The agent updates her information based on unfolding observations by overweighting information that confirms her initial opinion and under-weighting information that contradicts her initial opinion. This is a Bayesian updating framework, except for the mistake (bias), in the encoding of data. The gambler's fallacy is a special case with the misperception that successive samples are drawn without replacement.

Another deviation from Bayesian reasoning is *irreversibility*. According to Bayesian rationality, if you discover that a piece of information is mistaken, the memory should erase it so that it will have no impact on future judgment (i.e., information reversibility). However, the brain is such that when new information merges with old information, it is impossible to undo the effect of the old information even when such information is mistaken. Information in the brain is long-lasting. For example, when juries are instructed to ignore certain statements after they have been heard, it is hard for them to fully ignore the statements when making their final judgment. Information sticks where it hits.

*Hindsight bias* is the opposite of confirmation bias: it is inferring *more* than is justified from the observation of a recent event. It reflects our tendency to rapidly rewrite our memory of the past to fit what we have just learned. Rapid rewriting creates "hindsight bias"; that is, the ex post recollection of the ex ante probability of an event will be biased in the direction of the event's realization. The problem is that revising our beliefs rapidly reflects more about how little we knew before the event, and so we should not reach too quickly to make strong beliefs out of recent events. Hindsight bias is on display every day in sport events and news coverage. It is an important force in political life and in organizational life.

### 3.2.7  Confidence Bias

Many studies show that people are often *overconfident*. Think about cycling without a helmet or driving without a seat belt because we believe our skills we help us avoid an accident. Overconfidence is the tendency to overestimate one's own (relative) abilities and expect the resulting outcomes to be better than they will be. Similarly *overoptimism* is the overestimation of general prospects.

There are many ways to model overconfidence with different practical implications. One possibility is that people overestimate the output they can generate, or they overestimate the marginal productivity of their effort. In either case they may end up striving less hard than if they were not overconfident. A different possibility is that overconfident persons think they are more skilled and talented that they really are. Drivers are overconfident about overall driving ability. This is especially true for young drivers and men, but much less so for women and mature drivers.

Whether overconfident persons will exert too little effort or too much will depend on whether effort and skill are complementary. Obviously overconfidence is not uniform in the population. For example, some studies show that women are less (over)confident than men, which might explain why women feel the need to work harder at school to achieve success.

### 3.2.8  Framing Bias

Any theory of rational choice must stipulate that the same problem will be evaluated in the same way regardless of how the problem is described; thus different but equivalent descriptions should lead to the same choice. Framing effects violate this bedrock normative condition of "description invariance." As already discussed in the saving example, a framing effect occurs when different but equivalent descriptions of a decision problem lead to systematically different decisions, and this is commonly taken as evidence for incoherence in individual choices.

The best-known framing problem is *risk framing*. It was first described by the so-called Asian disease problem (Tversky and Kahneman 1981). Participants to the experiment are first told the following story: "The United State is preparing for the outbreak of an unusual Asian disease, which is expected to kill 600 people. One possible program to combat the disease has been proposed." Then some participants to the experiment are presented with two options. *A*: If this program is adopted, 200 people will be saved. *B*: If this program is adopted, there is a one-third probability that 600 people will be saved and a two-thirds probability that no people will be saved. Other

participants are presented with two other options. *C*: If this program is adopted, 400 people will die. *D*: If this program is adopted, there is a one-third probability that nobody will die and a two-thirds probability that 600 people will die. The robust experimental finding is that subjects tend to prefer the "sure thing" *A* when given options *A* and *B*, but tend to prefer the gamble *D* when offered options *C* and *D*. Note, however, that options *A* and *C* are equivalent, as are options *B* and *D*. Subjects thus appear to be risk-averse for gains preferring option *A* to *B*, and risk-seeking for losses, preferring option *D* to *C*. This is a central feature of *prospect theory* (Kahneman and Tversky 1979). In prospect theory, it is the decision maker's private framing of the problem in terms of gains or losses that determines her evaluation of the options; the framing manipulation is thus viewed as a public tool for influencing this private frame.

There are many other framing problems. In *attribute framing* a single attribute of a single object is described in terms of either a positively valued attribute or an equivalent negatively valued attribute. The subject is then required to provide some evaluation of the described object. The typical and robust finding is that objects described in terms of a positive attribute are generally evaluated more favorably than equivalent objects described in terms of negative attribute. For example, in one study, beef described as "75 percent lean" was given higher ratings than beef described as "25 percent fat." In *goal framing* subjects are urged to engage in some activity (e.g., wearing seat belts). This plea involves a description of either the advantages of participating in the activity or the corresponding disadvantages of not participating. The most common result is that subjects are more likely to engage in the activity when the disadvantages of not engaging, rather than the advantages of engaging, are emphasized.

The framing effects reflect more generally the fact that human perception and cognition is heavily influenced by contrast. This is clear in the Titchener illusion of circles, where a circle looks larger when surrounded by smaller circles than when it is surrounded by larger circles. Since choices involve basic perceptions, it would then be surprising if choices were not sensitive to contrast as well. Similarly the comparison of an outcome with unrealized outcomes (disappointment) or with outcomes from forgone choices (regret) imply that the appeal of choices depends on the set of choices they are part of. There are situations where the possibility for more choices implies that fewer choices are made, or even no choice altogether. It is well documented that offering a broader set of choices can lead consumers to buy less, whereas in the standard model more choices can only lead to more purchases. The explanation is that the broader opportunity set makes choice more difficult and stressful, so some individuals may prefer to avoid choosing. Similarly individuals

facing difficult choices may prefer to go for a default option such as the "menu of the day" (or the house wine) in a restaurant offering a very large choice of menus (or wines).

### 3.2.9   Conformism Bias

Conformism in social psychology refers to the inclination of an individual to change spontaneously, without any explicit order or request by someone else, her opinion (beliefs) and or action (choices) to conform to the socially prevailing opinions and actions. It is the individual tendency to change the intrinsic optimal choice toward the most prevailing choices within a group.

We could attribute this conformism to some general force of learning, such as imitation or suggestion, regarded as innate and instinctive. There is nothing wrong with that. However, it could also reflect a possible lack of autonomy leading to individual mistakes when making choices. The conformism effect was forcefully demonstrated in Asch's famous experiment. The experiment consisted of many trials in which a subject was placed in a group of people who were secretly accomplices of the experimenter. The subject was asked to estimate the length of a line by matching it with one of three lines. This estimate was provided after the other accomplices had successively expressed their opinions. The reality of conformism emerged when the other accomplices announced a clearly wrong comparison line, as then about one-third of the tested subjects revised their own (correct) opinions to conform to the wrong judgment of the group. This finding is important because it violates one of the basic postulate of standard economics, that is, the idea of full autonomy of the individual. Indeed, despite the lack of any uncertainty about the correct judgment, agents may renounce their preferences (or judgments) and conform to an erroneous choice of the others. This is conforming to social pressures.

Unlike sociopsychologists, economists are reluctant to take conformism as a primitive assumption but prefer to derive conformism endogenously. A first possible explanation for this conformism bias is that individuals may suffer from being "different" and that conformism is a reaction to this. Another explanation assumes agents care about "status," which can only be inferred from their actions. If all agents prefer to be perceived as "good" by others and if they all agree about what is a "good" type, then they will all make uniform choices through fear that they would appear "bad" by deviating from the norm.

