# Innovation and Diffusion of Medical Treatment*

Barton H. Hamilton
Olin Business School, Washington University in St. Louis


Andrés Hincapié
Department of Economics, University of North Carolina at Chapel Hill


Robert A. Miller
Tepper School of Business, Carnegie Mellon University


Nicholas W. Papageorge
Department of Economics, Johns Hopkins University, IZA and NBER

ABSTRACT: We develop and estimate a dynamic structural model of demand in a setting where product characteristics evolve in response to aggregate consumer choices. The direction and speed of innovation is inefficient because individuals do not account for their influence on innovation, which creates an externality. We apply the model to drugs invented to combat HIV, which differ in their efficacy and their propensity to cause side effects. We find that the externalities are quantitatively important and that a temporary subsidy would have increased average social welfare by improving average health and would have reduced inequality in lifetime utility across health groups.

---

# 1 Introduction

Economists have long recognized that innovation, including the entry of new products and the exit of obsolete ones, is not only determined by science and luck, but also responds to latent consumer demand (Hicks, 1932). Sometimes referred to as demand-pull innovation (Schmookler, 1966; Scherer, 1982), the responsiveness of innovation to demand generates an externality because the benefits an individual indirectly confers upon all (other) future individuals through his effect on innovative activity are not reflected in the price he pays for the product in the decentralized economy (Jovanovic and MacDonald, 1994; Waldfogel, 2003; Finkelstein, 2004).

This paper develops a dynamic structural model of demand for medical treatment that captures how product innovation is endogenous to patient choices. Agents in the model can choose a treatment on the market, join a clinical trial to access an experimental treatment, or refrain from consuming any treatment at all. Patient choices reflect a tradeoff between multiple product attributes: efficacy, which improves long-run health; and current-period side effects. To capture how innovation affects patients, the model incorporates an evolving set of available treatments. In each period, new products come available and older, low-quality products exit the market. The resulting evolution of the choice set is governed by a stochastic process that is a function of aggregate consumer behavior, including the share of patients in clinical trials. When making private medication choices, individual consumers do not take into account their impact on aggregate demand, which generates the externality.

Incorporating an aggregate process to capture an endogenously evolving choice set complicates an otherwise standard model of dynamic demand.[1] Doing so has two important benefits. First, it allows us to characterize a demand-driven externality in a structural model, which means we can evaluate counterfactual policies that could mitigate resulting inefficiencies. Second, agents in the model form beliefs about future available treatments using this stochastic process. Thus, we are able to model consumer beliefs in a market with uncertainty over future choice sets without resorting to simplifying assumptions, such as perfect foresight or fully

---

[1]The key complication is that the problem is no longer stationary.

unanticipated shocks, which can lead to potential biases in estimated parameters and inaccurate policy conclusions.

We apply the model to investigate the demand for a rapidly evolving set of treatments for HIV (human immunodeficiency virus). We identify the model using a biennial panel from four American cities that tracks a replenished panel of individuals along with the path of innovations for over twenty years, from when the market for HIV treatments emerged around 1984 until it matured. During this period, frequent, incremental innovations in medication were punctuated by sporadic breakthroughs. The data include an objective, continuous measure of health based on immune system status obtained from blood tests administered with the survey every six months. We use our estimates to quantify the magnitude of the demand externality in the market for HIV drugs.

Using the estimated model, we perform counterfactuals focused on the externality arising from endogenous innovation. Our results reveal that individuals' preferences tilt the path of innovation towards treatments with fewer side effects, away from the invention of more efficacious treatments. Moreover, a strong distaste for experimentation slows the diffusion of new, superior products as well as the development of future treatments in clinical trials. As a measure of the externality acting through demand for experimental products, we compute the marginal increase in aggregate welfare generated by a social planner who assigns the marginal patient in the decentralized economy to clinical trials at two different point in time (1991 and 1996). In year 1996 the marginal person (who does not want to join a trial) experiences a loss of roughly $600 by participating. However, because trial participation spurs innovation by increasing the expected quality and the expected number of new products, the net social gain is about $2,000 per individual. We also consider a more realistic policy in which the planner does not have enough information about individuals' preferences to assign the marginal participant to a trial. We show that an optimal Pigouvian subsidy pays all trial participants $16,000 and leads to welfare gains of roughly $5,000 for each individual in the economy. Broadly, our results indicate that providing monetary incentives for trial participation can improve welfare by accelerating the process of innovation.

