# A Theory of Addiction[†]

Faruk Gul

and

Wolfgang Pesendorfer

Princeton University

January 2001

## Abstract

We construct an infinite horizon consumption model and use it to define and analyze addiction. Consumption is *compulsive* if it differs from what the individual would have chosen had commitment been available. A good is *addictive* if its consumption leads to more compulsive consumption of the same good in the future. We analyze two types of drug policies. A policy is *prohibitive* if it decreases the maximally feasible drug consumption. We show that prohibitive policies make agents better off and - if they are not binding - lead to higher drug demand. A *price policy* is one that increases the opportunity cost of drug consumption without changing the maximally feasible drug consumption. We show that price policies make the agent worse-off and decrease drug demand if the drug is a normal good.

# 1. Introduction

Substantial resources are spent to reduce the availability of and the demand for drugs. These efforts are justified by the belief that addiction is a serious health and social problem. This belief is supported by distressing descriptions of the life of a typical drug addict and a large number of deaths attributed to nicotine, alcohol, opiate or cocaine/amphetamine addiction. There are however, many other goods whose consumption is dangerous or associated with an unattractive life style. With the exception of a few psychothropic substances these properties are not considered sufficient reasons for banning a substance, let alone spending billions on enforcing the ban. What, if anything, is special about drugs that could justify restricting its supply and demand?

Standard economic analysis uses the individuals' choice behavior as a welfare criterion. Alternative $x$ is deemed to be better for the agent than alternative $y$ if and only if given the opportunity, the agent would choose $x$ over $y$. While typical in economic analysis, the identification of welfare and choice is certainly not the norm in discussions of addiction. Instead, addiction is often viewed as a disease that inflicts the agent's decision-making ability.[1] It is believed that after being struck by the disease, a person can no longer be trusted to make the right decision for his "true" self.[2] The role of intervention is to "cure" (i.e. induce abstinence) or at least "control" (i.e. reduce consumption) the disease.

Viewing addiction as a disease creates a wedge between choice and welfare. This wedge makes room for desirable interventions that modify the addict's choices but also creates the need for a new welfare criterion. Consider a costly treatment that, if successful, will remove the agent's drug dependency (i.e. cure the disease). If the probability of success is sufficiently high then the treatment is desirable regardless of whether the agent thinks so or not. Conversely, if the probability of success is sufficiently small then the treatment is undesirable. How can the planner determine whether the probability of success justifies the cost of the treatment?

---

[1] "Is alcoholism a disease? Yes. Alcoholism is a chronic, often progressive disease with symptoms that include a strong need to drink despite negative consequences, such as serious job or health problems." (cited from: National Institute on Alcohol Abuse and Alcoholism. http://silk.nih.gov/silk/niaaa1/questions/q-a.htm#question2)

[2] There are numerous criticisms of the disease model of drug addiction (see for example, Davies (1992)).

In this paper, we provide a model of addiction that is consistent with the view that addicts may benefit from interventions that modify their choices. At the same time, the model offers clear guidance for welfare comparisons. Building on previous work (Gul and Pesendorfer (2000a)), we assume that the agent may have a preference for commitment; that is, his welfare may go up when some alternatives are eliminated from his set of choices. We refer to options that the agent would rather not have as *temptations*. A temptation lowers the agent's utility either because it distorts his choice or because it necessitates costly self-control. In the latter case, the agent does not choose the tempting alternative but its availability makes him worse off. Thus, our model allows welfare to depend both on what the individual chooses and on the set of options from which the choice is made.

To see how our model works, consider an agent who must choose from the choice problem $z$. Each element of $z$ is of the form $(c, x)$, where $c$ is a consumption vector that includes the drug and $x$ is the (continuation) choice problem for the next period. We capture the dynamic nature of addiction by allowing past consumption to affect current preferences. The agent's preferences are defined over choice problems and can be represented by the utility function $W$ where

$$W(s, z) = \max_{\{(c,x) \in z\}} \{u(s, c) + \delta W(s', x) + V(s, (c, x))\} - \max_{\{(c', y) \in z\}} V(s, (c', y))$$

Past consumption determines the state $s$ in the current period. Next period's state, $s'$, is determined jointly by the current state $s$ and current consumption $c$. The function $V$ represents the agent's temptation while $u + \delta W$ is his commitment utility; that is, $u + \delta W$ describes what the agent would do in the absence of temptation. If all options in $z$ are equally tempting, then the $V$-terms in the representation above drop out. Therefore, such consumption problems are evaluated according to $u + \delta W$. In particular, if $z$ consists of a single choice $(c, x)$; that is, if the agent were able to commit to $(c, x)$ in some previous period, then the overall utility of the current choice problem is the commitment utility $u + \delta W$, of the singe option $(c, x)$.

The individual's choice $(c, x)$ maximizes $u + \delta W + V$. This choice reflects the compromise between the commitment utility and temptation. We say that an individual is *compulsive* if his choice does not maximize the commitment utility and hence temptation

distorts his choice. A drug is *addictive* if an increase in drug consumption leads to more compulsive drug consumption in the future. Thus, we define addiction as a widening of the gap between the individual's choice and what he would have chosen before experiencing temptation.

As in standard models, choice and welfare are synonymous in our model. Therefore, we may elicit how much a "treatment program" is worth by confronting the individual with the appropriate choices. Suppose we give the individual the option to plan drug consumption one or more periods in advance. We can infer the social value of this commitment opportunity by asking the agent how much consumption he would be willing to give up in exchange for the commitment option.

Our model suggests that addicts should seek commitment opportunities. We observe such behavior in the form of enrollment in voluntary rehabilitation programs. For example, consider an addict who seeks treatment for alcohol addiction and is given the drug disulfiram. Disulfiram is a deterrent medication that is used to fight alcohol addiction. Disulfiram produces a sensitivity to alcohol which results in a highly unpleasant reaction when the patient under treatment ingests even small amounts of alcohol. This effect lasts up to 2 weeks after ingestion of the last dose.[3] Hence, the patient is committed to abstaining from alcohol as long as the drug is effective (Chick 1992). Similarly, the opiate antagonist naltrexone blocks the opioid receptors in the brain and hence the euphoric effects of these drugs for up to 3 days after the last dose. Naltrexone is voluntarily used by some heroin and morphine addicts.

Further evidence for the demand for commitment devices are the recent efforts by pharmaceutical companies to develop *vaccines* for nicotine (Pentel, et al. (2000)) and cocaine.[4] The function of these vaccines is to prevent the drug from reaching the brain, so

---

[3] "Disulfiram plus even small amounts of alcohol produces flushing, throbbing in head and neck, throbbing headache, respiratory difficulty, nausea, copious vomiting, sweating, thirst, chest pain, palpitation, dyspnea, hyperventilation, tachycardia, hypotension, syncope, marked uneasiness, weakness, vertigo, blurred vision, and confusion. In severe reactions, there may be respiratory depression, cardiovascular collapse, arrhythmias, myocardial infarction, acute congestive heart failure, unconsciousness, convulsions, and death." (cited from: http://www.mentalhealth.com/drug/)

[4] "When injected in laboratory animals, the vaccine stimulates the immune system to produce antibodies that bind tightly to nicotine. The antibody-bound nicotine is too large to enter the brain, thereby preventing nicotine from producing its effects. The antibody-bound nicotine is eventually broken down to other harmless molecules." cited from http://pharmacology.about.com/health/pharmacology/library/99news/bl9n1217a.htm

as to eliminate its effects and provide commitment for individuals. A novel feature of these vaccines is their long term effectiveness, and hence their ability to provide commitment over many months.

The economics literature has typically identified addiction with inter-temporal complementarities. Becker and Murphy (1986) view the consumption of an addictive good much like an investment that increases the return of future consumption. The preferences analyzed by Becker and Murphy are "standard" in the sense that individuals can never benefit from the elimination of some alternatives. Therefore, an individual who voluntarily acquires costly commitment devices such as the drugs described above is inconsistent with the Becker and Murphy preferences.

In Becker and Murphy's treatment of addiction, drug consumption is never "bad" in terms of individual welfare, and hence their model leaves no room for a drug policy. However, Becker and Murphy do distinguish between addictions that are harmful and those that are beneficial: an addiction is harmful if it leads to a utility penalty in future periods. However, the mere fact that the agent chooses to become addicted implies that the addiction's *net* effect on utility is positive. Becker and Murphy's distinction between a harmful addiction and a beneficial habit is based on *when* utility is experienced. However, choice experiments cannot identify when an individual experiences the utility of a given choice. Therefore, this distinction does not have behavioral content. Empirically distinguishing harmful addictions and beneficial habits as defined by Becker and Murphy would require a direct measurement of utility flows.

O'Donoghue and Rabin (1997) offer a model of addiction that merges the approach of Becker and Murphy with hyperbolic discounting. In their model, the individual may consume more than his past selves would like because of a presence-bias in his preferences. As in our approach, this model implies that agents will utilize commitment opportunities (at least if they are sophisticated). However, their notion of a harmful addiction is based on that of Becker and Murphy and therefore relies on hedonistic utility. Moreover, the multi-selves view of the agent implies that the decision to get addicted benefits the current self but typically, harms future selves. Therefore, revealed preference information is no longer sufficient to identify what is good for the agent. In such cases it is difficult to devise a criterion for evaluating treatment and policy alternatives.

4

To analyze the welfare and demand effects of policy alternatives, we consider a special case of our model in which there is a single tempting good ("the drug") and consumption of the drug is addictive. Moreover, we assume that the drug is "bad"; that is, the commitment utility is decreasing in drug consumption. Thus, if the agent could commit to a consumption path, he would never consume the drug. However, drug consumption is tempting and, as a result, the agent may consume the drug when commitment is not possible.

A typical drug policy affects agents along two dimensions. First, the drug policy may have a prohibitive effect, that is, the policy may reduce the maximally feasible level of drug consumption. For example, drug enforcement efforts may occasionally interrupt the supply of drugs. Second, drug policies may have a price effect. That is, the cost of the drug may change as a result of the drug policy. Often drug consumption is a relatively small part of an agent's budget and opportunity cost of drug consumption goes up without affecting the maximally feasible drug consumption in the current period. We refer to such a policy as a *price policy*. Similarly, when the maximally feasible drug consumption is reduced without affecting the opportunity cost of drug consumption we say the policy is a *prohibitive policy*.

A prohibitive policy always makes the agent better-off. By contrast, a price policy always makes the agent worse-off. The reason is that a prohibitive policy offers some commitment for the agent without affecting the feasible consumption of goods other than the drug. On the other hand, a price policy offers no commitment because it does not change the most tempting alternative. Yet, for a given level of drug consumption the price policy reduces the consumption of other goods and hence leads to lower welfare.

We also examine the demand effects of price and prohibitive policies in simple stationary choice problems. Clearly, when a prohibitive policy is binding and makes the desired level of drug consumption infeasible it leads to lower drug consumption. We show however, that if a prohibitive policy is not binding, it leads to *higher* drug consumption. The reason is that by providing future commitment opportunities the prohibitive policy makes it less costly to get addicted. In contrast, a price policy decreases drug demand when the drug is a normal good. As in standard consumer theory the demand effect of a price policy is in general ambiguous.

Our analysis shows that the demand and welfare effects of drug policies may move in opposite directions even when a drug is bad; that is, when the agent would not consume the drug if costless commitment were available. The fact that a drug policy does not decrease drug demand does not imply that that policy is not successful in terms of welfare. Conversely, a policy that successfully reduces drug demand may be harmful.

The paper is organized as follows. Section 2 introduces the model of preferences and provides the definition of compulsive consumption. Section 3 defines and characterizes addiction. Section 4 examines the positive and normative implications of policies. Finally, Section 5 provides axioms for the utility functions used in the earlier sections.

## 2.   SSC Preferences and Compulsive Consumption

There are $l$ goods and $C = [0,1]^l$ is the set of possible consumption vectors. We consider an agent who is confronted with a dynamic choice problem. Every period $t = 1, 2, \ldots$ the agent must take an action. This action results in a consumption for period $t$ and constrains future actions.