A simple way of modeling conformism is to extend the utility function of an individual $i$ to include a penalty $p(x_i, m)$ that depends on the distance between her individual

choice $x_i$ and the "normal" choice $m$. The penalty function must be of the first order, for otherwise a small deviation away from the norm toward the intrinsic preferred choice would be beneficial as it would only cause a second-order loss in the penalty function and first-order gain in the intrinsic utility. Obviously the dispersion of choices is important because the smaller the dispersion, the greater the degree of conformism. A suitable penalty function that takes account of the dispersion of choices and the first-order condition is the index of individual conformity measuring the standardized distance from the mode:

$$p(x_i, m) = \frac{|x_i - m(x_i, x_{-i})|}{\sigma(x_i, x_{-i})}, \tag{3.5}$$

where $x_i$ is the choice of individual $i$, $m(x_i, x_{-i})$ is the mode and $\sigma(x_i, x_{-i})$ is the standard deviation around the mode (where $x_{-i}$ denotes choice variables pertaining to anyone other than individual $i$). This index is a simple and perhaps realistic representation of how individuals synthetically compare themselves to others.

### 3.2.10   Identity and Social Norms

There is a substantial psychological literature on group identification that is important for social interaction and market allocation. A special case is the "identity economics" coined by Akerlof and Kranton (2010). To get a grasp of this approach, consider the relative performance of boys and girls at school. It is a general feature and well-documented fact across various PISA (Program for International Student Assessment) studies that girls perform better on average than boys at age 15. Why is it so? This cannot be due to better salary prospects, since women are paid less and are more likely to work part time. Pursuing the analysis is in fact intriguing because we cannot attribute such difference to either family background or migration status nor to gender difference in those factors. We cannot attribute such gender difference to school difference either (because they attend the same schools in general). Finally, we can hardly claim gender difference in cognitive ability. So there is something else. Something less visible and obvious but still very important. What could it be? This is where identity and norm come in.

When we examine people's decisions from the perspective of their identities and social norms, we provide new insights into many different economic questions. Who people are and how they think of themselves is key to the decisions that they make. Their identities and norms are basic motivations. Tastes vary with social norms. This vision of tastes is important because norms are powerful sources of motivation. Norms

affect fine-grain decisions of the moment. Norms drive life-changing decisions as well: on matters as important as whether to quit school, whether to go to university, or whether to go to work.

The important determinant of whether an organization functions well is not only the monetary incentive system, as standard economic models would suggest, but also how well its members identify with the organization and with their activities within it. Their work must have some meaning for them to function effectively in a company. Workers may well be willing to trade a reduction in wages for "meaningful" work. This is probably a central feature in the occupational choice of many workers in public services with a real vocation of serving the public (e.g., in education, health, and justice services). This effect refers to *intrinsic motivation*: the satisfaction a worker gets from work for its own sake. An interesting phenomenon documented in psychology is the possibility that extrinsic incentives (e.g., money) can "crowd out" intrinsic motivation. Blood donation is a concrete illustration. In general, if workers do not identify with their job, they will seek to game the incentive system, rather than to meet the organization's goals. Likewise good schooling occurs not as a result of monetary rewards and costs but because students, parents, and teachers identify with their schools, and because that identification is associated with learning. Given this, education policy should look at what some successful programs have done to establish a school identity that motivates students and teachers to work according to a common purpose. If we focus on training teachers in how to inspire their students to identify with their school rather than teaching students to take standardized tests, we just might be able to reproduce within these schools great results.

As economists and policy makers we could be content to continue looking only at prices, incomes, and related statistics to explain people's decisions. In some situations, that might be enough to understand what is happening. But in other situations, we would miss major sources of motivation and thus could adopt useless, if not counterproductive, measures aimed at producing the outcomes we seek.

## 3.3   Behavioral Markets

In 1991 Vernon Smith attacked Daniel Kahneman (before both received Nobel Prize in 2002 for work in behavioral economics). Smith's claim was that anomalies at the individual level play no role at the aggregate level, in particular in competitive markets. Introducing the possibility of systematic imperfections in individual rationality, as we all are willing to accept, raises new questions in the study of markets. Among them

the central question is whether the market will erase or exploit limits on consumer rationality? Does competitive pressure eliminate irrational choices and induce agents to make rational choices?

Individuals make mistakes that markets do not fully correct. There are three factors limiting the extent to which the market can remedy biased decisions. First, decisions are not frequent and do not deliver clear feedback. Second, individuals are not specialized in making those decisions. Third, individuals are protected from market pressure and competition.

### 3.3.1   Money Pump

Davidson, McKinsey, and Suppes (1955) use the money pump argument to justify rational choices. The argument is as follows. Suppose that a consumer has nontransitive preferences for consumption bundles $a$, $b$, $c$, where $a \succ b$, $b \succ c$, and $c \succ a$. The consumer has cyclic preferences, preferring $a$ to $b$, $b$ to $c$, and $c$ to $a$, meaning that the consumer is willing to trade $a$ for $c$, next $c$ for $b$, and then $b$ for $a$, so getting back to the initial bundle. That is, the cyclic consumer is always willing to pay a small amount $\Delta$ of money to get $a$ instead of $b$, $b$ instead of $c$, and $c$ instead of $a$. So, by allowing the agent to cycle between the different bundles $c$, $b$, $a$, $c$, $b$, $a$, ... against successive small payments of $\Delta$ units of money, the market will "pump" an indefinite amount of money out of the consumer. This argument establishes that an irrational consumer with intransitive preferences is doomed to bankruptcy when operating in the market. It suggests that all intransitivities should be removed from a consumer's preference.

### 3.3.2   Complementary Mistakes

Whether individual mistakes are erased or exacerbated depends on whether behaviors are strategic substitutes or strategic complements. When behaviors are complements (like in the stock market with buying or selling decisions), a small number of irrational traders can force others to behave irrationally.

A good illustration is the guessing game. This guessing game is like a market bubble where market participants are not rational and the market not in equilibrium. Consider a classroom experiment where students are asked to pick a number between 0 and 100 and not to let others see their pick. The winner of the contest is the student who is closest to two-thirds of the average number picked by all students (hence the name of Guessing Game). Ties will be broken randomly. In this game there is no dominant

strategy (i.e., a unique best guess independently of the guess of others). However, there is a unique Nash equilibrium when every player is rational *and* expects each other to behave rationally. To compute it, ask what guess would be irrational and eliminate this guess. For instance, any guess above 66.67 is irrational for every player, since it cannot possibly be two-thirds of the average guess. These guesses can reasonably be eliminated for every player, but then since no player is expected to guess above 66.67 and two-thirds of 66.67 is approximately 44.45, any guess above 44.45 is also irrational. This process of iterated elimination of weakly dominated guesses will continue until all guesses above 0 have been eliminated.

Now when the experiment is performed among ordinary students, it is usually found that the winning guess is much higher than 0. Some students guessed close to 100 indicating they did not understand the game at all. A large number of students guessed 33.3 (i.e., two-thirds of 50), indicating an expectation that other players will guess randomly. A small but significant numbers of students guessed 22.2 (i.e., two-thirds of 33.3), indicating a second iteration based on an assumption that others would guess 33.3. In many experiments the average guess was around 33, and the winning guess was around 22. So we can see that the rational equilibrium does not predict well in this strategic environment. Interestingly, even perfectly rational players participating in such a game should not guess 0 unless they know that others players are rational. If a rational player believes that others are not rational, she will not follow the chain of elimination described above and she could rationally guess above 0.