This study contributes to a literature on dynamic demand under uncertainty. Fol-

2

lowing Petrin (2002), each product in our model is a bundle of characteristics—in our case, efficacy and side effects.[2] Moreover, similar to Gowrisankaran and Rysman (2012), product characteristics have dynamic impacts on consumers. Our work is more directly related to a set of studies that have augmented dynamic models of demand to incorporate the realities of health-related choices, including learning (Crawford and Shum (2005), Chan and Hamilton (2006), Fernandez (2013), Darden (2017), Dickstein (2018)), the relationships between side effects, health and work (Papageorge, 2016), risky behavior and equilibrium effects (Chan et al., 2016) and consumer experimentation with new products (Fernandez (2013), Chan and Hamilton (2006)). Similar to some of these studies, we consider uncertainty regarding commercialized products. One important departure from earlier work is to explicitly incorporate uncertainty also over current experimental products and products that may emerge in the future.

Another important departure from earlier work is to explicitly model the externality arising when a rapidly evolving choice set is endogenous to consumer demand. Several papers have investigated how demand affects innovation, e.g., through market size. For example, Finkelstein (2004) shows that policies promoting vaccine use accelerate the development of vaccines and Acemoglu and Linn (2004) relate market size to pharmaceutical innovation. Dranove et al. (2014) identify a "social value" of pharmaceutical innovation, showing that Medicare Part D spurred the development of some drugs. A common theme in this literature is that a demand externality arises if consumer behavior drives innovation. Waldfogel (2003) uses the term "preference externalities" to describe the mechanism through which market shares can influence products, thus benefitting individuals with similar tastes.[3] Bolton and Harris (1999) argue that a free-riding problem emerges if experimentation accelerates innovation. In our context, if consumer experimentation provides social benefits by spurring innovation, rational individuals may choose to do so less

---

[2]Studies pioneering the 'characteristics approach' include Stigler (1945), Lancaster (1966) and Rosen (1974).

[3]Demand externalities have been discussed in a variety of scenarios, including sorting into neighborhoods (Bayer and McMillan, 2012) and the emergence of food deserts (Allcott et al., 2017). In the context of obesity, Bhattacharya and Packalen (2012) provide evidence that individual efforts to prevent obesity can shrink the market size for obesity treatments, which slows technological progress.

than is socially optimal. While earlier work has examined how changes in markets or latent demand affect innovation, our approach structurally models this relationship, which allows us to evaluate counterfactual policies that could mitigate the resulting externality, such as paying individuals to participate in clinical trials.

While our focus in this study is on demand-driven innovation, we also relate to a literature on supply-side determinants of innovation. We highlight two related papers. Goettler and Gordon (2011) develop a model in which market structure (monopoly versus duopoly) affects innovation in the market for microprocessors.[4] Igami (2017) studies the market for hard disk drives as it transitions from one product generation to the next (5.25- to 3.5-inch). In his model firms have perfect foresight over exogenously evolving demand, and play a dynamic game in which innovation amounts to introducing the single new product generation. In both papers the unidimensional state of the art becomes the starting point for future innovations. Given their settings, the supply side is specified to allow for strategic interaction among firms and the models are well-suited for examining counterfactual policies related to firm behavior, such as market structure.

Our setting is different and thus requires a different approach. Roughly 11 firms produced HIV medications, making it difficult to assume monopolistic or duopolistic competition. Moreover, in our context, products are not uni-dimensional in quality, the size of innovations does not remain constant and perfect foresight is a poor approximation. Finally, we have rich data on demand and little information on firms. Given our data and context, we treat the entry of new products as coming from a process that is related to existing technology and to aggregate demand. However, our framework supports multidimensional products, multiple product entry, variable changes in technology (both incremental innovations and breakthroughs) and new products that are not necessarily technological improvements.[5] All of these factors play a critical role in determining demand by heterogeneous consumers. On the consumer side we allow for substantial heterogeneity on observables (objective

---

[4]They find that the presence of a second firm can slow innovation because no firm expects to capture all profits.

[5]This is a feature in our data and an equilibrium that emerges naturally in models where individuals are not fully informed about new product characteristics (Miller, 1988).

health status, race, age, education, previous consumption and labor participation).

While the stochastic process of innovation we specify captures many features of our setting, it is important to note that it uses reduced-form objects to approximate a complex innovation environment, thus limiting the types of counterfactuals we can perform. For example, our model would be ill-suited for analyses of policies that would affect the demand externality primarily through firm behavior (e.g., changes to market structure as in Goettler and Gordon (2011) and Igami (2017)). Instead, using a stochastic process allowing for broad variation in the kinds of innovations that can occur and given a rich model of demand we focus on how changes to patient behavior (e.g., through remuneration for participation in a clinical trial) affect consumer welfare through impacts on the speed and direction of innovation.