A deterministic dynamic choice problem can be described recursively as a set of alternatives, each yielding a current consumption and a continuation choice problem.[5] We use $\bar{Z}$ to denote deterministic choice problems. Each $z \in \bar{Z}$ is a (compact) set of alternatives of the form $(c, x)$ where $c$ denotes the current consumption and $x \in \bar{Z}$ denotes the continuation problem. A broader class of choice problems, $Z$, allows for uncertainty. In that case, the agent chooses among lotteries over current consumption and continuation choice problems. We use $x, y$ or $z$ to denote generic choice problems (elements of $Z$ or $\bar{Z}$). Generic choices (elements of a given $z$) are denoted $\mu, \nu$ or $\eta$ and constitute probability distributions over $C \times Z$. The degenerate lottery that yields with certainty the current consumption $c$ and the continuation problem $x$ is denoted $(c, x)$.[6]

The set of choice problems $Z$ serves as the domain of preferences for the agent. This allows us to describe agents who struggle with temptation. For example, the agent may strictly prefer a choice problem in which some alternatives are unavailable because these

---

[5]  See Gul and Pesendorfer (2000b) for a detailed discussion of dynamic choice problems.
[6]  For most of the analysis, we restrict to deterministic choice problems. This is done for notational simplicity. However, some of our results utilize lotteries and hence we need to consider choice problems that include uncertainty.

alternatives present temptations that are hard to resist. Even when the agent makes the same ultimate choice from two distinct choice problems he may have a strict preference for one choice problem because making the same choice from the other requires more self-control. Below, we represent the individual's preferences by a utility function. This utility function is analogous to the *indirect utility function* in standard consumer theory. The difference is that the traditional indirect utility function is defined only for choice problems that can be represented by a budget set while our utility function is defined for all choice problems.[7]

The preferences analyzed in this paper depend on the agent's past consumption. To capture this dependence, we index the individual's preferences by $s \in S$, the state in the initial period of the choice problem. The state $s$ represents the relevant consumption history prior to the initial period of analysis. We assume that there is a finite number $K$ such that consumption in only the last $K$ periods influences the agents preferences. Therefore, $S := C^K$. We refer to the indexed family of preferences $\succeq := \{\succeq_s\}_{s \in S}$ simply as *the agent* or the preference $\succeq$. For any state $s = (c^1, \ldots, c^K)$ and $c \in C$, let $sc$ denote the state $(c^2, \ldots, c^K, c)$. We say that the utility function $W : S \times Z \to \mathbb{R}$ represents the preference $\succeq$ if, for all $s$, $x \succeq_s y$ iff $W(s, x) \geq W(s, y)$.

Section 5 provides axioms under which $\succeq$ can be represented by a continuous function $W$ of the following form. There are continuous utility functions $u : S \times C \to \mathbb{R}$, $V : S \times (C \times Z) \to \mathbb{R}$ and a discount factor $\delta \in (0, 1)$ such that for $z \in \bar{Z}$

$$W(s, z) = \max_{(c,x) \in z} \{u(s, c) + \delta W(sc, x) + V(s, (c, x))\} - \max_{(c',y) \in z} V(s, (c', y)) \tag{1}$$

To understand this representation, first consider a choice problems that offers commitment; that is, a choice problems with one alternative, $\{(c, x)\}$. In that case, the $V$-terms drop out and $W(s, \{(c, x)\}) = u(s, c) + \delta W(sc, x)$. Therefore, we say that the function $U := u + \delta W$ represents the agent's commitment utility.

Next, consider a choice problems with two alternatives, $\{(c, x), (c', y)\}$. Assume that $U(s, (c, x)) > U(s, (c', y))$ and $V(s, (c, x)) < V(s, (c', y))$. Then, if follows from equation (1)

---

[7] For a detailed definition of the class of choice problems captured by $Z$, see Gul and Pesendorfer (2000b).

that $W(\{(c,x)\} > W(\{(c,x),(c',y)\})$. Hence, the agent strictly prefers the choice problem where only $(c,x)$ is available to the choice problem where, in addition, $(c',y)$ is available. We interpret this as a situation where the agent suffers from the temptation presented by $(c',y)$ and conclude that $V$ measures this temptation. The agent compromises between commitment utility $U$ and temptation utility $V$ by maximizing $U + V$. That is, $U + V$ governs the agent's choice from $z$.

Let $(c,x)$ be the alternative that maximizes $U(s,\cdot) + V(s,\cdot)$ in $z$. By choosing $(c,x)$ the agent enjoys the utility $U(s,(c,x))$ but also incurs a self-control cost equal to $-[V(s,(c,x)) - \max_{(c',y)\in z} V(s,(c',y))]$. This cost is zero if $(c,x)$ also maximizes $V$ in $z$ and positive otherwise. By maximizing $U + V$ the agent ensures an optimal trade-off between the commitment utility and avoiding self-control costs.

We refer to preferences that can be represented by a utility function that satisfies equation (1) as *stationary self-control* (SSC) preferences. Straightforward application of results from dynamic programming imply that for every $(u,V,\delta)$ with $u$, $V$ continuous, there is a unique $W$ that satisfies equation (1). We say that $(u,V,\delta)$ represents the SSC preference $\succeq$ if the the unique $W$ that satisfies equation 1 represents $\succeq$. An SSC preference $\succeq$ is *regular* if for all $s$, $U(s,\cdot)$ is not constant and $V(s,\cdot)$ is not an affine transformation of $U(s,\cdot)$.

The three main concepts of this paper are *preference for commitment*, *compulsive consumption* and *self-control*. To illustrate these concepts, consider a deterministic choice problem $z \in \bar{Z}$ and assume that $(c,x)$ is the unique maximizer of the commitment utility $U(s,\cdot)$ in $z$ whereas $(c',y)$ is the unique maximizer of the temptation utility $V$ in $z$. We refer to alternatives that have higher temptation utility than $(c,x)$ as *temptations* and hence $(c',y)$ is a temptation in $z$. From equation (1) we can infer that removing the temptation $(c',y)$ from the set $z$ increases the agent's welfare. Hence, temptations create a *preference for commitment*; that is, situations where the agent strictly prefers the choice problem $x$ over $z$ even though $x$ offers fewer choices; that is, $x \subset z$. The agent is *compulsive* if he does not choose the maximizer of the commitment utility $(c,x)$. For example, the agent is compulsive if he succumbs to temptation and chooses $(c',y)$. The agent exercises *self-control* if he does not choose the most tempting alternative $(c',y)$. Note that if the

8

agent chooses neither $(c, x)$ nor $(c', y)$ he chooses a compulsive consumption and exercises self-control at the same time.

Below, we formally define preference for commitment ($\mathbf{P}$), compulsive consumption ($\mathbf{C}$) and self-control ($\mathbf{S}$). To define these concepts for all choice problems including those with stochastic choices we need the following notation. For any probability distribution $\mu$ on $C \times Z$ we define $U(s, \mu) := \int U(s, (c, z)) d\mu(c, z)$ to be the expected commitment utility and $V(s, \mu) = \int V(s, (c, z)) d\mu(s, z)$ to be the expected temptation utility. For any function $f : S \times Z \to \mathbb{R}$, define $\mathcal{C}_f(s, z) := \{\mu \in z : f(s, \mu) \geq f(s, \nu), \forall \nu \in z\}$ to be the $f(s, \cdot)$ maximizers in $z$. Hence at state $s$, $\mathcal{C}_{U+V}(s, z)$ denotes the agent's optimal choices from $z$, while $\mathcal{C}_U(s, z)$ denotes the commitment utility maximizers and $\mathcal{C}_V(s, z)$ denotes the most tempting alternatives in $z$.

**Definition:** *Let $(u, V, \delta)$ represent the regular SSC preference $\succeq$ and let $U$ be the corresponding commitment utility. Then,*

*(i) $\succeq_s$ has $\mathbf{P}$ at $z$ iff $\mathcal{C}_U(s, z) \cap \mathcal{C}_V(s, z) = \emptyset$.*

*(ii) $\succeq_s$ has $\mathbf{S}$ at $z$ iff $\mathcal{C}_{U+V}(s, z) \cap \mathcal{C}_V(s, z) = \emptyset$.*

*(iii) $\succeq_s$ is $\mathbf{C}$ at $z$ iff $\mathcal{C}_{U+V}(s, z) \backslash \mathcal{C}_U(s, z) \neq \emptyset$.*

Part (i) of the definition says that the agent has a preference for commitment ($\mathbf{P}$) if the choice problem contains temptations. Part (ii) says that the agent has self-control ($\mathbf{S}$) if he does not choose one of the alternatives that maximizes the temptation utility. Part (iii) says that the agent is compulsive ($\mathbf{C}$) if at least one of his optimal choices does not maximizes the commitment utility.

In the appendix we show that the commitment and temptation utilities associated with a regular SSC preference are unique up to a common affine transformation. Consequently, for regular SSC preferences, $\mathbf{P}$, $\mathbf{S}$ and $\mathbf{C}$ are properties of the preference and not of the particular representation.

Below, we offer criteria for ranking states with respect to preference for commitment, self-control and compulsive consumption. We say that the preference has more preference for commitment at $s$ than at $s'$ if $\succeq_s$ has preference for commitment at $z$ implies $\succeq_{s'}$ has preference for commitment at $z$. A similar definition yields a ranking of states with respect to self-control and compulsive consumption.

**Definition:** *An SSC preference $\succeq$ has more* **P** *at $s$ than at $s'$ (denoted $s\mathbf{P}s'$) if $\succeq_{s'}$ has* **P** *at $z$ implies $\succeq_s$ has* **P** *at $z$. Similarly, $s\mathbf{S}s'$ ($s\mathbf{C}s'$) if $\succeq_{s'}$ has* **S** *(is* **C***) at $z$ implies $\succeq_s$ has* **S** *(is* **C***) at $z$.*

Preference for commitment, self-control and compulsive consumption lead to partial orders over states, where a given pair of states $s$ and $s'$ may not be ranked. However, Proposition 1 below establishes that if the agent has more preference for commitment and less self-control at $s$ than at $s'$ then he is more compulsive at $s$ and than at $s'$. Less self-control implies fewer instances where the most tempting alternative is not chosen. Hence, the agent's choice shifts towards alternatives preferred by the temptation utility. Greater preference for commitment implies a greater divergence between commitment and temptation utilities. Together these two effects imply that there are more instances where the choice does not maximize the commitment utility.

**Proposition 1:** *For any regular SSC preference $s\mathbf{P}s'$ and $s'\mathbf{S}s$ implies $s\mathbf{C}s'$.*

The notion of compulsive consumption plays a central role in the clinical definition of addiction and in the definition of addiction we present in the next section. What distinguishes addiction from other types of compulsive behavior is the fact that the compulsive consumption of an addictive substance is "caused" (or made worse) by past consumption of the *same* substance. In order to focus on this dependency of compulsive consumption on the past consumption, in the next section, we restrict attention to a subset of SSC preferences, called simple SSC preferences.

## 3. Simple SSC Preferences and Addiction

In this section we study a subset of SSC preferences, referred to as *simple* SSC preferences. The subset is characterized by restrictions on the temptation utility and by assumptions on how past consumption can influence preferences.

For a simple SSC preference, the temptation utility depends on current consumption of good $l$ only. Thus, temptation is myopic and focused entirely on one good, referred to below as "the drug". This assumption allows us to steer clear of issues related to "cross-addiction", where consumption of one substance affects future preferences for another substance.

Agents with simple SSC preferences can rank all states according to their preference for commitment ($\mathbf{P}$); that is, for every two states $s, s'$, $s\mathbf{P}s'$ or $s'\mathbf{P}s$. Recall that in general, preference for commitment induces a partial order on the states. Here, we assume that this order is complete. This assumption allows us to parameterizes preference for commitment, self-control and compulsive consumption. As we show in section 5, this assumption also implies that the preferences in the next period can be affected by current consumption but not by consumption in past periods. In section 5 we provide an axiomatic characterization of simple SSC preferences.

It is often convenient to distinguish between the drug and the non-drug goods. Therefore, we sometimes write $(b, d) \in C$ instead of $c$, where $b$ denotes the first $l - 1$ coordinates of $c$ and $d$ is the drug coordinate of $c$.

A simple SSC preference can be represented by a function $W$ of the form

$$W(s_l, z) = \max_{(b,d,x) \in z} \{u_0(b, d) + \sigma(s_l)v_0(d) + \delta W(d, x)\} - \max_{(b',d',y) \in z} (\pi(s_l) + \sigma(s_l))v_0(d') \quad (2)$$

where $s_l$ is the $l$'th coordinate of $s$. Without risk of confusion, we will omit the subscript $l$, but it will be understood that for simple SSC preferences, the state $s$ is last period's drug consumption. The discount factor is $\delta \in (0, 1)$; $u_0$ is a continuous real valued function on $C$; $v_0, \pi, \sigma$ are continuous real valued functions on $[0, 1]$, $v_0$ is strictly increasing $\pi \geq 0$ and $\pi + \sigma > 0$.

For the simple SSC preference represented in equation (2), the commitment utility is given by $u_0 - \pi v_0 + \delta W$ and the temptation utility is given by $V = (\pi + \sigma)v_0$. Setting

$$U = u + \delta W = u_0 - \pi v_0 + \delta W$$

$$V = (\pi + \sigma)v_0$$

we can verify that a simple SSC preference is indeed a SSC preference. Note that the temptation utility depends only on current drug consumption. A simple SSC preference is regular iff there is no state $s$ such that $U(s, \cdot)$ is constant.[8] We identify a simple SSC preference with the functions and discount factor $(u_0, v_0, \pi, \sigma, \delta)$ used in its representation.