Keynes (1936) believed that similar behavior was at work in the stock market and could explain a market bubble. This is the case when the price of shares is not based on what people think their fundamental value is; rather, the price is based on what they think everyone else thinks their value is, or what everybody else would predict the average value is.

The nonequilibrium approach assumes a "cognitive hierarchy" in which more "rational" players best-respond to the perception that others do less thinking. Nagel (1995), based on an original idea of Hervé Moulin, suggested that people based their guesses on different levels of rationality. She found many guesses on different levels of rationality: level 0 rationality (guessing 50); level 1 rationality with best response to level 0 rationality (guessing 33 in response to 50); level 2 rationality, with best response to level 1 rationality (guessing 22 in response to 33); and so on. These cognitive hierarchy approaches are more precise than Nash equilibrium because they always predict a single statistical distribution of play, and are generally more accurate than equilibrium in predicting behavior. Interestingly, if the experiment is re-run with the same players, the results will now come much closer to the rational prediction.

### 3.3.3   Rationality Tug-of-War

Will the market and firms take advantage of limited consumer rationality or will the market and firms help consumers? Whether markets will correct irrationality depends on factors like whether consumers know their own limits and hence are receptive to advice, and whether there is more profit in protecting consumers or taking advantage of them.

As an illustration consider the Gabaix–Laibson (2006) model. There is a market of products with "add-on" prices (e.g., bank transaction fees that can be easily hidden). If consumers do not know about the hidden add-on price, then competitive firms will offer a low price on base goods (below marginal cost) and will charge high markups on add-ons. Sophisticated consumers who know the add-on price, but can easily substitute away from the add-on (e.g., avoiding bank ATM fees), will prefer products with expensive add-ons to benefit from the low base-good price. The naive consumers are subsidizing the sophisticated consumers. As a consequence competition does not lead to add-on prices being revealed (to protect the consumers) because a firm that reveals its add-ons will not attract either naive consumer (mistakenly thinking price is too high) or sophisticates benefiting from low base-good price.

## 3.4   Behavioral Policy

The standard economic approach assumes that people make appropriate decisions when they are well informed. The role of the government is then to combat ignorance or misinformation. So, if the government can provide relevant and reliable information more effectively than private markets, it will be a potentially beneficial intervention for consumers to do so. Information and education campaigns to combat ignorance are desirable to help people make informed choices.

Assuming information is not an issue, the reasons for government intervention still involve market failure. It may be appropriate, for example, for the government to tax pollution and subsidize charitable contributions to ensure adequate levels of their provision (see part III on departures from efficiency). However, once market failures are corrected, under the standard approach there will be nothing wrong with the choices people make, given the constraints they face (obviously that does not rule out government intervention to enforce property rights and redistribute resources!).

In practice, benevolent policy makers worry that people make inappropriate choices. One important example is that governments are concerned that people save "too little"

and that this is a systematic mistake given their present bias. Such decision-making failures are another justification for intervention.

### 3.4.1   Internalities versus Externalities

Behavioral public policy permits the possibility of decision-making failures, and this opens up the possibility of enhancing individual welfare by correcting or preventing "bad" choices. Behavioral policy is an extension of standard public policy. The common feature is that public policy can change behavior by changing relative prices, budgets, and information. The distinctive feature of behavioral policy is that there are additional channels through which policy can change behavior and welfare, even if the policy leaves prices, budgets, and information unchanged!

The central idea is that if people make systematic mistakes, or biased decisions they regret later on, behavioral policy uses those mistakes and biases that hurt people (1) to help them protect themselves and, at the same time, (2) to respect their autonomy of choice.

What do all these "invalid" choices have in common? People are facing the wrong prices.

With traditional public policy, prices are wrong or misaligned because of *externalities.* Externalities are costs that people impose on others but do not internalize. So, in the presence of externalities, the prices people pay for things do not reflect the "true cost" to *others*: the market price is wrong. The point will be explored in depth in chapter 8.

With behavioral public policy, prices are wrong because of *internalities.* Internalities are costs that people impose on themselves but do not internalize. So, in the presence of internalities, the prices people pay for things do not reflect true costs to *themselves.* This is so because people face internal conflicts in their choices: the question "Do you want a piece of chocolate?" triggers an internal clash between temptation and reason for someone on a diet program. People might have present bias, reference-dependent bias, and so on. It is not enough simply begging people to do the right thing. The government must realign prices and incentives so that it is in people's own interest to do the right thing.

Behavioral policy uses people's biases to help them by making the healthy option the default (status quo bias) and by giving immediate rewards for healthy choices (present bias). Behavioral policy supercharges economic incentives with deposit contracts and regret lotteries to act as commitment devices against people's time inconsistency. Behavioral public policy is consistent with the traditional justification for public intervention (the enforcement of property rights, correction of market failures, redistribution

of income). It also introduces justifications for intervention, notably by allowing public policy that raises welfare by limiting the possibilities for decision-making failure and its consequences.

### 3.4.2   Automatic Enrollment

Consider the public policy of automatic enrollment to a saving plan or organ donation scheme with a small cost to opt out. The default option is the outcome resulting from inaction. A rational consumer is not influenced by the default option because the cost to opt out is small. But behavioral consumers are influenced by the default option because of status quo bias. In practice, the choice of the default option matters: for organ donation schemes and retirement saving plans there is considerable evidence that the default option affects participation rates, even though such a default neither affects opportunities (low cost to opt out) nor provides new information. Johnson and Goldstein (*Science* 2003) report a significant effect of automatic enrollment into organ donation: the consent rate is 85.9 percent in Sweden with automatic enrollment against 4.3 percent in Denmark where there is no automatic enrollment. The introduction of automatic enrollment in 401(K) pension plans in the United State also provides an effective default option (costlessly changeable by employees). The default contribution rate is 3 percent of the salary into the 401(K) plan and 100 percent of allocation to the money market. Before the automatic enrollment, the participation rate was 40 percent with a dispersion of contribution rates. After the automatic enrollment, the participation rate was 88 percent without dispersion of contribution rates, which are mostly concentrated around the default rate of 3 percent. This automatic enrollment policy looks nonpaternalistic because the desirable default can improve the welfare of those who mindlessly adhere to the default without restricting the options of those who do not. The policy faces little opposition because it is beneficial to those who believe there is a status quo bias, and it is harmless to those who believe there is no status quo bias.

### 3.4.3   The SMarT Plan

People accumulate an inadequate level of savings because of the different biases involved. The present bias implies a willingness to save, but only tomorrow. The opportunity-cost bias implies a willingness to save but only out of pay increases. The opportunity-cost bias reflects the tendency for people to treat "out-of-pocket" costs differently from "opportunity costs." People tend to underweight (or neglect)

opportunity costs relative to out-of-pocket costs. The status quo bias implies a willingness to adhere to a saving plan once implemented. Combining these different biases, Benartzi and Thaler (2004) designed a saving plan to help workers save more. The plan was implemented in a workplace with employees invited to join a saving plan. They could elect in advance a portion of their current income, and also a possibly different portion of future income, to be saved. This program induced large increases in saving. Although workers chose to save little out of current income, they *committed* to save a large portion of future income. In a short period of time, the average saving rate increased from 4.4 to 8.8 percent.