Finally, we contribute to research on structural estimation by providing a simulation-based econometric method to estimate models of endogenous innovation. Methodologically, our empirical strategy builds on Hotz and Miller (1993), Hotz et al. (1994) and Altuğ and Miller (1998) in using conditional choice probabilities (henceforth, *CCPs*) and forward simulation techniques to incorporate how individuals form expectations about future innovations. In our context, the individual's choice set evolves stochastically as a function of endogenous product exit and entry. The latter is determined by the innovation process which contains two components: unexpected, aggregate supply shocks and a systematic component, endogenous to aggregate demand, captured by a multi-dimensional reference point for innovation.

The remainder of this paper is organized as follows. Section 2 provides a brief historical background, describes our data set, and motivates the model structure with patterns in the data. Section 3 describes the model. Section 4 analyzes identification and describes the estimation strategy. Section 5 presents parameter estimates. Section 6 provides model predictions about the likelihood of technological progress. Section 7 introduces our counterfactual regimes. Section 8 concludes.

## 2  Data

Our empirical application focuses on the market for HIV treatments which came into existence around 1984 with the beginning of the HIV pandemic, causing over

613,000 deaths in the U.S. by 2008.[6] HIV infection leads to a reduction in the ability of the immune system to fight off routine infections, a condition known as AIDS (acquired immunodeficiency syndrome). In developed countries, where access to medication is widespread and often subsidized, technological advancement has transformed HIV infection into a manageable condition with treatments whose side effects are fairly mild. This was not always the case. In the early years of the epidemic, available treatments were not only largely ineffective, but also had uncomfortable, painful and even deadly side effects. Over time many innovations appeared, most of them small, and some worse than existing technology—being more toxic without being more effective. In the mid-nineties, a new set of treatments collectively known as HAART (highly active anti-retroviral treatment) was introduced, transforming HIV from a virtual death sentence into a chronic condition.[7] Within two years, mortality rates fell by over 80% among HIV infected (HIV+) men (Bhaskaran et al., 2008). However, HAART also involved drugs that were highly toxic, driving some people to refrain from using them to avoid often intolerable side effects. Innovations occurring after the mid-nineties had fewer side effects, but were generally no more effective than earlier versions of HAART.

## 2.1 The MACS Data Set

We use public data from the Multi-center AIDS Cohort Study (MACS). The MACS is an ongoing longitudinal investigation (beginning in 1984) of HIV infection in men who have sex with men (MSM) conducted at four sites: Baltimore, Chicago, Pittsburgh and Los Angeles. At each semi-annual visit, survey data are collected on HIV+ men's treatment decisions, out-of-pocket treatment expenditures, and physical ailments (which can reflect drug side effects), along with sociodemographic

---

[6]For comparison, over the same period in the U.S., there were 508,000 homicides and U.S. deaths in World War II were just under 420,000. Currently, there are roughly 50,000 new infections and 13,000 deaths per year in the U.S. that are attributed to HIV/AIDS. Globally, the number of deaths due to HIV/AIDS stands at roughly 35,000,000.

[7]There is no vaccine or cure for HIV or AIDS, but HAART is the current standard treatment. In general, 1996 is marked as the year when two crucial clinical guidelines that comprise HAART came to be commonly acknowledged. First, protease inhibitors (made widely available towards the end of 1995) would be an effective HIV treatment. Second, several anti-retroviral drugs taken simultaneously could indefinitely delay the onset of AIDS.

information, such as labor supply, income, race, and education. In addition, blood tests are administered at each visit to objectively measure health status. Our measure of health, the *CD4 count*, is based on the immune system and is defined as the number of white blood cells per cubic millimeter of blood. Absent HIV infection, a normal count ranges between 500 and 1500. For HIV+ individuals, a count below 500 indicates that the immune system has begun to deteriorate. However, such individuals may remain asymptomatic. When the CD4 count drops below about 300, a patient is said to suffer from AIDS and his immune system becomes unable to fight off routine infections, which compromises his survival probability.[8] Few data sets contain such objective, continuous measures of health and detailed treatment data along with economic information, making the MACS data set uniquely well-suited for our analysis of demand-pull innovation in the market for medical treatments. However, our data set is not ideal as it does not contain information on treatment prices. In our empirical work we approximate the cost of market treatments using out-of-pocket expenditures after controlling for objective health and other observables.

The full MACS data set we start with contains information on 6,972 subjects at 49 semi-annual visits for a total of 111,271 observations in the form of subject-visit dyads. We limit our attention to HIV+ individuals, leaving us with 47,753 observations. Due to a lack of data on gross income and out-of-pocket expenditures at earlier visits, we use two samples, a larger sample (20,142 observations) covering visits 6 to 49 which only includes health status and product usage, and a smaller sample (16,851 observations) that starts at visit 14 (roughly, late 1990) containing all variables. The construction of both samples is described in Appendix A. The smaller sample comprises 1,719 males, 68 percent white, 22 percent black and the rest Hispanic; 86 percent received some secondary education or more, and 23 percent attended graduate school.[9] Underscoring the gravity of HIV infection, about 40 percent of the HIV+ subjects we observe at least once die prior to the end of the

---

[8]The CD4 cutoff below which AIDS occurs varies between 200 and 350.