---

[8] Recall that an SSC preference is regular if for all $s$, $U(s, \cdot)$ is not constant and $V(s, \cdot)$ is not an affine transformation of $U$. For a simple SSC preference $V$ depends on current drug consumption only and is strictly increasing in current drug consumption. When $U(s, \cdot)$ is not constant it depends on future consumption and hence, $V$ is not an affine transformation of $U(s, \cdot)$.

For simple SSC preferences, the impact of past consumption is measured by the functions $\pi$ and $\sigma$. An increase in $\sigma$ implies that $U + V$ puts more weight on $v_0$. Since $U + V$ governs behavior, this suggests a loss of self-control. An increase in $\pi$ implies that increase in the gap between the commitment utility $U$ and the temptation utility $V$ and hence suggests an increase in the preference for commitment. Proposition 2 confirms this intuition.

**Proposition 2:**   Let $(u_0, v_0, \pi, \sigma, \delta)$ be a regular, simple SSC preference. Then,

$(i)$ $s\mathbf{P}s'$ iff $\pi(s) \geq \pi(s')$

$(ii)$ $s'\mathbf{S}s$ iff $\sigma(s) \geq \sigma(s')$

$(iii)$ $s\mathbf{C}s'$ iff $\pi(s) \geq \pi(s')$ and $\sigma(s) \geq \sigma(s')$.

Given $(i)$ and $(ii)$, the only if part of $(iii)$ follows from Proposition 1. Proposition 2 establishes that for regular, simple SSC preferences, more $\mathbf{P}$ and less $\mathbf{S}$ is in fact equivalent to more $\mathbf{C}$.

Psychologists and health professionals commonly refer to an individual as addicted if, after repeated self-administration of a drug, the individual develops a pattern of compulsive drug seeking and drug-taking behavior.[9] The clinical definition emphasizes a lack of control on the part of addicted subjects and suggest a conflict between what the addict *ought to consume* and what he *actually consumes*.

In our model, the agent is compulsive when the choice (the $U + V$ maximizer) is different from the $U$ optimal alternative. Thus, an agent is compulsive if behavior would change were commitment possible. Similar to the clinical definition above, we define an increase in drug consumption to be addictive if higher current drug consumption makes the more compulsive; that is, following the increase in drug consumption there are more situations in which the agent makes a choice that does not maximize $U$. Note that the state in the following period is equal to the drug consumption is the current period. Hence, if current period drug consumption increases from $d$ to $d + \epsilon$ then the next period's state increases from $s = d$ to $s + \epsilon$.

---

[9] See Robinson and Berridge (1993), pg. 248.

**Definition:** *An $\epsilon > 0$ increase in drug consumption is addictive at $d$ if, for $s = d$ and $s' = d + \epsilon$, $s'\mathbf{C}s$ and $\succeq_{s'} \neq \succeq_s$. The drug is addictive if every $\epsilon > 0$ increase is addictive at every $d$.*

Note that for any simple SSC preference, $\pi(s) = \pi(s')$ and $\sigma(s) = \sigma(s')$ implies that the agent's preference in states $s$ and $s'$ are the same. Hence, the following characterization of addiction follows immediately from Proposition 2.

**Corollary 2:** *Let $(u_0, v_0, \pi, \sigma, \delta)$ be a regular, simple SSC preference. An $\epsilon > 0$ increase in drug consumption is addictive iff, for $s = d$ and $s' = d + \epsilon$,*

$(i)$ $\pi(s') \geq \pi(s)$, $\sigma(s) \geq \sigma(s')$ and

$(ii)$ $\pi(s') + \sigma(s') > \pi(s) + \sigma(s)$.

Recall that $\pi$ measures the individual's preference for commitment and $\sigma$ measures self-control. Hence, Corollary 2 decomposes the effect of an addictive consumption into two components: an increase in the preference for commitment and a reduction of self-control. Our next objective is to relate each of these two components of addiction to behavior.

We denote with $D(s, z)$ the current period drug demand when the agent faces the deterministic choice problem $z \in \bar{Z}$ and the state is $s$. Formally, $d \in D(s, z)$ if there exists an a non-drug consumption $b$ and a continuation problem $x$ such that $(b, d, x)$ is an optimal choice from $z$. We write $D(s, x) \geq D(s', y)$ if $d \in D(s, x), d' \in D(s', y)$ implies $d \geq d'$. Proposition 3 shows that lower self-control (higher $\sigma$) leads to higher drug demand in every decision problem.

**Proposition 3:** *If $\sigma(s) \geq \sigma(s')$ then $D(s, z) \geq D(s', z)$ for all $z \in \bar{Z}$.*

Recall that the agent's choice is governed by $u_0(b, d) + \sigma(s)v_0(d) + \delta W(d, x)$. Since $v_0$ is increasing in $d$, it follows that drug demand increases as $\sigma$ increases. Psychologists use the term *reinforcement* to describe the fact that an increase in current drug consumption leads to an increase in future drug consumption. If $\epsilon > 0$ and $\sigma(d + \epsilon) > \sigma(d)$ then the $\epsilon$ increase is reinforcing. In particular, an addictive increase in drug consumption is always reinforcing.

Reinforcement is necessary but not sufficient for an increase to be addictive. Addiction also means an increase in the preference for commitment, as measured by $\pi$. To translate

the change of $\pi$ into behavior we consider situations in which the individual chooses between alternatives that all yield the same current consumption. This choice is unaffected by temptation since all options that yield the same current consumption are equally tempting.

Note that $(c, (c', x))$ denotes an alternative that yields $c$ in the current period, $c'$ next period and the choice problem $x$ two periods hence. Therefore, for $z \in \bar{Z}$ the choice problem $\{c\} \times z = \{(c, (c', x)) \mid (c', x) \in z\}$ has the property that current consumption is fixed at $c$ and the choice from $z$ is made in the current period. This choice leads to a drug consumption in the next period. We call this drug consumption *advance drug demand* from $z$. It is easy to see from equation (2) that advanced demand does not depend on the current state $s$ and depends on current consumption $c = (b, d)$ only through $d$. Hence, we use $D^A(d, z)$ to denote advanced drug demand. That is; $d^* \in D^A(d, z)$ if and only if there exists $b, b^*, x$ such that $(b, d, (b^*, d^*, x))$ is an optimal choice from $z$. We write $D(d, x) \geq D(d', y)$ if $d^* \in D(s, x), d^{**} \in D(s', y)$ implies $d^* \geq d^{**}$.

Proposition 4 shows that if current (drug) consumption increases $\pi$ then the advance drug demand decreases. In other words, an increase in preference for commitment leads to a decrease in the advance demand for the drug.

**Proposition 4:** If $\pi(d) \geq \pi(d')$ then $D^A(d', z) \geq D^A(d, z)$.

The commitment utility $u_0(b, d) - \pi(s)v_0(d) + \delta W(d, x)$ governs the agent's advance demand. Since $v_0$ is decreasing in $d$ an increase in $\pi$ implies a decrease in advance demand.

Together Propositions 3 and 4 show that an addictive increase in drug consumption leads to an *increase in drug demand* and a *decrease in advance drug demand*. Advance drug demand captures the agent's behavior when he can commit to a choice prior to the consumption period. For choice problems that do not offer commitment, addiction leads to a widening of the gap between what the agent chooses and what he would have chosen had commitment been available. To put it differently, addiction increases both the agents demand for drugs and his need for commitment mechanisms for reducing consumption.

As an illustration consider an individual who must decide whether to attend a party. The party offers the opportunity to consume an addictive drug and a chance to meet friends. If the individual stays home he is committed to a drug-free evening and does not

meet friends. Suppose that the quality of the party is determined by the number of friends attending.

In this example, we can compare the decisions of the individual when he is addicted (high past drug consumption) and when he is not addicted (low past drug consumption). Consider a party that the individual would attend independent of prior drug consumption. When addicted, the agent has lower self-control and higher preference for commitment. A lower self-control implies that he consumes more of the drug at the party. By contrast, a higher preference for commitment has no effect on drug consumption. Addiction also has an effect on the decision to attend the party. Both the loss of self-control and the increase in preference for commitment make the party less attractive. Hence, if the non-addicted agent is just indifferent between attending the party and staying home, then the addicted agent will stay home. The increase in preference for commitment implies that every level of drug consumption is less desirable for the commitment utility. The loss in self-control implies a greater cost of self-control for any given level of drug consumption. Both make commitment to a drug-free evening more desirable.

## 4.   Behavioral and Policy Implications of Addiction

In this section, we analyze interventions that may impact the behavior and welfare of individuals struggling with addiction. Such policies may be available to the addict in the form of voluntary rehabilitation programs or may be imposed on him by government intervention.

A typical anti-drug intervention can affect addicts along two dimensions. First, it may reduce the maximum feasible consumption of the drug. We call such policies *prohibitive*. For example, banning a drug use may reduce the maximum feasible consumption to zero. The use of a deterrent drug is an example of a voluntary prohibitive policy. Second, a policy may increase the opportunity cost of drug consumption. Fines or taxes levied on addictive goods fall into this category. When such measures do not significantly change the maximum feasible drug consumption we call them *price* policies.

We assume that $C = [0,1]^2$ and consider a simple SSC preference $(u_0, v_0, \pi, \sigma, \delta)$, where $u_0, v_0, \pi$ and $\sigma$ are twice continuously differentiable. We also assume that $u_0$ is

strictly increasing in its first argument. Hence, the first good is indeed a good. We say that the drug is bad if $u_0$ is strictly decreasing in its second argument. Note that if the drug is bad then the preference is regular.

Suppose that the individual is endowed with one unit of good 1 in each period and faces the following stationary consumption problem. Each period the agent chooses consumption $(b, d)$ in the set $B$ where

$$B = \{(b, d) \in [0, 1]^2 \mid b + pd \le 1, d \le 1\}$$

We assume that the price of the drug, $p$, is less than 1. The individual cannot borrow or lend and can at most consume 1 unit of each good in every period.

A drug policy is a pair $(\tau, q)$ where $\tau \ge 0$ is a per unit tax on the drug and $q \in [0, 1]$ is the maximum feasible drug consumption. Let

$$B(\tau, q) = \{(b, d) \in [0, 1]^2 \mid b + (p + \tau)d \le 1, d \le q\}$$

denote the individual's opportunity set under the policy $(\tau, q)$ and $x(\tau, q)$ denote the corresponding stationary choice problem:

$$x(\tau, q) := \{(c, x(\tau, q)) \mid c \in B(\tau, q)\}$$

Note that any policy $(0, q)$ with $q < 1$ is a prohibitive policy since it reduces the maximum feasible drug consumption but does not affect the opportunity cost of drugs. A price policy is a policy $(\tau, 1)$ with $p + \tau \le 1$. In this case, the maximum feasible drug consumption remains at 1 in every period but the opportunity cost of the drug is increased to $p + \tau$. If the tax is high enough, in particular, if $p + \tau > 1$ then the policy $(\tau, 1)$ also has a prohibitive effect since it decreases the maximal drug consumption to $\frac{1}{p+\tau}$.

Propositions 5 and 6 examine the welfare effects of prohibitive and price policies. Proposition 5 demonstrates that a prohibitive policy on a bad, addictive drug increases the agent's welfare. Proposition 6 shows that a price policy decreases the agent's welfare.

**Proposition 5:** *If the drug is bad and addictive then $W(s, x(0, q)) > W(s, x(0, q'))$ for $q' > q$.*

**Proof:** Let $d'^0 = \hat{d}^0 = s$ and $\{(b^t, d^t)_{t \geq 1}\}$ denote the optimal consumption plan for the choice problem $x(0, q)$ at state $s$. Similarly, let $\{(b'^t, d'^t)_{t \geq 1}\}$ denote the optimal consumption plan for the choice problem $x(0, q')$ at state $s$. Define $\hat{b}^t = 1 - pd'^t$ and $\hat{d}^t = \min\{d'^t, q\}$ for all $t \geq 1$. Clearly, $\hat{d}^t \leq d'^t$ for all $t \geq 1$. Since the drug is bad and addictive, this implies $u_0(\hat{b}^t, \hat{d}^t) - \pi(\hat{d}^{t-1})v_0(\hat{d}^t) \geq u_0(b'^t, d'^t) - \pi(d'^{t-1})v_0(d'^t)$ and $[\pi(\hat{d}^{t-1}) + \sigma(\hat{d}^{t-1})][v_0(\hat{d}^t) - v_0(q)] \geq [\pi(d'^{t-1}) + \sigma(d'^{t-1})][v_0(d'^t) - v_0(q')]$. Moreover, at least one of the two preceding inequalities is strict. To see this, note that if $v_0(\hat{d}^t) - v_0(q) < 0$ or $v_0(d'^t) - v_0(q') < 0$ then the second inequality is strict. If $v_0(\hat{d}^t) - v_0(q) = v_0(d'^t) - v_0(q') = 0$ then $\hat{b}^t > b'^t$ so the first inequality is strict. Hence,

$$
\begin{aligned}
W(s, x(0, q)) &\geq \sum_{t=0}^{\infty} \delta^t [u_0(\hat{b}^t, \hat{d}^t) + \sigma(\hat{d}^{t-1})v_0(\hat{d}^t) - (\pi(\hat{d}^{t-1}) + \sigma(\hat{d}^{t-1}))v_0(q)] \\
&> \sum_{t=0}^{\infty} \delta^t [(u_0(b'^t, d'^t) + \sigma(d'^{t-1})v_0(d'^t) - (\pi(d'^{t-1}) + \sigma(d'^{t-1}))v_0(q')] \\
&= W(s, x(0, q'))
\end{aligned}
$$

$\square$

A prohibitive policy has two effects; it reduces self-control costs and it may render the previous level of drug consumption infeasible. The reduction in self-control costs always increases welfare. If the drug is bad the commitment utility maximizing level of drug consumption is zero. Hence, the reduction in consumption increases utility in the current period. Moreover, if the good is addictive, this reduction in consumption leads to a decrease in level of addiction which reduces future self-control costs. Thus, a purely prohibitive policy on a bad, addictive drug always increases welfare.