### 3.4.4   Complementarity

Behavioral public policy should complement, not replace, more substantive public policies. For example, if standard public policy suggests creating a price differential between healthy and nonhealthy food, or between bottled water and a soda drink, behavioral analysis could help us better understand the consumer response to various forms of public intervention. Behavioral analysis could suggest whether consumers would respond better to a subsidy on bottled water or to a tax on regular soda. It would also suggests how to complement a tax-subsidy policy with more effective labeling of healthy/nonhealthy food. But that's the most it can do according to the behavioral economists themselves.

The limits to behavioral policy arise from social interaction. Indeed a key difference between psychologists and economists is that psychologists are interested in individual behavior while economists are interested in explaining the outcome of social interaction among many individuals. Behavioral economics focuses on human dysfunction (interpreted as making irrational choices) with a perspective to help people become more functional. But, in general, most people are sufficiently functional most of the time. Hence the focus of economists on people who are "rational" remains a useful benchmark to form policy choices. Behavioral analysis can contribute to the improvement of existing policies, but it offers no realistic prospect of replacing policies. The narrow and complex models of behavior used in psychology cannot easily be used to study the behavior of many people interacting. However, it is the social interaction of people that will eventually determine the final outcome of any policy intervention. The need to study groups that consist of large number of people requires constraints on economics that are not present for the psychologist. Economists need simple and broad models of behavior. Hence economists focus on rational and selfish behavior that provides a reasonable description over a broad range of social settings.

## 3.5   Behavioral Welfare

Standard welfare analysis is concerned about the evaluation of public policy. To evaluate whether a policy is good or bad, we must evaluate its effect on individual welfare levels. The effect of policies on welfare are calculated in two stages: (1) the effect of the policy on behavior and (2) the effect of the change in behavior on welfare.

The standard welfare approach to compute the effect of policy on behavior is to assume individuals will rationally respond to policy change. Then the preferences revealed by choices are used to compute the effect of change in behavior on welfare. The welfare effect of policy intervention is measured by the extent to which individual preferences are satisfied. The standard preference revelation axioms assume that the choices people make are valid and thus correctly reveal their true preferences. There is no conflict between actions and intentions. The assumption that preferences are revealed by choices can be traced back to Vilfredo Pareto (or Wilfried Pareto at birth). Pareto recommended the use of choices ("objective facts") to reveal preferences ("subjective fact"). This is a philosophical stance and not a robust empirical regularity. Pareto justified his assumption by restricting attention only to repeated actions so that rational choices emerge as the consequence of learning. For instance, credit card account holders learn to pay their bills on time by first suffering the payment of a late fee when they do not. But because Pareto clearly limited the domain of revealed preference to "repeated actions" in which learning has taught people what they want, he leaves out important economic decisions that are rare, partly irreversible, or difficult to learn about from trial and error: educational choices, occupational choices, retirement and saving plans, fertility and mate choices, housing choices, and so forth. In principle, we could consider that people learn by observing others, but people are generally far more responsive to their own experiences than to the experience of others.

### 3.5.1   New Welfare Criterion

The real challenge with the behavioral approach is that people can make mistakes. Choice mistakes represent a conflict between actions and intentions. So the question is then how are (true) preferences revealed when choices reveal mistakes rather than preferences, that is when actions do not reveal true intentions? There are many cases, as described in earlier sections of this chapter, in which even the choices of mature consumers do not reveal a true preference but rather reflect the combined influence of true preferences and choice mistakes.

When people cannot act according to plan or when people do not save enough for their retirement, we cannot assume that they are acting in their own best interest, and that those decisions (or nondecisions) are revealing their true preferences. If preferences are only imperfectly revealed by choices, what yardstick can we use to evaluate policy recommendations? The danger is to replace individual preferences by some ad hoc external preferences to legitimate policy choices. This is the danger of paternalism such as when parents teach their children how to behave based on their own preferences. It may be legitimate to do this for children, but it is not legitimate for adults in full control of their autonomy (unless they voluntarily accept to abandon their freedom of choice). Often there exists a compromise, so economists may use choices to identify (true) preferences but take care to acknowledge the possible wedge between revealed preferences and true preferences.

### 3.5.2   Choice-Based Welfare Analysis

In standard theory, agents decide on choice, $x$, from a set of possible choices, $X$. The goal of policy is to identify the optimal choice $x^*$, $x^* \in X$.

In behavioral models, agents choose from *generalized choice sets* $G = (X, d)$ where $d$ is an *ancillary condition* that affects choice but by assumption does not affect (true) preferences (e.g., salience, framing, default option). Let $C(X, d)$ denote choice made in a given generalized choice sets $G$. Choice mistakes and inconsistent choices imply that different ancillary conditions $d \neq d'$ lead to different choices even if the choice set $X$ is unchanged

$$C(X, \ d) \neq C(X, \ d') \qquad \text{for } d \neq d'. \tag{3.6}$$

We can thus define the revealed preference relation $P$ as $x P x'$ if $x$ is always chosen over $x'$ for any ancillary conditions $d$. Using the revealed preference relation $P$, it is possible to identify the choice set that maximizes welfare instead of a single point.

With sufficiently many observed choices, it is effectively possible to obtain bounds on welfare. To illustrate this, consider three different saving plans with varying benefits and corresponding contributions rates: high ($H$), middle ($M$), and low ($L$). Suppose that we have collected observations from two different framing conditions, $d$ and $d'$. In frame $d$, revealed preferences are $H > M > L$ and the consumer chooses saving plan $H$, whereas in frame $d'$ preferences are $M > H > L$ and the consumer chooses saving plan $M$. We do not need to understand why the frame affects the choice of saving plan to make a welfare statement about the optimal policy. Indeed $L$ cannot be optimal given the observed choices because it is never chosen no matter what the framing condition

is. Therefore the optimal policy must be bounded between $M$ and $H$. That delivers bounds on welfare based entirely on choice observations.

The revealed preference relation can identify the set of choices that maximizes welfare (but not the unique optimal choice). Welfare bounds are tight when choices are less sensitive to framing conditions, that is, when behavioral problems are small. However, welfare bounds and the set of optimal policies are large when behavioral problems are large. In the previous example, if there exists another frame $d''$ such that $L > M > H$, then $L$ would be chosen and the welfare bounds would include $L$, $M$, and $H$. That is, any saving plans could be the optimal policy. So this approach is not restrictive enough to generate policy prescriptions when ancillary conditions exist that lead to vast changes in choices. There are two alternatives solutions: preferences refinements and the structural model.

### 3.5.3   Refinement and Structural Modeling

The idea behind preference refinement is to discard certain ancillary conditions that have become too "contaminated" for welfare analysis. For instance, by dropping the nonvalid framing $d''$, we can eliminate plan $L$ from the set of optimal policies. With fewer ancillary conditions we have more restrictive bounds on welfare and policy. There is a good argument that supports this alternative for redistributive policies and the importance of the status quo. According to reference dependence, people tend to concern themselves more with income change (gains and losses) than income levels. Moreover feeling the impact of loss is larger than feeling the impact of gain (loss aversion). As a result people will give higher subjective weight to avoiding a loss than experiencing a gain. Status quo comes into play if the redistributive policies cause the rich to lose and the poor to gain. Due to loss aversion more prominent weight is given to the rich relative to the poor when a redistributive policy is being evaluated.