[9]Participation in clinical trials for experimental treatments has been shown to be lower among African-Americans, which may reflect different costs associated with treatments or differences in expected health outcomes (Harris et al., 1996). It may also reflect distrust in the medical system due to the Tuskegee experiment (Alsan and Wanamaker, 2018).

sample period.

Table 1 shows that the share of observations with positive physical ailments is 0.43 and the average CD4 count is 475, in the smaller sample. The share of observations with positive labor supply is 0.63. There is substantial variation in labor supply; 74 percent (68 percent) of unique individuals are observed working (not working) at least once.[10] The share of observations with positive market product consumption and trial product consumption are 0.65 and 0.07, respectively. There is also variation in treatment consumption; 83 percent of unique individuals are observed using a market product at least once and 24 percent opt for early access by participating in a clinical trial at least once during the sample period, suggesting a willingness to experiment with products of uncertain quality.

**TABLE 1:** Summary Statistics: Subjects-Visits. Visits 14-47 (1990-2007)

|  | Sample | Pre Haart | Post Haart |
|---|---|---|---|
| Obs | 16851 | 6972 | 9879 |
| Ailments | 0.43 | 0.45 | 0.41 |
| Market Product | 0.65 | 0.49 | 0.76 |
| Trial Product | 0.07 | 0.09 | 0.05 |
| Work | 0.63 | 0.70 | 0.58 |
| Age | 44.48 | 40.89 | 47.01 |
|  | (8.03) | (6.99) | (7.75) |
| CD4 | 475 | 407 | 524 |
|  | (297) | (298) | (287) |
| Gross Income | 17567 | 19036 | 16531 |
|  | (8787) | (8733) | (8677) |
| Out-of-pocket Expenditures | 266 | 179 | 327 |
|  | (706) | (598) | (767) |

Notes: Standard deviation in parentheses. Gross income and out-of-pocket expenditures are semestral and measured in real dollars of 2000. Pre HAART era corresponds to visit $\leq$ 24 or roughly before 1996.

## 2.2 Key Empirical Patterns

Next, we discuss five key patterns in the data that we incorporate into our structural model of demand under endogenous innovation.

**Individuals respond to technological change.** A distinguishing feature of the market for HIV treatments is that innovations in product quality have life-saving

---

[10]This is consistent with results in Papageorge (2016) who studies labor supply and medication usage with the MACS data.

effects. Figure 1(a) shows that prior to the introduction of HAART, death rates were much higher despite a multitude of new treatments becoming available. After HAART, death rates plunge, and continue to fall until 2007, as smaller innovations occurred that made drugs incrementally more effective and less toxic. Table 1 above shows that improvements in survival coincide with improvements in immune system health as measured by the CD4 count. Improvements in health and survival occur as our sample ages and becomes less likely to participate in the labor market (12 percentage points less after 1995), which is reflected in the reduction of unconditional average semestral gross income from about \$19,036 in the pre-HAART era to \$16,531 after 1995.



(a) Survival        (b) Market Treatment        (c) Trial Treatment

**FIGURE 1:** Survival and Consumer Demand over Time

Notes: Left panel shows the probability of dying between periods $t$ and $t+1$ conditional on surviving until $t$. More than 1500 surveyed individuals died for AIDS-related causes during our analysis period. The middle and right panels show consumption by health status.

Table 1 above also shows that improvements in product quality induce individuals to consume more HIV treatments. The share of individuals consuming a market product went from 0.49 in the pre-HAART era to 0.76 after HAART was introduced, and individuals' out-of-pocket expenditures went from \$179 to \$327 per semester. Figure 1(b) shows that consumption of market treatments differed across health levels prior to the introduction of HAART. Individuals with low CD4 counts were more likely to use available medications, which were relatively ineffective, while healthier individuals often avoided treatment altogether. Demand for treatment increased and converged across health levels in response to the introduction

9

of more effective products after HAART.

Beyond market treatments, HIV+ individuals often have the option to consume experimental products in clinical trials. The most dramatic feature of Figure 1(c) is the spike in trial treatment around the time HAART was introduced. Early trial participation is driven largely by individuals with low CD4 counts, suggesting that less healthy individuals may be more willing to consume experimental products of uncertain qualities.[11] Once effective treatments are available, trial participation is no longer driven by sick people willing to face uncertainty in exchange for early access to a product of potentially higher quality.