It is easy to see how Proposition 5 may fail for a drug that is not bad. In that case, the argument is similar to the standard economic argument for why a quota or a ban may reduce welfare. To see how Proposition 5 fails when the drug is not addictive, consider an agent who is in state $s = .5$ in period 1. Suppose that abstaining $(d = 0)$ or binging $(d = 1)$ for one period will cause all temptation to go away in the next period but consuming intermediate levels will cause temptation to persist. Moreover, assume that the cost of self-control in the current state is very high. In that case, it may be optimal for the agent to binge in the current period and abstain thereafter. A policy that reduces

the maximally feasible level of drug consumption from 1 to $q = .5$ may reduce the agents welfare by forcing him to either incur the (reduced but still) high cost of self-control in the current period or remain addicted.

Proposition 5 showed that a reduction in temptation increases welfare in our model. Proposition 6 shows that a policy that does not reduce temptation cannot increase the agent's welfare. Recall that a price policy is an increase in the opportunity cost of drug consumption that does not reduce the maximally feasible drug consumption. When the tax increases from $\tau$ to $\tau'$ and $p + \tau' \leq 1$ then even after the tax increase, the maximally feasible drug consumption in the current period is unchanged at 1.

**Proposition 6:** *If the drug is bad and $p + \tau' \leq 1$ then $W(s, x(\tau, 1)) \geq W(s, x(\tau', 1))$ for $\tau' > \tau$.*

**Proof:** Let $d^0 = s$ and $\{(b^t, d^t)_{t \geq 1}\}$ be the optimal consumption plan for the problem $x(\tau', 1)$. Since $\{(b^t, d^t)_{t \geq 1}\}$ is a feasible choice from $x(\tau, 1)$ we have

$$W(s, x(\tau, 1)) \geq \sum_{t=0}^{\infty} \delta^t \left( u_0(b, d) + \sigma(d^{t-1}) v_0(d^t) - (\pi(d^{t-1}) + \sigma(d^{t-1})) v_0(1) \right)$$
$$= W(s, x(\tau', 1))$$

$\square$

Since a price policy does not affect the maximally feasible drug consumption it does not reduce self-control costs. Therefore, a price policy cannot improve the agent's welfare. The proof of Proposition 6 gives a simple revealed preference argument: since temptation is unaffected by the tax increase and the set of alternatives is smaller the agent cannot be better off. The key hypothesis in Proposition 6 is that the maximally feasible drug consumption is smaller than the maximal amount of drugs the individual can afford in the current period. This hypothesis is likely to be satisfied if the drugs under consideration are inexpensive or cannot be consumed in large doses. For example, for most smokers the maximally feasible cigarette consumption is a small fraction of the individual's budget.

Next, we analyze the impact of policies on the demand for drugs. Current period drug demand in state $s$ under the policy $(\tau, q)$ is denoted $D(s, x(\tau, q))$. Consider a purely prohibitive policy $(0, q)$. If the prohibitive policy is binding, that is, if $D(s, x(0, q)) = q$

then a reduction in the maximum allowed drug consumption $q$ will obviously lead to a reduction in drug demand. Proposition 7 shows that if the policy is not binding then a reduction in $q$ will lead to an *increase* in drug demand.

**Proposition 7:** *If the drug is bad and addictive and $D(s, x(0, q))$ is a differentiable function of $q$ then $\partial D(s, x(0, q)) / \partial q < 0$.*

**Proof:** If $(b^*, d^*)$ is the optimal choice in the current period and $(b^{**}, d^{**})$ is the optimal choice in the next period, then

$$W(s, x(0, q)) = u_0(1 - pd^*, d^*) + \sigma(s)v_0(d^*) + \delta W(d^*, x(0, q)) - (\pi(s) + \sigma(s))v(q)$$

Since the optimal consumption is interior, the necessary first order condition is

$$
\begin{aligned}
0 = &- p\frac{\partial u_0(1 - pd^*, d^*)}{\partial b} + \frac{\partial u_0(1 - pd^*, d^*)}{\partial d} + \sigma(s)\frac{\partial v_0(d^*)}{\partial d} + \delta\frac{\partial W(d^*, x(0, q))}{\partial d} \\
= &- p\frac{\partial u_0(1 - pd^*, d^*)}{\partial b} + \frac{\partial u_0(1 - pd^*, d^*)}{\partial d} + \sigma(s)\frac{\partial v_0(d^*)}{\partial d} \\
&+ \delta\left[v_0(d^{**})\frac{\partial \sigma(d^*)}{\partial d} - v_0(q)\frac{\partial(\pi(d^*) + \sigma(d^*))}{\partial d}\right] \\
\equiv &A(d^*)
\end{aligned}
$$

Taking the total derivative we get

$$\mathrm{d}d\frac{\partial A(d^*)}{\partial d} - \mathrm{d}q\frac{\partial(\pi(d^*) + \sigma(d^*))}{\partial d}\frac{\partial v_0(q)}{\partial q} = 0$$

From the second order condition, we infer that $\frac{\partial A(d^*)}{\partial d} \leq 0$. Since $\frac{\mathrm{d}d}{\mathrm{d}q}$ is well-defined it follows that $\frac{\partial A(d^*)}{\partial d} < 0$. By corollary 2, $\pi + \sigma$ is strictly increasing. Since $v_0$ is also strictly increasing the desired result follows. $\qquad\square$.

A prohibitive policy effectively reduces the cost of drug consumption by reducing the future cost of self-control associated with current drug consumption. For this reason, drug demand increases as the prohibitive policy becomes more stringent. In contrast, the analysis of the demand effect of a pure price policy $((\tau, 1), p + \tau \leq 1)$ is no different than the analysis of demand effects of price changes in a standard consumer problem: if current drug consumption is a normal good then drug demand is decreasing in $\tau$.[10]

---

[10] In general, just as in standard demand theory, the response to an increase in $\tau$ is ambiguous.

In the following simple example, we analyze further the effect of a change in $\tau$. Consider an agent with constant self-control $\sigma(d) = \alpha$ and increasing preference for commitment, $\pi(d) = \alpha d$. Let

$$u_0(b, d) = b - d$$

$$v_0(b, d) = \log d$$

The individual faces the choice problem $x(\tau, q)$ and therefore solves the following maximization problem:

$$W(s, x(\tau, q)) := \max_{\{d^t\}} \sum_{t=1}^{\infty} \delta^{t-1} \left[ 1 - (p + \tau)d^t - d^t + \alpha \log d^t - \alpha(1 + d^{t-1}) \max_{B(\tau, q)} \log d \right]$$

First, consider the case where $q < \frac{1}{p+\tau}$. Then, the maximally feasible drug consumption is $q$ and a change in $\tau$ constitutes a pure price policy. Drug demand is

$$D(s, x(\tau, q)) = \min \left\{ \frac{a}{1 + p + \tau + \delta\alpha \log q}, q \right\}$$

When drug consumption is not constrained by $q$, the only effect of an increase in $\tau$ is to increase the opportunity cost of consuming drugs. In our example current drug consumption is a normal good and therefore drug demand is decreasing in $\tau$. We call this reduction in drug demand the *price effect* of the drug policy.

Second, consider the case where $q \geq \frac{1}{p+\tau}$. Then, the maximally feasible drug consumption is $\frac{1}{p+\tau}$. An increase in $\tau$ raises the opportunity cost of the drug but also reduces temptation (i.e., the maximally feasible drug consumption). Solving the maximization problem above yields,

$$D(s, x(\tau, q)) = \min \left\{ \frac{\alpha}{1 + p + \tau + \delta a \log(1/(p + \tau))}, 1 \right\}$$

Hence, the demand effect of an increase in $\tau$ is ambiguous. To see this, assume that drug demand is less than $q$. Inspecting the demand function above, we find that

$$\frac{\partial}{\partial \tau} D(s, x(\tau, q)) = (\delta\alpha - (p + \tau))D(s, x(\tau, q))$$

Therefore, if $\delta\alpha > p + \tau$, the increase in $\tau$ leads to an *increase* in drug demand. The increase in $\tau$ reduces the self-control costs associated with addiction and results in higher

20

drug demand. This *self-control effect* effect dominates the price effect. Finally, if $\delta\alpha < p+\tau$ drug demand decreases. In this case, the price effect dominates the self-control effect.

As in the tax example above, most actual policies affect both the opportunity cost and the maximally feasible drug consumption. As our analysis shows, the prohibitive aspects of a policy benefit the individual as long as the drug is bad and addictive. However, when a prohibitive policy is not binding it also results in an increase in drug demand. On the other hand, an increase in opportunity cost of the drug harms individual and, in the case of a normal good, will lead to lower drug demand.

## 5. Representation Theorems

Preferences are defined over $Z$, the set of choice problems. Let $\Delta$ denote the set of probability measures on $C \times Z$. A choice problem $z \in Z$ can be identified with a compact subset of $\Delta$. More precisely, for any subset $X$ of a metric space, we let $\Delta(X)$ denote the set of all probability measures on the Borel $\sigma-$algebra of $X$ and $\mathcal{K}(X)$ denote the set of all nonempty compact subsets of $X$. Each $z \in Z$ can be identified with an element in $\mathcal{K}(\Delta(C \times Z))$ and conversely each element in $\mathcal{K}(\Delta(C \times Z))$ identifies a choice problem $z \in Z$. For formal definitions of $Z$ and the map that associates each element of $Z$ with its equivalent recursive description as an element of $\mathcal{K}(\Delta(C \times Z))$, we refer the reader to Gul and Pesendorfer (2000b). In what follows only the recursive definition is used and hence without risk of confusion we identify the sets $Z$ and $\mathcal{K}(\Delta(C \times Z))$. In Gul and Pesendorfer (2000b) we note that since $C$ is a compact metric space then $Z, \Delta(C \times Z)$ and $\mathcal{K}(\Delta(C \times Z))$ are compact metric spaces as well.

The individual's preferences are indexed by $s \in S$, the state in the initial period of the choice problem. The state $s$ represents the relevant consumption history prior to the initial period. We assume that there is a finite number $K$ such that consumption in only the last $K$ periods influences the agents preferences and therefore $S := C^K$. Without loss of generality, we assume that $K$ is the minimal length of the individual's consumption history that allows us to describe $\succeq$.[11] We refer to the indexed family of preferences $\succeq := \{\succeq_s\}_{s \in S}$

---

[11] That is, there is a pair of states, $(s = (c^1, \cdots, c^K), \hat{s} = (\hat{c}^1, \cdots, \hat{c}^K))$ that differ only in their first component $(c^1 \neq \hat{c}^1, c^t = \hat{c}^t, t \geq 2)$ and lead to different preferences $(\succeq_s \neq \succeq_{\hat{s}})$.

simply as *the agent* or the preference $\succeq$. Recall that for any state $s = (c^1, \ldots, c^K)$ $sc$ denotes the state $(c^2, \ldots, c^K, c)$. We impose the following axioms on $\succeq_s$ for every $s \in S$.

**Axiom 1:**  *(Preference Relation)* $\succeq_s$ *is a complete and transitive binary relation.*

**Axiom 2:**  *(Strong Continuity) The sets* $\{x \mid x \succeq_s z\}$ *and* $\{x \mid z \succeq_s x\}$ *are closed in* $Z$.

**Axiom 3:**  *(Independence)* $\{\mu\} \succ_s \{\nu\}$ *implies* $\{\alpha\mu + (1-\alpha)\eta\} \succ_s \{\alpha\nu + (1-\alpha)\eta\}$ $\forall \alpha \in (0, 1)$.