Structural modeling takes a different approach. The idea of mapping observed choices directly into statements about welfare is abandoned. Instead, the observed choices are interpreted using a behavioral model that seeks to explain their deviations from rationality. The objective is to discover preferences by building a behavioral model that can explain how ancillary conditions affect the choices, and then to use this model to predict which choices reveal true preferences. By this approach, it is assumed that preferences can be identified from the observed choices on the basis of a structural decision-making model. For example, one could construct a present-bias $(\beta, \delta)$ model that explains saving plan choices, and then calculate the optimal policy within such a model using normative discounting with $\beta = 1$.

### 3.5.4   Application: Global Warming

Global warming involves intertemporal trade-offs that raise important normative questions for public policy. We discuss the economics of climate policy more fully in chapter 26. Here we focus on the fact that any discussion of policy must invoke trade-offs between different time periods (bear a cost now to mitigate emissions in order to benefit from lower temperatures later) and between different generations (should the current generation emit pollution that changes the climate for future generations?). Indeed, whether or not people are rational, there are good reasons why revealed choices cannot be used to infer the true intertemporal preferences of an agent, and why the choices cannot be relied upon to make normative judgments. We consider in turn the rational model and the irrational model.

For people who are *rational*, there is no reason for the policy maker to give any normative value to the discount rate revealed by individual choices. Basically, why should the payoff at time $t$ have less weight than payoffs at time $t + \Delta$? A good explanation may relate to the risk of mortality. But then how can we account for the fact that young adults are markedly less patient than middle-aged and older adults? Yet, discounting to account for the mortality risk is too small to account for the revealed exponential discounting rate of around 5 percent per year. This is a matter of considerable importance for curbing global warming. Indeed the policy recommendations for addressing global warming are heavily dependent on the choice of the normative discount rate. The Stern Report (2006) used a normative discount rate of 0.1 percent per year to make its recommendations, namely a normative discount rate about 1/50th of the revealed discount rate of 5 percent. It is therefore not surprising to find such a big gap between the call for immediate and massive actions to curb $CO_2$ emissions in the Stern Report and what people seem willing to accept. The rational approach to intertemporal choice assumes a constant discount rate whereby agents make choices that are consistent over time. The exponential discount function with constant discount rate, $U = u(x_0) + \delta u(x_1) + \delta^2 u(x_2) + \dots$, as originally stated by Ramsey (1928), is the *only* discount function that generates dynamically consistent choices: preferences held at some point in time do not change with the passage of time (unless obviously new information arrives). No preference reversal is often invoked as a rationality requirement.

If people are *irrational*, their choices may contradict their intentions. As was the case with the self-control problem, nonconstant discount rates imply dynamically inconsistent choices. For the global warming problem, suppose that an agent can make some investment at a cost of $C$ (i.e., pollution abatement costs) to gain delayed benefits

of $B$ (i.e., curb global warming). For simplicity, suppose that the benefit occurs one period after the investment cost. The agent has a $(\beta, \delta)$ preference, also called a quasi-hyperbolic discounting function whereby cost and benefits at times $0, 1, 2, \ldots, n$, are discounted respectively by rates $1, \beta\delta, \beta\delta^2, \beta\delta^3, \ldots, \beta\delta^n$, with $0 < \beta, \delta \leq 1$. When $\beta = 1$, the model is identical to the constant exponential discounting model. When $\beta < 1$, this model replicates the "hyperbolic discounting" pattern with more periodic discounting in the short run than in the long run. Consider the case $\beta = 1/2$ and $\delta = 1$. Let $C = 100$ and $B = 180$ so that the undiscounted benefit is larger than the undiscounted cost. When evaluated from an earlier perspective $t$ periods before the cost has to be incurred, $t \geq 1$, this investment is desirable because

$$-\beta\delta^t C + \beta\delta^{t+1} B = -\left(\frac{1}{2}\right)100 + \left(\frac{1}{2}\right)180 = 90 > 0. \tag{3.7}$$

But the investment becomes nondesirable when the agent is asked to act immediately:

$$-C + \beta\delta B = -100 + \left(\frac{1}{2}\right)180 = -10 < 0. \tag{3.8}$$

So the agent is faced with conflicting preferences.

When asked to make a binding commitment in advance, the agent will chose to invest. When such precommitment is not feasible, the agent will not invest, since she always reneges from her previous plan when the moment arrives to act. In such a case revealed choices cannot be used to infer the true intertemporal preference of the agent, and cannot be relied upon to make normative judgments. Additional normative assumptions are needed. For instance, we could assume that $\beta = 1$ in order to evaluate the policy choice from the perspective of the rational self of this agent with a no–self-control problem. Bernheim and Rangel (2005) provide a formal justification of this normative criterion based on aggregation principles when the consumer's horizon is sufficiently long. The problem is similar to the welfare aggregation involving many individuals to rank policies: here we aggregate over multiple selves. The idea is that person at time $t$ is a different person at time $t + 1$ due to the preference reversal. As in the problem with multiple consumers, it is possible to apply a multi-person welfare analysis. If the consumers' horizon is sufficiently long, the aggregation is over many different selves. Then the reason for using normative criterion $\beta = 1$ is that the consumer evaluates trade-offs between any two periods $t$ and $t + 1$ by exactly the same discount rate in all periods but one; then the influence of anyone self must decline to zero as the number of selves becomes large. To put it differently, if we aggregate preferences according to the frequency with which rationality prevails, we end up using rational preferences

for normative analysis because the "momentary lapses" of reason approach zero when considered over a very long horizon. Obviously momentary lapses of reason matter for positive analysis because they have long-lasting effects.

## 3.6   Other-Regarding Preferences

Standard economics assumes that agents are rational and selfish. We have already discussed numerous deviations from rationality. We now discussed deviations from selfish behavior. Experimental economics has confirmed the predictions of competitive markets, namely that even with a limited number of participants, experimental markets clear at competitive prices. The experimental results give support to the predictions of the competitive equilibrium model analyzed in chapter 2. The equilibrium outcome was based on selfish optimization by agents interested only in their own material consumption and profits. In contrast in a *strategic environment* (where an individual choice can affect someone else), agents do not seem to act according to the standard model of selfish optimization. For instance, in the simplest ultimatum game that we describe next, the deviations from the standard model are systematic. The deviation is not from rationality but rather from the standard assumption of selfish behavior.

### 3.6.1   Ultimatum Game

In the ultimatum game two players bargain over the distribution of a surplus of fixed size 1. The first player (proposer) chooses any share of the surplus $s \in [0, 1]$. The second player (responder) either accepts or rejects the proposal. If the responder accepts the proposal $s$, then the responder's payoff is $r(s) = s$ and the proposer's payoff is $p(s) = 1 - s$. If the responder rejects the proposal, both receive nothing $p = r = 0$. It is a Nash equilibrium for the proposer to offer $s > 0$ and for the responder to accept any offer greater or equal to $s$. However, such an equilibrium is not fully rational (i.e., subgame-perfect) because it relies on the (noncredible) threat that the responder will reject the positive offers less than $s$. So the unique (subgame-perfect) Nash equilibrium involves the proposer making an offer $s = 0$ (if the set of possible proposals is continuous) and the responder will accept the offer as matter of indifference. So in equilibrium $s^* = 0$, and the payoffs are $r(s^*) = 0$ and $p(s^*) = 1$.