**Product characteristics are multidimensional.** Treatment consumption was never universal. While healthy patients often avoided treatment before the introduction of life-saving innovations around 1995, some individuals at risk continued to go untreated thereafter. Figure 1(b) shows that treatment consumption climbs to roughly 80% after the introduction of HAART. In part, this happens because products are costly, though out-of-pocket costs for medical care do not differ much across treatment choices.
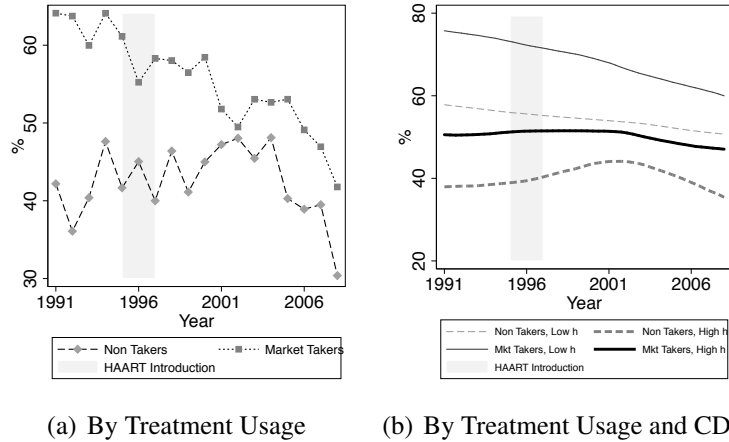
Another possibility is that patients avoid effective medications due to side effects. Figure 2(a) provides support for this view, showing that individuals who consume a market product suffer more physical ailments (e.g., nausea or cramping) and Figure 2(b) shows that this result holds even after controlling for underlying immune system health. Moreover, as products become less toxic (causing fewer side effects) over time, the gap in ailments between those who are treated and those who are not decreases.

Consumption patterns are thus consistent with the idea that treatments are multi-attribute products: drug effectiveness at improving underlying health and drug propensity to cause side effects, which compromise quality of life. In the presence of multiple attributes, it is often the case that products cannot be perfectly ranked from best to worst. Rather, patients face a tradeoff between investing in their un-

---

[11]In the years just prior to HAART introduction the efficacy of products in the market had increased, pushing up the reference point for innovation and thus attracting more individuals into clinical trials. See Figure S1 in Appendix A.

derlying health versus maintaining their quality of life by avoiding treatments with
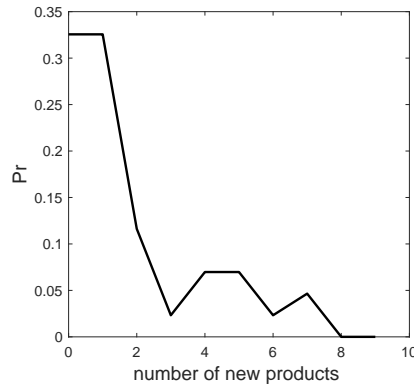harsh side effects.



(a) By Treatment Usage      (b) By Treatment Usage and CD4

**FIGURE 2:** Physical Ailments by Treatment Usage and CD4

Notes: Figure contains mean of ailments indicator over time. "Mkt Takers" refers to individuals consuming a market
treatment. "High h" refers to individuals with CD4 counts of 250 and above.

**The number of new treatments fluctuates over time.** We define a product (or
treatment) as a combination of single-product components.[12] (See Appendix A.)
This means that both AZT and the combination of AZT+3TC+Saquinavir are ex-
amples of products in our framework. This definition results from noting that the
interactions between components matter, and hence the sum of effects of consuming
each drug individually does not equal the effect of a treatment formed by the sum
of the drugs. Additionally, this definition corresponds to the nature of the market,
where large treatment innovations such as HAART are themselves combinations of
product components. By this definition 86 products were introduced to the market
over the sample period with substantial variation in the number of new treatments

---

[12]The number of firms introducing product components as well as the ownership of the firms
changes over time. For example, the first product component (AZT) was introduced by Burroughs-
Wellcome in 1987 which became Glaxo-Wellcome in 1995, GlaxoSmithKline in 2000, and trans-
ferred its HIV assets to the joint venture ViiV created with Pfizer in 2009. By the mid 1990's at least
6 firms had introduced product components and had valid patents (Glaxo Wellcome, Bristol-Myers
Squibb, Hoffmann-La Roche, Abbott, Merck, and Boehringer Ingelheim). In total, there were at
least 11 firms that introduced product components during the period we study.