Axioms 1-3 are standard. In axiom 4 we deviate from standard choice theory and allow the possibility that the adding options to a choice problem makes the consumer strictly worse off. For a detailed discussion of Axiom 4, we refer the reader to our earlier paper (Gul and Pesendorfer 2000a).

**Axiom 4:**  *(Set Betweenness)* $x \succeq_s y$ *implies* $x \succeq_s x \cup y \succeq_s y$.

Next, we make a separability assumption. For $z \in Z$ let $cz \in Z$ denote the choice problem $\{(c, z)\}$, that is, the degenerate choice problem that yields $c$ in the current period and the continuation problem $z$. Thus $c_1 c_2 \ldots c_K z$ is a degenerate choice problem that yields the consumption $(c_1, ..., c_K)$ in the first $K$ periods and the continuation problem $z$ in period $K+1$. For $s = (c_1, \ldots, c_K)$ we write $sz$ instead of $c_1 c_2 \ldots c_K z$. Axiom 5 considers choice problems of the form $\{(c, sz)\}$ and requires that preferences are not affected by the correlation between current consumption $c$ and the $K + 1$ period continuation problem $z$.

**Axiom 5:**  *(Separability)* $\{\frac{1}{2}(c, sz) + \frac{1}{2}(c', sz')\} \sim_{s'} \{\frac{1}{2}(c, sz') + \frac{1}{2}(c', sz)\}$.

Axiom 6 requires preferences to be stationary. Consider the degenerate lotteries, $(c, x)$ and $(c, y)$, each leading to the same period 1 consumption $c$. Stationarity requires that $\{(c, x)\}$ is preferred to $\{(c, y)\}$ in state $s$ if and only if the continuation problem $x$ is preferred to the continuation problem $y$ in state $sc$.

**Axiom 6:**  *(Stationarity)* $\{(c, x)\} \succeq_s \{(c, y)\}$ *iff* $x \succeq_{sc} y$.

Note that Axiom 6 implies that the conditional preferences at time $K + 1$ after consuming $s$ in the first $K$ periods is the same as the initial preference $\succeq_s$. Together, Axioms 5 and 6 restrict the manner in which past consumption influences future preferences. Axiom 5 ensures that correlation between consumption prior to period $t - K$ and the choice problem in period $t$ does not affect preferences whereas Axiom 6 ensures that the realization of consumption prior to period $t - K$ does not affect preferences in period $t$.

Axiom 7 requires individuals to be indifferent as to the timing of resolution of uncertainty. In a standard, expected utility environment this indifference is implicit in the assumption that the domain of preference is the set of lotteries over consumption paths. Our domain of preferences are choice problems and in this richer structure a separate assumption is required to rule out agents that are not indifferent to the timing of resolution of uncertainty as described by Kreps and Porteus (1978).

Consider the lotteries $\mu = \alpha(c, x) + (1-\alpha)(c, y)$ and $\nu = (c, \alpha x + (1-\alpha)y)$. The lottery $\mu$ returns the consumption $c$ together with the continuation problem $x$ with probability $\alpha$ and the consumption $c$ with the continuation problem $y$ with probability $1 - \alpha$. By contrast, $\nu$ returns $c$ together with the continuation problem $\alpha x + (1-\alpha)y$ with probability 1. Hence, $\mu$ resolves the uncertainty about $x$ and $y$ in the current period whereas $\nu$ resolves this uncertainty in the future. If $\{\mu\} \sim_s \{\nu\}$ then the agent is indifferent as to the timing of the resolution of uncertainty.

**Axiom 7:**   (Indifference to Timing) $\{\alpha(c, x) + (1 - \alpha)(c, y)\} \sim_s \{(c, \alpha x + (1 - \alpha)y)\}$.

**Definition:**   $\succeq_s$ is regular if there exists $x, x', y, y' \in Z$ such that $x' \subset x, y' \subset y, x \succ_s x'$ and $y' \succ_s y$. $\succeq$ is regular if each $\succeq_s$ is regular.[12]

Theorem 1 below establishes that all regular preferences that satisfy Axioms $1 - 7$ can be represented as a discounted sum of state-dependent utilities minus state-dependent self-control costs. We say that the function $W : S \times Z \to \mathbb{R}$ represents $\succeq$ when $x \succeq_s y$ iff $W(s, x) \geq W(s, y)$ for all $s$. For any $\mu \in \Delta$, let $\mu^1$ denote the marginal of $\mu$ on $C$. Axioms 1-7 yield the following representation.

**Theorem 1:**   If $\succeq$ is regular and satisfies Axioms $1 - 7$, then there exists $\delta \in (0, 1)$, continuous functions $u : S \times \Delta(C) \to \mathbb{R}$, $V : S \times \Delta \to \mathbb{R}$, $W : S \times Z \to \mathbb{R}$ with $u(s, \cdot), V(s, \cdot), W(s, \cdot)$ linear for all $s$, such that

$$W(s, z) = \max_{\mu \in z}\{u(s, \mu^1) + \delta \int W(sc, z)d\mu(c, z) + V(s, \mu)\} - \max_{\nu \in z} V(s, \nu)$$

---

[12]   In section 2, we presented the definition of a regular SSC preference. It can be shown that the current general definition and the one offered in section 2 are equivalent for SSC preferences.

for all $s \in S, \nu \in \Delta$ and $W$ represents $\succeq$. For any $\delta \in (0,1)$, continuous $u, V$ such that $u(s, \cdot)$ and $V(s, \cdot)$ are linear for all $s \in S$, there exists a unique function $W$ that satisfies the equation above and the $\succeq$ represented by this $W$ satisfies Axioms $1 - 7$.

The two main steps of the proof of Theorem 1 entail showing that a preference relation (over choice problems) that satisfies continuity, independence, set betweenness, stationarity and indifference to timing of resolution of uncertainty has a representation of the form

$$W(s, z) = \max_{\mu \in z}\{U(s, \mu) + V(s, \mu)\} - \max_{\nu} V(s, \nu)$$

and then using stationarity and separability to show that $U$ is of the form $U = u + \delta W$. In Gul and Pesendorfer (2000b) we offer a related proof under stronger stationarity and separability axioms, yielding a representation of state-independent preferences

Next, we provide assumptions under which preferences can be represented by a utility function $W$ that satisfies equation (2) (i.e., characterize simple SSC preferences).

Assumption $I$ below requires that the agent is tempted only by the prospect of immediate consumption. The assumption considers two situations in which the agent uses self-control and chooses the same alternative $\mu$. If the tempting alternative $\nu$ in one situation, yields the same consumption in the current period as the tempting alternative in the second situation $\eta$ they are equally tempting. That is, the agent is indifferent between the two situations.

**Assumption I:** $U(s, \mu) + V(s, \mu) > U(s, \nu) + V(s, \nu)$, $U(s, \mu) + V(s, \mu) > U(s, \eta) + V(s, \eta)$, $V(s, \nu) > V(s, \mu)$, $V(s, \eta) > V(s, \mu)$ and $\nu^1 = \eta^1$ implies $\{\mu, \nu\} \sim_s \{\mu, \eta\}$.

Assumption $N$ below ensures that goods other than $d$ are *neutral* i.e., cause no temptation and have no dynamic effects. That is, only good $d$ is tempting and only past consumption of $d$ affects future rankings of choice problems.

**Assumption N:** Let $c = (b, d)$, $c' = (b', d')$ and $U(s, (c, x)) > U(s, (c', x'))$. If $d = d'$ then $\succeq_{sc} = \succeq_{sc'}$. Moreover, $\succeq_s$ has **P** at $\{(c, x), (c', x')\}$ iff $d' > d$.

Recall that $\{(c, z)\} \succ_s \{(c, z), (c', z')\}$ means that $(c', z')$ has greater commitment utility $U$ and less temptation utility $V$ than $(c', z')$. Hence, $\succeq_s$ has a preference for commitment at $\{(c, z), (c', z')\}$. Therefore, the first statement ensures that there is no **P** (i.e.,

no temptation) so long as the options differ only with respect to current consumption of non-drugs. The second statement means that future preferences are the same so long as the current state and current consumption of drugs are the same. Finally, the third statement means that higher current drug consumption is always tempting. That is, if the alternative with higher current drug consumption yields lower commitment utility then it creates temptation, that is decreases the utility of the current situation.

Our next assumption requires that the agent can rank states according to $\mathbf{P}$, his preference preference for commitment. Thus, for every $s, s'$ either $s\mathbf{P}s'$ or $s'\mathbf{P}s$.

**Assumption P:**  $\mathbf{P}$ *is complete.*

**Theorem 2:**  *Let $\succeq$ be a regular SSC preference satisfying I, N and T. Then, (i) $S = [0,1]$, and (ii) there are continuous functions $v_0, \pi, \sigma : [0,1] \to I\!R$, $u_0 : C \to I\!R$, and $\delta \in (0,1)$, such that for $z \in \bar{Z}$*

$$W(s, z) = \max_{(b,d,x)\in z} \{u_0(b,d) + \sigma(s)v_0(d) + \delta W(d,x)\} - \max_{(b',d',y)\in z} (\pi(s) + \sigma(s))v_0(d')$$

*and $W$ represents $\succeq$. (iii) $v_0$ is strictly increasing; $\pi \geq 0$, $\pi + \sigma > 0$; and $s$ is the previous period's drug consumption.*

**Proof:**  See Appendix.

To get some intuition about the proof of Theorem 2, start with the representation of SSC preferences provided in Theorem 1. Then, assumption $I$ ensures that $V(s, \cdot)$ depends only on current consumption $c$. Assumption $N$ guarantees that $V(s, \cdot)$ depends only on current drug consumption $d$ and is strictly increasing in $d$.

Since $U$ measures commitment utility while $V$ measures temptation, if $U(s', \cdot), V(s', \cdot)$ were both convex combinations of $U(s, \cdot), V(s, \cdot)$, the gap between commitment utility and temptation utility would be less in state $s'$ than in $s$. In this case, the agent would have greater preference for commitment at $s$ than at $s'$. In Gul and Pesendorfer (2000a) we show that $U, V$ moving closer in this way is necessary and sufficient for $s\mathbf{P}s'$. Since $V(s, \cdot)$ does not depend on next period's choice problem and $U(s, \cdot)$ does, it follows that the convex combination of $U(s, \cdot), V(s, \cdot)$ that defines $V(s', \cdot)$ puts zero weight on $U(s, \cdot)$. Hence, $V(s', \cdot)$ is of the form $V(s', \cdot) = a(s)v_0(d)$.

To conclude the argument, we need to verify that only last period's consumption can influence current preferences. This result is an implication of Assumption P. Suppose, contrary to Theorem 2, that the consumption in the past two periods affects current preferences. The following example illustrates why this would violate Assumption P.

As in Theorem 2, the example considers an SSC preference that can be represented by $(u_0, v_0, \pi, \sigma, \delta)$. However, $\pi$ depends on the consumption in the past two periods. There is one good, $d$ and $s = (d^{-1}, d^0) \in [0,1]^2$. Further, let $u_0(d) = 2d$, $v_0(d) = d$, $\pi(s) = 3d^{-1} + d^0, \sigma(s) = 1$ and $\delta = 2/3$. We compare the preferences at states $s = (1, 0)$ and $s' = (0, 1)$. Let $z$ be the choice problem in which the agent is committed to zero consumption, that is $z = \{0, z\}$. Consider the choice problem $x = \{(0, z), (1, z)\}$. In state $s'$ the commitment utility is increasing in current drug consumption $(2 - \pi(s') > 0)$ and hence the agent has no $\mathbf{P}$ at $x$. By contrast, in state $s$ the commitment utility is decreasing in current drug consumption $(0 > 2 - \pi(s))$ and hence the agent has $\mathbf{P}$ at $x$. Now consider the choice problem $y = \{(1/2, x), (0, z)\}$. In state $s$ the commitment utility of $(1/2, x)$ is $-1/2 + \delta(3/2) = 1/2$ and the commitment utility of $(0, z)$ is zero. Therefore, in state $s$ the agent has no $\mathbf{P}$ at $y$. In state $s'$ the commitment utility of $(1/2, x)$ is $1/2 + \delta(-1) = -1/6$ and the commitment utility of $(0, z)$ is zero. Therefore, in state $s'$ the agent has $\mathbf{P}$ at $y$. We therefore conclude that $\succeq_s$ and $\succeq_{s'}$ cannot be ranked according to $\mathbf{P}$. A similar violation of assumption P can be obtained where $V$ to depend on the last two periods consumption.

## 6.  Conclusion

Most studies on drug abuse emphasize that addiction should be considered a disease.[13] In our approach drug abuse is identified with the discrepency between what the agent would want to commit to, as reflected by maximizing $U$, and what he ends-up consuming by maximizing $U + V$. We provide straightforward choice experiments for measuring this discrepency. Our approach is silent on the question of whether addiction is a disease or a part of the "normal" variation of preferences across individuals.