The ultimatum game is simple with sharp predictions: the proposer demands essentially everything and the responder accepts. However, those sharp predictions are

systemically wrong in the sense that they are violated in most experiments. Low proposals giving less than $\frac{1}{5}$ to the responder are rare ($s < 0.2$) and proposals giving more than $\frac{1}{2}$ to the responder are also rare ($s > 0.5$). Equal or almost equal splits are frequent ($s \simeq 0.5$). Proposals are also rejected in some experiments, with the probability of rejection increasing as the responder's share of the surplus $s$ decreases.

In a *competitive environment* the results from the play of ultimatum game experiments are again compatible with the predictions of self-interest. Consider a variation of the ultimatum game with a unique responder but $n > 1$ proposers (i.e., competition among $n$ proposers). The proposers simultaneously make an offer to the unique responder. So the list of offers is $s = (s_1, \ldots, s_i, \ldots, s_n)$. If the responder accepts the offer $s_i$ from proposer $i$, then the responder earns $s_i$, the proposer $i$ earns $1 - s_i$, and other proposers earn nothing. If the responder rejects all offers, then everyone earns nothing. The subgame perfect (Nash) equilibrium involves the responder receiving almost all of the surplus with $s^* \simeq 1$. And experiments confirm this equilibrium prediction of self-interest.

### 3.6.2   Social Preferences

In social organizations, people make friends and enemies, and compare themselves to others. Workers may thus sacrifice some potential earnings to help their friends and harm their enemies, or to create better social comparisons. It is also well documented that the perception of fairness is a key determinant of strike action. If people compare their own wages to those who work in similar activities, then the results may create turnover costs or costs in a social organization.

There are many different forms of social preferences, displaying Selfishness, Altruism, Fairness, or Envy. They can collectively be called SAFE preferences and they depend on how the others' material consumption enters your own utility function. For the sake of clarity, consider an exchange economy with only one good (which we call income) and two individuals. Each individual has preferences, not only over her own income but also over the income of the other. Preferences are complete, transitive, and continuous with the resulting utility for individual $i$ being $v_i(y_i, y_j)$ (where $i \neq j$). Utility is increasing in own income, so $\partial v_i / \partial y_i > 0$. By this formulation, four types of social preferences are possible:

1. *Selfishness*   Utility is independent of the income of the other:

$$\frac{\partial v_i(y_i, y_j)}{\partial y_j} = 0 \qquad \text{for all } y_i, y_j. \tag{3.9}$$

2. *Altruism*     Utility is increasing in the income of the other. Altruism is a form of unconditional kindness to others:

$$\frac{\partial v_i(y_i, y_j)}{\partial y_j} > 0 \qquad \text{for all } y_i, \ y_j. \tag{3.10}$$

3. *Envy*     Utility is decreasing in the income of the other. An envious person always values the material payoff of other negatively. It is a form of unconditional enviousness to others:

$$\frac{\partial v_i(y_i, y_j)}{\partial y_j} < 0 \qquad \text{for all } y_i, \ y_j. \tag{3.11}$$

4. *Fairness*     Utility is either increasing or decreasing in the income of the other, if the other is respectively poorer or richer. It is a form of inequity aversion that can exhibit both altruism or envy to other depending on relative position:

$$\frac{\partial v_i(y_i, \ y_j)}{\partial y_j} \begin{cases} \leq 0 & \text{for all } y_i \leq y_j, \\ > 0 & \text{for all } y_i > y_j. \end{cases} \tag{3.12}$$

If agents have extended preferences depending both on their own monetary payoffs and on the payoffs of others, then the equilibrium outcomes can be reconciled with experimental findings. In fact such extended preferences do a good job in organizing experimental results that are at odds with standard predictions (e.g., as in the ultimatum game described above).

### 3.6.3   Market Impact

What is the impact of other-regarding preferences on the market behavior? Consider the competitive economy of chapter 2 with price-taking consumers and firms but with other-regarding preferences. There are two forms of preference interdependence: (1) consumers do care about the consumption levels of others (consumption externalities), and (2) they do care about the budget possibilities of others (income externalities).

Consequently the utility of individual $i$ will depend both on her own material consumption $x_i$ and on the consumption choices of others $x_{-i}$ as well as her own budget possibilities $b_i$ compared to the budget sets of others $b_{-i}$. Denote by $u_i(x_i, x_{-i}, b)$ the utility of individual $i$ from consumption profile $x = (x_i, x_{-i})$ and the budget profile $b = (b_i, b_{-i})$. Preferences are assumed to be strictly convex and strictly monotone in own consumption. Price-taking behavior in competitive markets implies that own

consumption decisions have no impact on prices and nobody is rationed at the prevailing prices. Also the consumption decisions and budget possibilities of others are taken as given when making own consumption choices. So the price-taking consumer $i$ chooses her own consumption that solves

$$\max_{x_i \in b_i} u_i(x_i, \ x_{-i}, \ b). \tag{3.13}$$

The optimal consumption of individual $i$, $x_i^*(x_{-i}, \ b_i, \ b_{-i})$, is a function of her own budget $b_i$ and the consumption and budget sets of others $(x_{-i}, \ b_{-i})$. If the own consumption choices are independent of the consumption choices and budget sets of others, then we say that consumers behave *as if* they were selfish. This is the case under the following separability assumption:

**Definition 3.1**   (Dufwenberg et al. 2011)   The preferences of consumers $i$ are separable if for all $x$, $x'$ and all $b$, $b'$,

$$u_i(x_i, \ x_{-i}, \ b) \geq u(x_i', \ x_{-i}, \ b) \iff u_i(x_i, \ x_{-i}', \ b') \geq u(x_i', \ x_{-i}', \ b'). \tag{3.14}$$

This separability assumption is required to make a meaningful comparison between the competitive equilibria in an economy with and without other-regarding preferences. Hence agents who care directly about the consumption of others cannot be directly distinguished from selfish agents in their consumption behavior. They look as if they are selfish even though they are not. In a Walrasian equilibrium each firm $j$ maximizes its profits $\pi_j^*$ for given price $p^*$, each consumer $i$ chooses her utility maximizing consumption bundle $x_i^*$ for given budget sets $b^*$, and the budget sets are compatible with equilibrium price $p^*$ in the sense that $b_i^* = \left\{ x_i : p^* x_i \leq p^* \omega_i + \Sigma_j \theta_{ij} \pi_j^*(p^*) \right\}$ with $\omega_i$ the initial endowment of individual $i$ and $\theta_{ij}$ the profit share of individual $i$ in firm $j$. The major theorem concerning the comparison of equilibria now follows.

**Theorem 3.2**   (Dufwenberg et al. 2011)   If all agents have separable preferences that are strictly monotone in own consumption, any Walrasian equilibrium of an economy with other-regarding preferences is a Walrasian equilibrium of the standard economy with selfish preferences.

This result implies that the competitive market outcome is consistent with behavior far more general than selfish optimization. There is a simple intuition for this result. If an agent's decision does not influence the market price or the volume of trade, then he has no opportunity to change the material consumptions of others in the economy. As a

result in a competitive environment agents typically behave as if they care only about their own material consumption, even though they have more general preferences.

There is a dual observation that is more familiar. In certain noncompetitive environments people may act as if they were altruistic, even if they care only about their own material consumption. This is the case in many environments intended to induce cooperation among selfish agents (e.g., think of the Kyoto Protocol on climate change). The fact that market equilibrium may not be affected by other-regarding preferences does not mean that the market's outcome will be efficient (contrarily to the standard competitive economy). In general, the market outcome will be inefficient. Also this result does not exclude the possibility that the market outcome can be affected by other-regarding preferences when there are some forms of market failures such as those described in part III on departures from efficiency.