introduced each period.[13] Figure 3 shows that the unconditional probability of observing more than one product being introduced in a given period is more than 30%, suggesting that product introduction has an intensive margin.
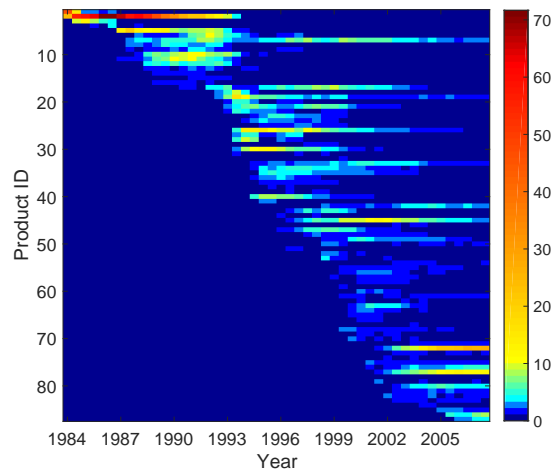


**FIGURE 3:** Empirical Distribution of Number of New Treatments.

**Market concentration shifts as innovation progresses and the market matures.**
Substantial variation in the number of new products (Figure 3) along with consumer preferences for multiple dimensions of drug quality is reflected in both innovation and market concentration. Figure 4 shows innovation and diffusion of new products over time using a heat map—dark colors corresponds to low (or zero) market share and warmer colors indicate higher market shares. Early on there are a few products with high shares. As time passes new products strip market share from incumbents and less popular products exit. Low market shares are common in the years following HAART introduction around 1995, when many new treatments were introduced, most of which were effective but with strong side effects. Consumers often switched among options depending on their health in an attempt to balance health accumulation with a preference for treatments without side effects. As the market matured, effective treatments with fewer side effects entered the market, removing the tradeoff between the two treatment qualities and increasing market concentration once again.

---

[13]Table S1 in Appendix A presents our market products including the individual drugs they are composed of as well as their entry and exit time as observed in our data.

**FIGURE 4:** Diffusion of Products Over Time

Notes: HIV treatments from 1984 to 2008. Each ID—or row—represents a product. Color indicates the share of the market that the product captures. Shares are conditional on individuals who consumed a product.
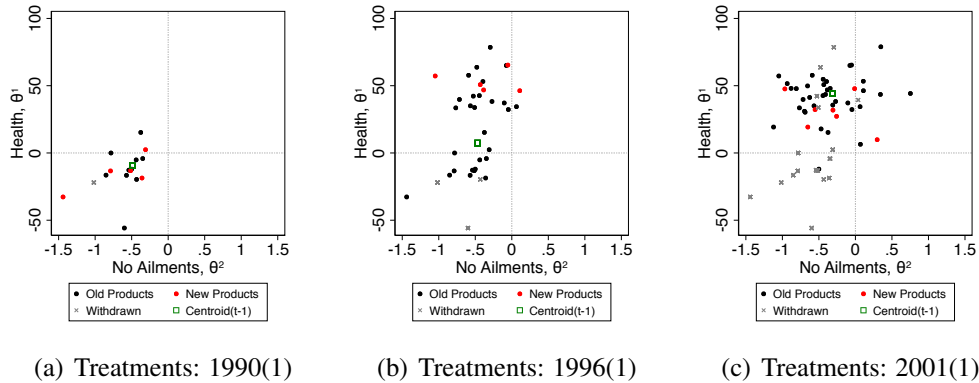
**Innovation reflects current technology and demand.** Empirical patterns until now provide evidence that demand responds to technological innovation in a market where consumers face a tradeoff between treatment effectiveness and side effects that manifest as physical ailments. Healthy consumers respond to this tradeoff by avoiding treatments with harsh side effects when treatments are not very effective. Sicker patients are more willing to suffer side effects even when treatment effectiveness is low and are willing to use experimental products in clinical trials when available alternatives are of poor quality.

Next, we provide evidence that consumers not only respond to innovation, but that innovation responds to consumer demand. We begin by illustrating the process of innovation using snapshots of the evolution of the market in Figure 5.[14] Each snapshot plots treatment characteristics (effectiveness and lack of side effects) indicating new, old, and withdrawn products as well as the lagged centroid, a summary measure of current market technology defined as the share-weighted average of product characteristics.[15] New products are introduced around the centroid sug-

---

[14]The evolution of the market is best illustrated in our animated appendix: https://www.dropbox.com/s/2icr4dxrpx9metk/treatmentevolutionNew.mp4?dl=0.

[15]We measure efficacy as the marginal contribution of a treatment to CD4 count and lack of side effects as the marginal contribution of a treatment to the log odds ratio of not causing ailments versus

13

gesting that future technologies are based on prevalent technologies today. Over time, the path of technology advances on the efficacy dimension first and then on the side effects dimension.[16]



(a) Treatments: 1990(1)     (b) Treatments: 1996(1)     (c) Treatments: 2001(1)

**FIGURE 5:** Treatment Evolution

Notes: Figure shows snapshots of the evolution of the state of the product market at the different stages. Products are two-dimensional. On the *x*-axis is a measure of a treatments ability to not cause side effects. On the *y*-axis is a measure of its contribution to underlying health. Dimensions are measured in different scales. Incumbent products are shown in black. New products are shown in red. Withdrawn products are shown as *x*. The green square is a measure of the prevalent technology in the previous period.