While our approach is compatible with the disease concept of addiction, there are important differences between the two. Consider the following example: the opiate antagonist naltrexone blocks the opioid receptors in the brain and hence the euphoric effects

---

[13]  To emphasize the organic basis of the condition the term "disease of the brain" is often used.

of these drugs for up to 3 days after the last dose. Naltrexone is used in the treatment of heroin and morphine. However, with the exception of highly motivated addicts such as parolees, probationers and health care professionals, most addicts receiving naltrexone tend to stop taking their medicine and relapse. Addicts often report that they stop taking naltrexone because it prevents "getting high". Doctors call this as a "compliance problem" with naltrexone. For them, this is simply a limitation on the usefulness naltrexone, the same way that toxicity might be a limitation on the usefulness of some other medication.

In our model, there can be two reasons for an addict to discontinue naltrexone and resume heroin consumption: either 3 days is not the right time horizon for commitment or the addict does not wish to commit. The former would suggest a need for longer acting drugs while the latter would mean that there is neither a need nor any room for treatment of this addict. In fact, by our definition, an individual who is unwilling to commit to reducing his drug consumption, for any length of time, at any future date is not an addict. Hence, where the disease model of addiction finds a compliance problem our model suggests that there may be no problem at all.

The fact that naltrexone continues to be used by the most motivated addicts, those who are more likely to abstain even without commitment, suggests a reduction of the cost of self-control as a possible motive taking naltrexone.

Economists interpret behavior as a reflection of the agents' stable interests and desires. In standard economic analysis there is no room for the notion of a behavioral problem, except to the extent that the behavior is a problem for someone else. Consequently, there is no role for therapy aimed at controlling problem behavior. In contrast, psychologists often view behavior to be independent of and even an impediment to the agent's welfare. Our model of temptation and self-control provides a potential bridge between these two approaches. Like standard models in economics, we take as given agents' interests and desires (i.e. utility functions) and accept the hypothesis that behavior is motivated by these interests and desires (i.e. utility maximization). But, we extend the domain of utility functions to include temptation. Without the aid of some outside agency, it is difficult and often very costly for the individual to commit, that is; reduce temptation. Hence, our model leaves room for welfare enhancing treatments and policy. In our interpretation, the

role of treatment and policy is to develop commitment devices and opportunities for the agent.

Our model provides a framework for the analysis of both the purposeful actions (e.g. decisions made in the stock market) studied by most economists as well as the compulsive and detrimental behavior (e.g. addiction) studied by many psychologists and health care professionals. We have analyzed the interaction of these two types of behavior and evaluated policy alternatives. Our focus was on psychoactive drugs but the model presented in this paper can also be applied to other types of compulsive behavior such as over-eating and other forms of dependency.

# 7. Appendix

## 7.1 Proof of Proposition 1

Lemma 12 in Gul and Pesendorfer (2000a) implies that whenever $(a)$ $U(s, \cdot), U(s', \cdot)$ are not constant, $(b)$ $V(s, \cdot)$ is not an affine transformation of $U(s, \cdot)$ and $(c)$ $V(s', \cdot)$ is not an affine transformation of $U(s', \cdot)$; $s\mathbf{P}s'$ if and only if there is a non-negative full rank matrix $\Theta$ and a $\lambda \in \mathbb{R}^2$ such that

$$\begin{pmatrix} U(s', \cdot) \\ V(s', \cdot) \end{pmatrix} = \Theta \begin{pmatrix} U(s, \cdot) \\ V(s, \cdot) \end{pmatrix} + \lambda \tag{3}$$

Similar arguments establish that if $(a) - (c)$ hold then $s'\mathbf{S}s$ iff there is a positive full rank matrix $\hat{\Theta}$ and a $\hat{\lambda} \in \mathbb{R}^2$ such that

$$\begin{pmatrix} U(s, \cdot) + V(s, \cdot) \\ V(s, \cdot) \end{pmatrix} = \hat{\Theta} \begin{pmatrix} U(s', \cdot) + V(s', \cdot) \\ V(s', \cdot) \end{pmatrix} + \hat{\lambda} \tag{4}$$

and $s\mathbf{C}s'$ iff there is a positive full rank matrix $\tilde{\Theta}$ and a $\tilde{\lambda} \in \mathbb{R}^2$ such that

$$\begin{pmatrix} U(s', \cdot) \\ U(s', \cdot) + V(s', \cdot) \end{pmatrix} = \tilde{\Theta} \begin{pmatrix} U(s, \cdot) \\ U(s, \cdot) + V(s, \cdot) \end{pmatrix} + \tilde{\lambda} \tag{5}$$

Regularity of $\succeq$ implies that conditions $(a) - (c)$ are satisfied. Therefore, equation (3) implies that equation (5) holds for some full rank $\tilde{\Theta}$ and $\tilde{\lambda}$. Routine calculations establish that $\tilde{\Theta}$ is non-negative. $\qquad\square$

## 7.2 Proof of Proposition 2

As in the proof of Proposition 2 above, we note that $s\mathbf{P}s$ if and only if there is a positive full rank matrix $\Theta$ and a $\lambda \in \mathbb{R}^2$ such that

$$\begin{pmatrix} U(s', \cdot) \\ v(s', \cdot) \end{pmatrix} = \Theta \begin{pmatrix} U(s, \cdot) \\ v(s, \cdot) \end{pmatrix} + \lambda$$

Note that

$$U(s', \cdot) = U(s, \cdot) + \frac{\pi(s) - \pi(s')}{\pi(s) + \sigma(s)} v(s, \cdot)$$

and

$$v(s', \cdot) = \frac{\pi(s') + \sigma(s')}{\pi(s) + \sigma(s)} v(s, \cdot)$$

29

Since $\pi + \sigma > 0$ it follows that $s\mathbf{P}s'$ if and only if $\pi(s) \geq \pi(s')$.

Similarly, $s'\mathbf{S}s$ if and only if there is a positive full rank matrix $\Theta$ and a $\lambda \in \mathbb{R}^2$ such that

$$\begin{pmatrix} U(s',\cdot) + v(s',\cdot) \\ v(s',\cdot) \end{pmatrix} = \Theta \begin{pmatrix} U(s,\cdot) + v(s,\cdot) \\ v(s,\cdot) \end{pmatrix} + \lambda$$

Note that

$$U(s,\cdot) + v(s,\cdot) = U(s',\cdot) + v(s',\cdot) + \frac{\sigma(s) - \sigma(s)}{\pi(s') + \sigma(s')} v(s',\cdot)$$

and

$$v(s,\cdot) = \frac{\pi(s) + \sigma(s)}{\pi(s') + \sigma(s')} v(s',\cdot)$$

Since $\pi + \sigma > 0$ it follows that $s'\mathbf{S}s$ if and only if $\sigma(s) \geq \sigma(s')$.

Finally, $s\mathbf{C}s'$ if and only if

$$\begin{pmatrix} U(s',\cdot) \\ U(s',\cdot) + v(s',\cdot) \end{pmatrix} = \Theta \begin{pmatrix} U(s,\cdot) \\ U(s,\cdot) + v(s,\cdot) \end{pmatrix} + \lambda$$

Note that

$$U(s',\cdot) = \frac{\pi(s') + \sigma(s)}{\pi(s) + \sigma(s)} U(s,\cdot) + \frac{\pi(s) - \pi(s')}{\pi(s) + \sigma(s)} (U(s,\cdot) + v(s,\cdot))$$

and

$$U(s,\cdot) + v(s,\cdot) = \frac{\sigma(s) - \sigma(s')}{\pi(s) + \sigma(s)} (U(s,\cdot)) + \frac{\pi(s) + \sigma(s')}{\pi(s) + \sigma(s)} (U(s,\cdot) + v(s,\cdot))$$

Therefore, $s\mathbf{C}s'$ if and only if $\pi(s) \geq \pi(s')$ and $\sigma(s) \geq \sigma(s')$. $\qquad\square$

## 7.3  Proof of Proposition 3

Assume that $\sigma(s) \geq \sigma(s')$. Let $d'$ be the maximal element in $D(s', z)$ and let $(b', d', x') \in \mathcal{C}_{U+v}(s', z)$ be a corresponding choice. Then,

$$u_0(b', d') + \sigma(s')v_0(d') + \delta W(d', x') \geq u_0(b'', d'') + \sigma(s')v_0(d'') + \delta W(d'', x'')$$

for all $(b'', d'', x'') \in z$. Since $\sigma(s) \geq \sigma(s')$, for any $(b'', d'', x'') \in z$ such that $d'' < d'$

$$u_0(b', d') + \sigma(s)v_0(d') + \delta W(d', x') > u_0(b'', d'') + \sigma(s)v_0(d'') + \delta W(d'', x'')$$

Hence $D(s, z) \geq D(s', z)$. $\qquad\square$

## 7.4  Proof of Proposition 4

Let $\pi(s) \geq \pi(s')$. Recall that advanced demand at state $s$ refers to a choice made prior to the consumption $(b, s)$ at some state $s^0$. Let $d$ be the maximal element in $D^A(s, z)$ and let $((b^0, s), (b, d, x)) \in \mathcal{C}_{U+v}(s^0, z)$ be a corresponding choice. Then,

$$u_0(b^0, s) + \sigma(s^0)v_0(s) + \delta[u_0(b, d) - \pi(s)v_0(d)] + \delta^2 W(d, x) \geq$$
$$u_0(b^0, s) + \sigma(s^0)v_0(s) + \delta[u_0(b'', d'') - \pi(s)v_0(d'')] + \delta^2 W(d'', x'')$$

for all $((b, s), (b'', d'', x'')) \in z$. Then, for any $((b^0, s'), (b'', d'', x'')) \in z$ with $d'' < d$

$$u_0(b^0, s) + \sigma(s^0)v_0(s') + \delta[u_0(b, d) - \pi(s')v_0(d)] + \delta^2 W(d, x) \geq$$
$$u_0(b^0, s') + \sigma(s^0)v_0(s') + \delta[u_0(b'', d'') - \pi(s')v_0(d'')] + \delta^2 W(d'', x'')$$

Hence $D^A(s', z) \geq D^A(s, z)$. $\qquad\square$

## 8.  Proof of Theorems 1 and 2

### 8.1  Definitions

Let $Z^1 := \{\{\mu\} \mid \mu \in \Delta\}$. For $n > 1$, define $Z^n := \{\{\mu\} \in Z^1 \mid \mu^2(Z^{n-1}) = 1\}$. For $c \in C$, define $Z^1(c) := \{\{\mu\} \in Z^1 \mid \mu^1(c) = 1\}$. For $n > 1$ and $c^1, \ldots, c^n \in C$ define $Z^n(c^1, \ldots, c^n) := \{\{\mu\} \in Z^1(c^1) \mid \mu^2(Z^{n-1}(c^2, \ldots, c^n)) = 1\}$. Finally, we define $Z_1^{n+1}(c^2, \ldots, c^{n+1}) = \{\{\mu\} \in Z^{n+1} \mid \mu^2(Z^n(c^2, \ldots, c^{n+1})) = 1\}$.

There is an obvious homeomorphism between $Z^n(c_1, \ldots, c_n)$ and $Z$. This homeomorhism associates with each $\{\mu\} \in Z^n(c^1, \ldots, c^n)$, a particular $\{\nu\} \in Z$ by ignoring the degenerate distribution of consumptions that $\mu$ yields from period 1 to $n$. Similarly, there is a homeomorphism between $Z_1^{n+1}(c^2, \ldots, c^{n+1}))$ and $Z$. This homeomorhism associates with each $\{\mu\} \in Z_1^{n+1}$ a particular $\{\nu\} \in Z$ by ignoring the degenerate distribution of consumptions that $\mu$ yields from period 2 to $n + 1$. We use $\mu^{T+2}$ to denote the marginal distribution $\nu^2$ of the associated $\nu$.

Let $\alpha x + (1 - \alpha)y := \{\alpha\mu + (1 - \alpha)\nu \mid \mu \in x, \nu \in y\}$. We say that a function $f : Z \to \mathbb{R}$ is linear if $f(\alpha x + (1 - \alpha)y) = \alpha f(x) + (1 - \alpha)f(y)$ for all $x, y \in Z$.