## 3.7   Conclusions

In this chapter we did not seek to be completely comprehensive, but only to provide an introductory account of the main themes in behavioral economics, and to explore some of the implications for public policy. For that purpose there was no need and no attempt to pursue the theory into every one of its corners. There was no attempt either to be sophisticated in bringing risk or expectation into the theory, except in a very simple way. We also only introduced dynamic considerations in a nonformal manner because this is an issue we analyze in more depth in part VIII of the book.

To conclude this review of behavioral approach to economics, there is one puzzle we would like to share with the reader. This puzzle is reminiscent of the one raised by Little (1956) about the so-called welfare economics revolution. The puzzle is that the conclusions of behavioral economics are important and influential, especially among economists and possibly policy makers, but very few economists are clear as to what the word "behavioral" means, or what precisely the theory is about. Indeed it is only recently that the word "behavioral" has been employed. Time inconsistency, habit formation, satisficing, and social interaction figured in economic analysis long before they were swept under the umbrella of behavioral economics. It is by no means clear what this word means (except nonstandard preferences, beliefs, or choices). It is, to put it differently, not clear what behavioral economics is about. Despite this lack of clarity, the ideas we have discussed have influenced the opinions of many people. It obviously could not have had any such influence if its conclusions had been meaningless, or merely formal and recognized as such. Its conclusions certainly have some real (nonformal) meaning.

It is also fair to say that it is rather difficult to test the theory simply because there are many competing explanations for the mistakes people can make in their decisions. Since the subject matter is often something that arouses peoples' emotions such as eating disorders, illicit drug use, excessive alcohol consumption, smoking, gambling, compulsive consumption, and excessive debt, the result seems to be a lack of balance with the conclusions of the theory being either passionately attacked or passionately defended. Only the future will tell us if the theory is likely to have much direct influence on public policy. In the meantime the theory is likely to have a considerable indirect influence by molding the opinions of economics students and, possibly as a result, some of its fashionable conclusions passing into ordinary language and being taken for granted as though they were the most obvious scientific truth. The automatic enrollment policy is a good illustration of that.

## Further Reading

The starting point of behavioral economics that forcefully brings psychology into the economics of consumer choice is:

Thaler, R. 1980. Toward a positive theory of consumer choice. *Journal of Economic Behavior and Organisation* 1: 39–60.

For the money pump argument see:

Davidson, D., McKinsey, J. C. C., and Suppes, P. 1955. Outlines of a formal theory of value. *Philosophy of Science* 22: 140–60.

The self-control problem is analyzed in:

Strotz, R. H., 1956. Myopia and inconsistency in dynamic utility maximization. *Review of Economics Studies* 23: 165–80.

Thaler, R. H., and Shefrin, H. M. 1981. An economic theory of self-control. *Journal of Political Economy* 89: 392–406.

The $(\beta, \delta)$ model is successively developed in:

Phelps, E. S., and Pollak, R. A. 1968. On second-best national saving and game-equilibrium growth. *Review of Economic Studies* 35: 185–99.

Laibson, D. 1997. Golden eggs and hyperbolic discounting. *Quarterly Journal of Economics* 112: 443–47.

O'Donoghue, T., and Rabin, M. 1999. Doing it now or later. *American Economic Review* 89: 103–24.

The framing effects are in:

Kahneman, D., and Tversky, A. 2000. *Choices, Values, and Frames*. Cambridge: Cambridge University Press.

An overview of the ultimatum game and other bargaining predictions and experiments is in:

Roth, A. 1995. Bargaining experiments. In J. H. Kagel and A. E. Roth, eds., *Handbook of Experimental Economics*. Princeton: Princeton University Press: 253–348.

Conformism is in:

Bernheim, B. D. 1994. A theory of conformism. *Journal of Political Economy* 102: 841–77.

Jones, S. R. G. 1984. *The Economics of Conformism*. Oxford: Basil Blackwell.

The influence of social norms is in:

Elster, J. 1989. Social norms and economic theory. *Journal of Economic Perspectives.* 3: 99–117.

Akerlof, A., and Kranton, R. 2010. *Identity Economics: How Our Identities Shape Our Work, Wages, and Well-Being*. Princeton: Princeton University Press.

Two excellent reviews of the central issues that arise with deviations from rationality:

Bernheim, B. D., and Rangel, A. 2007. Toward choice-theoretic foundations for behavioral welfare economics. *American Economic Review Papers and Proceedings* 97: 464–70.

Della Vigna, S. 2009. Psychology and economics: evidence of the field. *Journal of Economic Literature* 47: 315–72.

Further reading on time discounting:

Ramsey, F. 1928. A mathematical theory of saving. *Economic Journal* 38: 543–49.

The first experiment of the guessing game is in:

Nagel, R. 1995. Unravelling in guessing games: An experimental study. *American Economic Review* 85: 1313–26.

General equilibrium treatment of the other-regarding preference is in:

Dufwenberg, M., Heidhues, P., Kirchsteiger, G., Riedel, F., and Sobel, J. 2011. Other-regarding preferences in general equilibrium. *Review of Economic Studies* 78: 640–66.

Overview of reciprocity and fairness is in:

Fehr, E., and Gächter, S. 2000. Fairness and retaliation: The economics of reciprocity. *Journal of Economic Perspectives* 14: 159–81.

Sobel, J. 2005. Interdependent preferences and reciprocity. *Journal of Economic Literature* 43: 392–436.

## Exercises

**3.1**    Distinguish between the behavioral model of choice and the rational model of choice. Briefly describe a model of each kind.

**3.2**    Intertemporal models in economics typically assume exponential discounting (a constant discount rate). Do such models predict an optimal consumption plan that does not change over time or a optimal consumption plan that changes over time?

**3.3**   You are requested to construct a model to predict the self-control problem of procrastination to start exercising. Assume that exercising has a cost today of $-6$ and a delayed benefit of 8. How would you model the fact that your early plan involves exercising tomorrow but not today? Show that when tomorrow arrives you will again want to postpone action. Discuss the inconsistency problem.

**3.4**   Let $c$ represent calories (or cigarettes), with $u$ the strictly concave (immediate) utility of consumption and $v$ the strictly convex (delayed) cost of consumption. Let the consumer have preferences described by the (intertemporal) utility function

$$U(c_{t-1}, c_t, c_{t+1}, c_{t+2}, \ldots) = \big[u(c_t) - v(c_{t-1})\big]$$
$$+ \beta\delta\big[u(c_{t+1}) - v(c_t)\big]$$
$$+ \beta\delta^2\big[u(c_{t+2}) - v(c_{t+1})\big],$$

where $0 < \beta \leq 1$ and $0 < \delta \leq 1$.

a. Calculate the optimal level of consumption from the perspective of date $t$.

b. Assuming that $\beta = 1$, show that $c_t = c_{t+1}$.

c. Assuming that $\beta < 1$, show that $c_t > c_{t+1}$.

d. For $\beta < 1$, what is the effect of increasing the self-control parameter $\beta$ upon the consumption $c_t$? Explain your finding.

**3.5**   Consider a naif with $\beta = \frac{1}{2}$ and $\delta = 1$. The naif has to finish a project by deadline $T$. At date $t$, the undiscounted project costs $\left(\frac{3}{2}\right)^t$ to implement. When will the naif undertake the project?