One possibility is that the observed path of technology is the result of supply-side technological constraints (e.g., shifts in the marginal costs of increasing treatment effectiveness relative to eliminating treatment side effects that make effectiveness innovations unprofitable). Although without cost production data we are unable to unequivocally rule out this possibility, we provide evidence against it. After HAART was invented, products came available that were more effective than earlier versions of HAART, but that had similar side effects profiles. Other products were less or equally effective compared to early versions of HAART, but had fewer side effects. Consumers preferred the latter, i.e., they demanded drugs well inside the frontier on the effectiveness dimension, choosing less effective treatments with fewer side effects. The more effective drugs with worse side effects eventually exited the market. In contrast, before HAART was invented, in an era when consumers faced low survival rates, they chose drugs on the effectiveness frontier even

---

causing ailments. See Appendix C for more details.

[16]The centroid is formally defined in the following section where we specify the structural demand model.

though less effective drugs with fewer side effects were available.

We begin to see evidence of these patterns in Figure 5. Comparing the second and third panels of the figure, notice that highly effective treatments exit the market while some treatments that are less effective, but with fewer side effects, remain. Figure S1 in Appendix A provides more direct empirical evidence that consumers demanded drugs well inside the frontier on both dimensions of treatment quality. This includes post-HAART usage of relatively ineffective treatments with fewer side effects even though more effective drugs were available, albeit with harsher side effects.[17]

These empirical patterns are not consistent with the idea that supply-side constraints alone can explain the process of technological innovation. For example, the fact that firms developed increasingly effective products after HAART was introduced means they were capable of innovating along the effectiveness dimension and that—at least in some cases—they found it cost-effective and feasible to produce such products. Consumer avoidance of these treatments, and the subsequent exit of these treatments from the market, is consistent with a role for consumer demand in driving technological innovation. The reasoning is that a profit motive drives firms to develop products with characteristics consumers prefer, and to withdraw products with characteristics that consumers avoid. The model we develop allows aggregate consumer demand to influence product entry and exit and to thereby influence the progression of technology.

## 3 A Model of Consumer Demand in a Market with Endogenous Innovation and Experimental Products

Based on the key empirical patterns discussed in the previous section, we now specify a model of demand for medical treatment. Consumers enter each period facing a new set of currently available treatments. They may also opt for an experimen-

---

[17]In particular, the figure depicts the average quality of treatments that were used versus the maximum quality available for each quality dimension (effectiveness and side effects). Before the mid 1990's (when HAART was not yet available) average efficacy consumed increased while average lack of side effects consumed remained flat, even though less toxic products were available. After the mid 1990's average lack of side effects consumed increased while average effectiveness consumed remained flat, even though more effective products were available.

tal treatment of unknown quality or to take no treatment at all. A tradeoff patients face is that medications can improve health, reduce symptoms and increase the likelihood of survival, but can also have uncomfortable side effects that reduce productivity and make work more difficult.[18] Patients have preferences over both effectiveness and side effects, which helps to explain increased willingness to suffer side effects as more effective treatments enter the market.

An important feature of our model is motivated by the evolving choice set shown in the previous section. We model new products as draws from a stochastic process, which is a function both of existing technology and of aggregate consumer behavior. Agents form beliefs over future available treatments using this stochastic process. Incorporating the evolving choice set complicates both the model and estimation because it introduces a non-stationary aggregate process into an otherwise fairly standard model of dynamic demand. There are several important benefits to this addition. First, it allows us to characterize consumer beliefs and capture consumer behavior in a market with uncertainty over future innovation without resorting to assuming perfect foresight or to viewing innovations as fully unanticipated shocks. Making these types of simplifying assumptions could lead to biases in estimated parameters and inaccurate conclusions from subsequent counterfactual policy analysis. Moreover, modeling beliefs over innovation helps to explain why patients are often observed consuming treatments well within the technological frontier, again without resorting to simplifying assumptions, such as lack of awareness that better products exist. In our model, consumers may optimally choose to delay switching to new and better treatments to avoid switching costs, in part because they anticipate even better future innovations that would make switching worthwhile.