31

## 8.2 Proof of Theorem 1

**Lemma 1:** *Assume $\succeq$ satisfies Axioms 5 and 6. If $W$ is linear in its second argument and represents $\succeq$ then $W(s,(\bar{c}^0\bar{c}^1\ldots\bar{c}^n s'\bar{z})) + W(s,(c^0c^1\ldots c^n s'z)) = W(s,(c^0c^1\ldots c^n s'\bar{z})) + W((\bar{c}^0\bar{c}^1\ldots\bar{c}^n s'z))$ for all $n$, $(\bar{c}^0\ldots\bar{c}^n)$, $(c^0\ldots c^n) \in C^{n+1}$, $s' \in C^K$, $z,\bar{z} \in Z$.*

**Proof:** By Axiom 5,

$$\frac{1}{2}(\bar{c}^0 c^1 \ldots c^n s'\bar{z}) + \frac{1}{2}(c^0 c^1 \ldots c^n s' z) \sim_s \frac{1}{2}(c^0 c^1 \ldots c^n s'\bar{z}) + \frac{1}{2}(\bar{c}^0 c^1 \ldots c^n s' z)$$

Assume that the lemma holds for $n-1$. Then, Axiom 6 implies that

$$\frac{1}{2}(\bar{c}^0 c^1 \ldots c^{n-1} s'\bar{z}) + \frac{1}{2}(\bar{c}^0 c^1 \ldots c^{n-1} s' z) \sim_s \frac{1}{2}(\bar{c}^0\bar{c}^1 \ldots \bar{c}^n s'\bar{z}) + \frac{1}{2}(\bar{c}^0\bar{c}^1 \ldots \bar{c}^{n-1} s' z)$$

Therefore, we conclude that

$$
\begin{aligned}
W(s,&(c^0 c^1 \ldots c^{n-1} s'\bar{z})) - W(s,(c^0 c^1 \ldots c^{n-1} s' z)) \\
&= W(s,(\bar{c}^0 c^1 \ldots c^{n-1} s'\bar{z})) - W(s,(\bar{c}^0 c^1 \ldots c^{n-1} s' z)) \\
&= W(s,(\bar{c}^0\bar{c}^1 \ldots \bar{c}^{n-1} s'\bar{z})) - W(s,(\bar{c}^0\bar{c}^1 \ldots \bar{c}^{n-1} s' z))
\end{aligned}
$$

and hence the Lemma holds for $n$. Observe that Axiom 5 implies that the Lemma holds for $n = 1$. $\square$

It is easy to show that if $\succeq$ satisfies Axioms 3, 6 and 7 then it also satisfies the following stronger version of the independence axiom:

**Axiom 3\*:** $x \succ_s y$, $\alpha \in (0,1)$ *implies* $\alpha x + (1-\alpha)z \succ_s \alpha y + (1-\alpha)z$.

Theorem 1 of Gul and Pesendorfer (2000) establishes that $\succeq_s$ satisfies Axioms 1, 2, 4 and 3\* if and only if there exist $\hat{W}, \hat{U}, \hat{V}$ such that

$$\hat{W}(s,z) := \max_{\mu \in z}\{\hat{U}(s,\mu) + \hat{V}(s,\mu)\} - \max_{\nu \in z}\hat{V}(s,\nu)$$

for all $z \in Z$ and $\hat{W}$ represents $\succeq$. Moreover, the functions $\hat{W}, \hat{U}, \hat{V}$ are continuous and linear in their second arguments. Fix $\bar{s}$ and define

$$W(s,y) := \hat{W}(\bar{s}, \{\mu\})$$

where $\{\mu\} \in Z^K(s)$ and $\mu^{K+1}(y) = 1$. Since $Z^K(s)$ is homeomorphic to $Z$ and $\hat{W}$ is continuous in its second argument, $W$ is continuous in both arguments.

**Claim 1:** $W$ represents $\succeq$. Moreover, there exist continuous functions $U, V$ such that

$$W(s, z) := \max_{\mu \in z}\{U(s, \mu) + V(s, \mu)\} - \max_{\nu \in z} V(s, \nu)$$

and $W, U, V$ are linear in their second arguments.

**Proof:** Axiom 6 implies $W(s, x) \geq W(s, y)$ iff $\hat{W}(s, x) \geq \hat{W}(s, y)$. Therefore, $W$ represents $\succeq$. Note that $\hat{W}$ is linear in its second argument. Let $\{\mu\}, \{\nu\}, \{\eta\} \in Z^K(s)$ with $\mu^{K+1}(x) = \nu^{K+1}(y) = \eta^{K+1}(\alpha x + (1-\alpha)y) = 1$. Axiom 7 and linearity of $\hat{W}$ in its second argument imply that

$$\begin{aligned}
W(s, \alpha x + (1 - \alpha)y) &= \hat{W}(\bar{s}, \{\eta\}) \\
&= \hat{W}(\bar{s}, \alpha\{\mu\} + (1 - \alpha)\{\nu\}) \\
&= \alpha\hat{W}(\bar{s}, \{\mu\}) + (1 - \alpha)\hat{W}(\bar{s}, \{\nu\}) \\
&= \alpha W(s, x) + (1 - \alpha)W(s, y)
\end{aligned}$$

Thus, $W$ is linear in its second argument. It follows that $W(s, z) = \alpha(s)\hat{W}(s, z) + \beta(s)$ for some $\alpha, \beta : S \to I\!\!R$ such that $\alpha(s) \geq 0$. Since $\succeq$ is regular, $\alpha(s) > 0$ for all $s$. Hence, $U = \alpha\hat{U} + \beta, V = \alpha\hat{V}$ and the $W$ have the desired properties. $\qquad\square$

**Claim 2:** Let $\{\mu_h\}, \{\mu_l\} \in Z^K(s)$ with $\mu_h^{K+1}(y_h) = \mu_l^{K+1}(y_l) = 1$. Then, $W(s', \{\mu_h\}) - W(s', \{\mu_l\}) = W(s'', \{\mu_h\}) - W(s'', \{\mu_l\})$ for all $s', s''$.

**Proof:** Let $\{\bar{\mu}_{hl}\}, \{\bar{\mu}_{hh}\} \in Z^{2K}(s', s)$ and $\{\bar{\mu}_{lh}\}, \{\bar{\mu}_{ll}\} \in Z^{2K}(s'', s)$ satisfy $\bar{\mu}_{hh}^{2K+1}(y_h) = \bar{\mu}_{lh}^{2K+1}(y_h) = 1$ and $\bar{\mu}_{ll}^{2K+1}(y_l) = \bar{\mu}_{hl}^{2K+1}(y_l) = 1$. Let $z = \{.5\bar{\mu}_{hh} + .5\bar{\mu}_{ll}\}$ and $x = \{.5\bar{\mu}_{hl} + .5\bar{\mu}_{lh}\}$. By Lemma 1, $x \sim_{\bar{s}} z$. Hence, $\hat{W}(\bar{s}, x) = \hat{W}(\bar{s}, z)$ and thus

$$\hat{W}(\bar{s}, \{\bar{\mu}_{hh}\}) - \hat{W}(\bar{s}, \{\bar{\mu}_{hl}\}) = \hat{W}(\bar{s}, \{\bar{\mu}_{lh}\}) - \hat{W}(\bar{s}, \{\bar{\mu}_{ll}\}) \tag{6}$$

Recall that

$$\hat{W}(\bar{s}, \{\bar{\mu}\}) = W(s', \{\mu\})$$

if $\{\bar{\mu}\} \in Z^{2K}(s',s), \{\mu\} \in Z^K(s)$ and $\mu^{K+1}(y) = \bar{\mu}^{2K+1}(y)$. Substituting $W$ for $\hat{W}$ in equation (6) then proves the claim. $\qquad\square$

**Claim 3:** There exist $\delta : S \times C \to (0,\infty)$ and $u : S \times C \to \mathbb{R}$ such that $U(s,\nu) = \int [u(s,c) + \delta(s,c)W(sc,z)]d\nu(c,z)$ for all $s \in S, \nu \in \Delta$.

**Proof:** Since $U(s,\cdot)$ is linear and continuous, it has an integral representation. That is;

$$U(s,\nu) = \int U(s, \mu_{(c,z)})d\nu(c,z)$$

By Axiom 6, $U(s, \mu_{(c,\cdot)})$ and $W(sc, \cdot)$ yield the same linear preferences over $Z$. By regularity, neither function is constant. It follows that $U(s, \mu_{(c,\cdot)})$ is a strictly positive affine transformation of $W(sc, \cdot)$. Hence, for some $u, \delta$,

$$U(s, \mu_{(c,\cdot)}) = u(s,c) + \delta(s,c)W(sc,y)$$

where $\delta(s,c) > 0$ for all $s \in S, c \in C$. Therefore,

$$U(s,\nu) = \int [u(s,c) + \delta(s,c)W(sc,y)]d\nu(c,z)$$

as desired $\qquad\square$

**Claim 4:** The function $\delta(\cdot)$ in Claim 3 is constant.

**Proof:** Let $k \in 1,...,K+1$ denote the smallest integer such that $c^n = \bar{c}^n$ for $n \le k$ implies $\delta(c^1,...,c^{K+1}) = \delta(\bar{c}^1,...,\bar{c}^{K+1})$. Let $(s, c^{K+1}) := (c^1,...,c^{K+1})$ and $(s_*, c_*^{K+1}) := (c_*^1,...,c_*^{K+1})$ where $c^n = c_*^n, \forall n \le k-1$.

Pick any $c \in C$. Let $s' = (c,\ldots,c,c^1,c^2,\ldots,c^{k-1})$ and fix any $\hat{s}$. By regularity there are $y_h, y_l \in Z$ such that $W(\hat{s}, y_h) > W(\hat{s}, y_l)$. Let $\{\bar{\mu}_{hl}\}, \{\bar{\mu}_{hh}\} \in Z^{2K-k+1}(c^k,\ldots,c^{K+1},\hat{s})$ and $\{\bar{\mu}_{lh}\}, \{\bar{\mu}_{ll}\} \in Z^{2K-k+1}(c_*^k,\ldots,c_*^{K+1},\hat{s})$ be such that $\bar{\mu}_{hh}^{2K-k+2}(y_h) = \bar{\mu}_{lh}^{2K-k+2}(y_h) = 1$ and $\bar{\mu}_{hl}^{2K-k+2}(y_l) = \bar{\mu}_{ll}^{2K-k+2}(y_l) = 1$. Let $x = \{.5\bar{\mu}_{hh} + .5\bar{\mu}_{ll}\}$ and $z = \{.5\bar{\mu}_{hl} + .5\bar{\mu}_{lh}\}$. By Lemma 1, $x \sim_{s'} z$. Hence, $W(s',x) = W(s',z)$.

34

Let $\{\mu_h\}, \{\mu_l\} \in Z^K(\hat{s})$ be such that $\mu_h^{K+1}(y_h) = \mu_l^{K+1}(y_l) = 1$. Applying Claim 3 repeatedly and using the fact that $\delta(s, c) = \delta(\bar{s}, \bar{c})$ whenever $c^n = \bar{c}^n$, $\forall n \leq k$ establishes $W(s', x) - W(s', z) = 0$ iff

$$\delta(s, c^{K+1})W(sc^{K+1}, \{\mu_h\}) + \delta(s_*, c_*^{K+1})W(s_*c_*^{K+1}, \{\mu_l\}) =$$
$$\delta(s, c^{K+1})W(sc^{K+1}, \{\mu_l\}) + \delta(s_*, c_*^{K+1})W(s_*c_*^{K+1}, \{\mu_h\})$$

Rearranging, this implies

$$\delta(s, c^{K+1})(W(sc^{K+1}, \{\mu_h\}) - W(sc^{K+1}, \{\mu_l\}) =$$
$$\delta(s_*, c_*^{K+1})(W(s_*c_*^{K+1}, \{\mu_h\}) - W(s_*c_*^{K+1}, \{\mu_l\})$$

By Axiom 6, $W(s, \{\mu_h\}) - W(s, \{\mu_l\}) > 0$. Hence, $\delta(s, c^{K+1}) = \delta(s_*, c_*^{K+1})$ by Claim 2. Therefore $k = 0$, which is the desired conclusion. $\square$

**Claim 5:** *Let $\delta \in \mathbb{R}$ denote the constant function in Claim 3. Then, $0 < \delta < 1$.*

**Proof:** That $\delta > 0$ has already been established. Pick any $c \in C$ and let $s = (c, c, \ldots, c)$. Let $z_c$ denote the unique $z = \{\mu\} \in Z$ such that $\mu(c, z) = 1$. Pick $y_1 \in Z$ such that $W(s, y_1) \neq W(s, z)$. By regularity, such a $y_1$ exists. Define $y_n \in Z$ inductively as follows: $y_n = \{\mu\}$ such that $\mu(c, y_{n-1}) = 1$. Note that $y_n$ converges to $z$. Hence, by continuity, $W(s, z) - W(s, y_n)$ must converge to 0. But, by Claims 3 and 4 $W(s, z) - W(s, y_n) = \delta^{n-1}(W(s, y_1) - W(s, z)) \neq 0$. Hence, $\delta < 1$. $\square$

Claims $1 - 5$ establish the existence of the desired representation.

To conclude the proof, let $\delta \in (0, 1)$ and $u : S \times \Delta(C) \to \mathbb{R}$ and $v : S \times C \to \mathbb{R}$ be continuous functions.