**3.6**   Consider the same project as in previous exercise but now with a sophisticate.

When will a sophisticate undertake the project?

**3.7**   "A horse! A horse! My kingdom for a horse!" (Richard III, scene iv.) Is this an example of extreme present bias?

**3.8**   Search and procrastination (Carroll et al. 2010). Let a "naive" consumer have discount parameters $0 < \beta < 1$ and $\delta = 1$ (daily discounting so that $\delta \simeq 1$). The daily loss from delay is $L$ (the daily benefit loss). The search cost $c_t$ is stochastic and drawn from a uniform distribution on the interval $[0, 1]$. Let $W$ describe the cost function today, with

$$W(c) = \begin{cases} c & \text{if act now,} \\ \beta[L + EV(c')] & \text{if wait,} \end{cases}$$

where $V$ represents the "exponentially discounted" ($\delta = 1$) cost function tomorrow

$$V(c) = \begin{cases} c & \text{if act tomorrow,} \\ L + EV(c') & \text{if wait tomorrow.} \end{cases}$$

a. Calculate the optimal stopping rule; that is, the equilibrium cost cutoff $c^*$ such that the "naive" agent is indifferent between acting and waiting at this cost cutoff. (*Hint*: Solve two equations with two unknowns $c^*$ and $EV$.)

b. How does $c^*$ change with $L$?

c. How does $c^*$ change with $\beta$, the short-term discount factor?

    d. Is $c^* > 0$? And is $c^* < 1$?

    e. If $L = \frac{1}{2}$ what is the probability of procrastination?.

**3.9**      Consider the same model as in previous exercise, but with a sophisticated agent. That is, assume that the short-term discount factor is $\beta = 1$.

    a. What is the equilibrium cost cutoff $c^{**}$?

    b. How does $c^{**}$ compare with $c^*$?

    c. If $L = \frac{1}{2}$, what is the probability of procrastination?

**3.10**      (Overconfidence). A worker chooses effort $e$, which has cost $c(e)$ that is increasing and convex in $e$. The productivity of effort also depends on a variable called skill $s$ so that output is

$$x = f(e, s) + \theta,$$

where $f(e, s)$ is increasing and concave in each component, and $\theta$ is a random term that can be called "luck." The random term $\theta$ is distributed according to $\mu(\theta)$. The firm pays a linear wage based on output

$$w(x) = w_0 + \alpha x.$$

Assuming separability between effort cost and wage benefit, the expected utility is

$$EV(e) = \int_\theta U(w(f(e, s) + \theta))\mu(\theta)d\theta - c(e).$$

    a. Suppose that skill does not matter and that $f(e, s) = e$. Compute the optimal effort choice.

    b. Now assume that the worker overestimates the marginal productivity of effort by a factor $\beta > 0$ so that $\hat{f}(e, s) = (1 + \beta)e$. Show that this overconfident worker may cut back on effort.

**3.11**      Consider the previous exercise with $f(e, s)$ increasing and concave and with a nonzero cross derivatives $f_{e,s} = \frac{\partial^2 f(e, s)}{\partial e \partial s} \neq 0$. $f_{e,s} > 0$ if $e$ and $s$ are complements, and $f_{e,s} < 0$ if $e$ and $s$ are substitutes.

    a. Compute the optimal effort choice conditional on the skill level.

    b. Now assume the worker overestimates her skill value by a factor $\Delta$ so that $\hat{s} = s + \Delta$. Show that whether a worker who is overconfident in this way will cut back on work (or do the opposite) depends on whether effort and skill are complementary.

**3.12**      (Beauty contest). "…professional investment may be likened to those newspaper competitions in which the competitors have to pick out the six prettiest faces from a hundred photographs, the prize being awarded to the competitor whose choice most nearly corresponds to the average preferences of the competitors as a whole; so that each competitor has to pick, not those faces which he himself finds prettiest, but those which he thinks likeliest to catch the fancy of the other competitors, all of whom are looking at the problem from the same point of view. It is not a case of choosing those which, to the best of one's judgment, are really the prettiest, nor even those which average opinion genuinely thinks the prettiest. We have reached the third degree where we devote our intelligences to anticipating what

average opinion expects the average opinion to be. And there are some, I believe, who practice the fourth, fifth and higher degrees." (John Maynard Keynes, *The General Theory of Employment, Interest, and Money*).

In reference to the quote above, consider the following guessing game. There are 100 participants, and each participant announces independently and simultaneously a number between 0 and 100 (only integers are allowed). There is a fee of $1 to participate to the guessing game, to be redistributed as a prize of $100 to the winner. The participant who guessed closest to a target wins the prize.

a. Suppose that the target is one-half of the average guess. Would you participate in this game? Explain why.

b. Suppose that the target is one-half of the average guess. What would be your best guess? Explain briefly.

c. What would we expect if everyone is rational and thinks everyone else is rational? Compare to your answer in part b.

**3.13**   (Ultimatum game). Consider a fixed income of $y$ to share between two players. Player 1 is the proposer and demands $x_1$ for himself. Player 2 who is the responder accepts or rejects the demand. If the demand is accepted, the payoffs are $v_1 = x_1$ and player 2's payoff is $v_2 = y - x_1$.

a. What is the subgame-perfect Nash equilibrium if both players are selfish?

b. What would happen if both players are fair?

**3.14**   The utility function of individual $i$ is given by $u_i(y_i, y_j) = \log y_i + \lambda \log y_j$, where $y_i$ is individual $i$'s income, $y_j$ is individual $j$'s income, and $\lambda$ is a parameter representing the other-regarding preference.

a. For what values of $\lambda$ is individual $i$ selfish?

b. For what values of $\lambda$ is individual $i$ an altruist?

c. For what values of $\lambda$ is individual $i$ fair?

d. For what values of $\lambda$ is individual $i$ envious?

**3.15**   (Envy model). Consider a 2-player model of envy with the utility function of individual $i$ given by $u_i(y_i, y_j)$, where $\partial u_i / \partial y_i > 0$ and $\partial u_i / \partial y_j < 0$. Show that any division of an object of fixed size in which the responder gets nothing will be necessarily rejected by the responder.

**3.16**   Consider a two-stage version of the ultimatum game with a shrinking amount to be distributed among two selfish players. The amount to be distributed in stage 1 is 1. Player 1 demands $x_1$ and player 2 accepts or rejects the demand. If player 2 accepts, he receives $1 - x_1$. If player 2 rejects, he can make a counterdemand $x_2$ and player 1 can accept or reject, but then the amount to be distributed is reduced to $\lambda < 1$. If player 1 accepts the counterdemand $x_2$, he receives $\lambda - x_2$. If player 1 rejects the counterdemand, both players receive 0.

a. Compute the subgame-perfect Nash equilibrium.

b. Show that the equilibrium demand by player 1 is decreasing in $\lambda$. Explain.

**3.17**   Repeat the previous two-stage ultimatum game with a shrinking amount to be distributed, but now assume a linear model of envy with utility of player $i$,

$$u_i(y_i, y_j) = y_i - \alpha y_j,$$

with the envy preference parameter $\alpha > 0$.

a. Compute the subgame-perfect Nash equilibrium

b. Show that the equilibrium demand by player 1 is decreasing in $\lambda$ if $\alpha < \lambda$. Explain.

c. Show that an unfavorable counteroffer is possible; that is, in equilibrium $x_2 < 1 - x_1$. Explain.

**3.18**    What is the status of behavioral models?