**Lemma 2 (A Fixed-Point Theorem):** *If $B$ is a closed subset of a Banach space with norm $\|\cdot\|$ and $T : B \to B$ is a contraction mapping (i.e., for some integer $m$ and scalar $\alpha \in (0, 1)$, $\|T^m(W) - T^m(W')\| \leq \alpha \|W - W'\|$ for all $W, W' \in B$), then there is a unique $W^* \in B$ such that $T(W^*) = W^*$.*

**Proof:** See [Bertsekas and Shreve (1978), p. 55] who note that the theorem in Ortega and Rheinholt (1970) can be generalized to Banach spaces. $\square$

Let $\mathcal{W}_b$ be the Banach space of all real-valued, bounded functions on $S \times Z$ (endowed with the sup norm). The operator $T : \mathcal{W}_b \to \mathcal{W}_b$, where

$$TW(s,z) = \max_{\mu \in z}\{u(s,\mu^1) + V(s,\mu) + \delta \int W(sc,x)d\mu(c,x)\} - \max_{\nu \in z} V(s,\nu)$$

is well-defined and is a contraction mapping. Hence, by Lemma 2, there exists a unique $W$ such that $T(W) = W$. To prove that $W$ is continuous, repeat the above argument for the subspace $\mathcal{W}_c \subset \mathcal{W}_b$ of all continuous, real valued functions on $S \times Z$. Note that $T(\mathcal{W}_c) \subset \mathcal{W}_c$. Hence, again by Lemma 2, $T$ has a fixed point $W^* \in \mathcal{W}_c$. Since $W$ is the unique fixed-point of $T$ in $\mathcal{W}_b$, we have $W^* = W$. Hence, $W$ is continuous.

For any $W, u, V, \delta$ such that

$$W(s,z) = \max_{\mu \in z}\{\int[u(s,c) + V(s,(c,x)) + \delta W(sc,x)]d\mu(c,x)\} - \max_{\nu \in z} V(s,\nu)$$

define $\succeq_s$ by $x \succeq_s y$ iff $W(s,x) \geq W(s,z)$. Verifying that $\succeq_s$ satisfies Axioms $1 - 7$ is straightforward. $\qquad \square$

## 8.3   Proof of Theorem 2

By Theorem 3, $\succeq$ can be represented by a continuous $W$ where

$$W(s,z) = \max_{\mu \in z}\{u(s,\mu^1) + v(s,\mu^1) + \delta \int W(sc,x)d\mu(c,x)\} - \max_{\nu \in z} v(s,\nu^1)$$

for some continuous $u, v$ and $\delta \in (0,1)$. Moreover, $W, u, V$ are linear in their second arguments. Let $U(s,\mu) = u(s,\mu^1) + \delta \int W(sc,x)d\mu(c',x)$

**Claim 6:** $V(s,\nu) = V(s,\hat\nu)$ whenever $\nu^1 = \hat\nu^1$.

**Proof:** If $V(s,\cdot) = \alpha U(s,\cdot) + \beta$ for some $\alpha \leq -1$, then $x \succeq_s y$ for all $x \subset y$ contradicting regularity. If $V(s,\cdot) = \alpha U(s,\cdot) + \beta$ for some $\alpha \geq 0$ then $x \succeq_s y$ for all $y \subset x \in Z$ again, contradicting regularity. Hence, for each $s \in S$ there are two possibilities: either $V(s,\cdot)$ is not an affine transformation of $U(s,\cdot)$ or there exists $\alpha \in (-1,0)$ such that $V(s,\cdot) = \alpha U(s,\cdot) + \beta$. In either case, there exist $\mu^s, \nu^s \in \Delta$ such that $U(s,\mu^s) + V(s,\mu^s) > U(s,\nu^s) + V(s,\mu^s)$ and $V(s,\mu^s) < V(s,\nu^s)$.

Take any $\nu, \hat{\nu} \in \Delta$ such that $\nu^1 = \hat{\nu}^1$. There exists $\alpha > 0$ small enough so that

$$U(s, \mu^s) + V(s, \mu^s) > U(s, \alpha\nu + (1-\alpha)\nu^s) + V(s, \alpha\nu + (1-\alpha)\nu^s)$$

$$U(s, \mu^s) + V(s, \mu^s) > U(s, \alpha\hat{\nu} + (1-\alpha)\nu^s) + V(s, \alpha\hat{\nu} + (1-\alpha)\nu^s)$$

$$V(s, \mu^s) < V(s, \alpha\nu + (1-\alpha)\nu^s)$$

$$V(s, \mu^s) < V(s, \alpha\hat{\nu} + (1-\alpha)\nu^s)$$

Then, linearity and Assumption I imply $\{\alpha\nu + (1-\alpha)\nu^s, \mu^s\} \sim_s \{\alpha\hat{\nu} + (1-\alpha)\nu^s, \mu^s\}$. Since $W$ represents $\succeq$ we have $V(s, \alpha\nu + (1-\alpha)\nu^s) = V(s, \alpha\hat{\nu} + (1-\alpha)\nu^s)$. Since $V$ is linear, we conclude $V(s, \nu) = V(s, \hat{\nu})$ as desired. $\qquad\square$

Regularity implies that neither $U(s, \cdot)$ nor $v(s, \cdot)$ is constant. Claim 6 then implies that $v(s, \cdot)$ is not an affine transformation of $U(s, \cdot)$. Hence, we may apply Theorem 7 of Gul and Pesendorfer (2000a) to yield the following implications:

**Fact 1:** *(Theorem 7 (Gul and Pesendorfer (2000a)) $s'\mathbf{P}s$ iff for some $\alpha_u, \alpha_v \in [0,1], \gamma > 0, \gamma_u, \gamma_v \in \mathbb{R}$*

$$\gamma U(s, \mu) = \alpha_u U(s', \mu) + (1-\alpha_u)v(s', \mu^1) + \gamma_u$$

$$\gamma v(s, \mu^1) = \alpha_v U(s', \mu) + (1-\alpha_v)v(s', \mu^1) + \gamma_v$$

*for all $\mu$.*

By Assumption P, $s'\mathbf{P}s$ or $s\mathbf{P}s'$. Without loss of generality assume $s'\mathbf{P}s$. By regularity there exists $c, x, y$ such that $U(s, (c, x)) > U(s, (c, y))$. Since $v(s, (c, x)) = v(s, (c, y))$ it follows that $\alpha_v = 0$. Pick any $s^0 \in S$. We conclude that for all $s \in S$

$$U(s, \mu) = \alpha_u(s)U(s^0, \mu) + \beta_u(s)v(s^0, \mu^1) + \gamma_u(s)$$

$$v(s, \mu^1) = \beta_v(s)v(s^0, \mu^1) + \gamma_v(s) \tag{7}$$

for some functions $\alpha_u, \beta_u, \beta_v, \gamma_u, \gamma_v$ such that $\alpha_u(s) > 0, \beta_v(s) > 0$ for all $s$. Note that $U$ and $v$ are continuous and hence $\alpha_u, \beta_u, \gamma_u, \gamma_v, \beta_v$ are continuous. Hence,

$$\int [u(s, c) + \delta W(sc, z)]d\nu(c, z) =$$

$$\int [\alpha_u(s)u(s^0, c) + \beta_u(s)v(s^0, c) + \gamma_u(s) + \alpha_u(s)\delta W(s^0 c, z)]d\nu(c, z) \tag{8}$$

The only terms on either side of (8) that depend on $\nu^2$ are $\delta W(sc, z)$ and $\alpha_u(s)\delta W(s^0 c, z)$. Since regularity implies that neither of these terms is constant it follows that

$$W(sc, \cdot) = \alpha_u(s)W(s^0 c, \cdot) + A(s, c)$$

Then, Claim 2 (in the proof of Theorem 1) implies that $\alpha_u(s) = 1$ for all $s$. It follows that $W(s^0 c, \cdot)$ represents $\succeq_{sc}$. Hence, $K = 1$. That is, $sc = c$ for all $s, c$. Henceforth, we write $c$ instead of $sc$.

Let $W_0(c, z) = W(c, z) - \gamma_u(c)$, $\hat{u}_0(c) = u(s^0, c) + \delta\gamma_u(c)$ for all $c$. Let $\hat{v}_0(c) = v(s^0, c)$ and $\hat{v}_0(\nu^1) = \int \hat{v}_0(c)d\nu^1(c)$. Then,

$$
\begin{aligned}
W_0(c, z) = W(c, z) - \gamma_u(c) &= \max_{\mu \in z}\{U(c, \mu) + v(c, \mu)\} - \max_{\nu \in z} v(c, \nu) - \gamma_u(c) \\
&= \max_{\mu \in z}\{U(s^0, \mu) + \beta_u(c)v(s^0, \mu) + \beta_v(c)v(s^0, \mu)\} - \max_{\nu \in z} \beta_v(c)v(s^0, \nu) \\
&= \max_{\mu \in z} \int [u(s^0, c') + \beta_u(c)v(s^0, c') + \beta_v(c)v(s^0, c') + \delta W(c', x)]d\mu(c', x)\} \\
&\quad - \max_{\nu \in z} \beta_v(c)\hat{v}_0(\nu) \\
&= \max_{\mu \in z} \int [\hat{u}_0(c') + \beta_u(c)\hat{v}_0(c') + \beta_v(c)\hat{v}_0(c') + \delta W_0(c', x)]d\mu(c', x)\} \\
&\quad - \max_{\nu \in z} \beta_v(c)\hat{v}_0(\nu^1)
\end{aligned}
$$

Let $\bar{\beta}_u := \max \beta_u(s), \underline{\beta}_u := \min \beta_u(s)$ and $\underline{\beta}_v := \min \beta_v(s)$. By continuity $\bar{\beta}_u, \underline{\beta}_u, \underline{\beta}_v$ are well-defined and since each $\succeq_s$ is regular $\underline{\beta}_v > 0$. Let

$$u_0 = \hat{u}_0 + \underline{\beta}_u \hat{v}_0$$

If $\bar{\beta}_u > \underline{\beta}_u$, let

$$\hat{\pi} = \frac{\beta_u - \underline{\beta}_u}{\bar{\beta}_u - \underline{\beta}_u}$$

$$\hat{\sigma} = \beta_v + \pi(\bar{\beta}_u - \underline{\beta}_u)$$

Otherwise, let

$$\hat{\pi} \equiv 0$$

$$\hat{\sigma} = \beta_v$$

It is easy to verify that $W_0$ represents $\succeq$. Fix any $b^0$ and let $v_0(d) := \hat{v}_0(b^0, d)$, $\pi_0(d) := \hat{\pi}_0(b^0, d)$ and $\sigma_0(d) := \hat{\sigma}_0(b^0, d)$. By Assumption N, $\hat{v}_0(b, d) = v_0(d)$, $\hat{\pi}_0(b, d) = \pi_0(d)$

and $\hat{\sigma}_0(b, d) = \sigma_0(d)$ for all $b \in [0, 1]^{l-1}$. Assumption N also implies that $v_0$ is strictly increasing. Hence, $u_0, v_0, \pi, \sigma, \delta$ satisfy all the desired properties. $\qquad\square$

# References

1. Bertsekas D. P. and S. E. Shreve, "Stochastic Optimal Control: The Discrete Time Case", Academic Press, New York 1978.

2. Becker, G. S., and K. Murphy, "A Theory of Rational Addiction", Journal of Political Economy, 1988, 675-700.

3. Chick, J., "Emergent Treatment Concepts", *Annual Review of Addictions Research and Treatment*, 1992, 297-312.

4. Davies, J. B. "The Myth of Addiction; the Application of the Pschological Theory of Attribution to Illicit Drug Use," Harwood Academic Publishers, Chur 1992.

5. Gul, F. and W. Pesendorfer, "Temptation and Self-Control", 2000a, forthcoming in: *Econometrica.*

6. O'Donoghue, T. and M. Rabin, (1998)"Addiction and Self-Control", in Addiction: Entries and Exits, Jon Elster, editor, Russel Sage Foundation, 1999.

7. Pentel, P.R.; Malin, D.H.; Ennifar, S.; Hieda, Y.; Keyler, D.E.; Lake, J.R.; Milstein, J.R.; Basham, L.E.; Coy, R.T.; Moon, J.W.D.; Naso, R.; Fattom, A. ,"A Nicotine Conjugate Vaccine Reduces Nicotine Distribution to Brain and Attenuates Its Behavioral and Cardiovascular Effects in Rats - Role of dose and dose interval", *Pharmacology Biochemistry & Behavior,* Vol: 65, Issue: 1, January (2000), 191-198.

8. Robinson T. E. and K. C. Berridge, "The neural basis of drug craving: an incentive-sensitization theory of addiction", *Brain Research Reviews* 18, (1993), 247-291.