Another benefit of incorporating the evolving choice set into the model is that it allows us to characterize a demand externality. We allow aggregate demand and the total share of consumers participating in clinical trials to affect how many new products enter the market and their characteristics (Acemoglu and Linn, 2004; Finkelstein, 2004). Agents do not take aggregate effects into account when making their individually rational choices. Hence, aggregate behavior can lead to innovation that

---

[18]Out-of-pocket payments, which we also incorporate, represent an additional cost, but they are small and vary little across treatments.

is socially sub-optimal in terms of speed or direction (i.e., which dimension of quality improves). By modeling how consumer behavior affects innovation, we are able to examine the magnitude of this type of externality and the welfare implications of policies that could help to address it.

## 3.1 Value Function

Formally, consumers inhabit a market that develops over a discrete number of periods $t \in 0, ..., T$. Let $\mathbf{P}_t$ be the set of existing commercialized market treatments available at $t$ and let $k \in \mathbb{N}$ denote a distinct market treatment. Consumer $i$ chooses between not purchasing a treatment, denoted $k = 0$, one of the market treatments available at $t$, $k \in \mathbf{P}_t$, and a treatment available in clinical trials (an experimental treatment) denoted $k = e_t$. By choosing treatment $k \in \mathbf{P}_t \cup \{0, e_t\}$ the consumer effectively chooses a bundle of treatment characteristics $\theta_k \in \mathbb{R}^2$ (effectiveness and side-effects). Let $b_{it}$ be an indicator for whether the consumer is alive at $t$ and let $z_{it}$ be the set of state variables, including the individual's characteristics (e.g., current health) along with aggregate variables summarizing the market, which we explain below. His flow utility from choosing alternative $k$ is the sum of a systematic component $u_k(z_{it})$, which depends on state variables and the characteristics of treatment $k$, and an idiosyncratic component contained in a vector $\varepsilon_{it}$ of treatment-specific *iid* Type-1 Extreme Value preference shocks. Individuals are forward-looking with discount factor $\beta \in (0, 1)$. Letting $d_{it}^e$ be the vector of optimal choices solving the consumer's maximization problem, the *ex-ante* value function for consumer $i$ at time $t$ with state $z_{it}$ can written be as:

$$V(z_{it}) \equiv E \left\{ \sum_{\tau=t}^{\infty} \sum_{k \in \mathbf{P}_\tau \cup \{0, e_\tau\}} \beta^{\tau-t} d_{ki\tau}^e b_{i\tau} \left[ u_k(z_{i\tau}) + \varepsilon_{ki\tau} \right] \middle| z_{it} \right\} \tag{1}$$

If we limit attention to individual-specific processes (e.g., preference shocks, health and survival), the value function in equation (1) looks deceivingly standard. However, embedded in equation (1) are consumer expectations over an aggregate process governing the evolution of their choice set, which we describe next.

17

## 3.2 Evolution of the Choice Set

**Market Shares, Product Quality, Entry and Exit.** Due to the complexity of the market (at least 11 pharmaceutical firms interacting with numerous government agencies and academic institutions) we do not explicitly model the suppliers' problem (e.g., the dynamic game in which producers of treatments may be engaged). Instead, we specify a stochastic process governing entry and exit of products. Entry is a function of prevailing technology along with market shares and proportion of consumers in clinical trials, which means it is endogenous to aggregate consumer behavior. This setup captures how market size and consumer experimentation can accelerate innovation. Our approach is in line with Acemoglu and Linn (2004), who relate market size to innovation, and with Finkelstein (2004) who shows that policies increasing use of certain types of products can direct innovation towards similar products. Our setup is also related to Bolton and Harris (1999), who link innovation to consumer willingness to experiment with new products. Exit is also endogenous to demand in that it occurs when few patients choose a treatment. Consumers form expectations over the process of innovation, which affects their current-period treatment choices.[19] In what follows, we describe the three objects that characterize the stochastic process: a distribution $g_\theta$ of characteristics of new and experimental treatments; a distribution $g_N$ of the number of new treatments arriving each period; and an exit rule.

Starting with $g_\theta$, the characteristics of new and experimental treatments are modeled as draws from a distribution around a *centroid*, which is a point in the characteristics space that summarizes prevailing technology. The centroid is essentially a weighted average of the characteristics of previously available treatments, where weights are market shares:

$$\omega_t \equiv \sum_{k \in \mathbf{P}_{t-1}} \tilde{s}_{kt-1} \theta_k, \qquad \tilde{s}_{kt-1} \equiv \frac{s_{kt-1}}{\sum_{k' \in \mathbf{P}_{t-1}} s_{k't-1}}. \tag{2}$$

---

[19] As we explain below, an evolving consumer choice set can also lead to a large number of products on the market. To make the problem more tractable, we reduce the size of the choice set by relying on repeat consumption and using a clustering algorithm that essentially groups similar products together. Appendix B.1 provides a more detailed explanation of the law of motion of the set of available treatments.

Here, $s_{kt}$ denotes the market share of treatment $k$ at $t$ and $s_{et}$ is the share of consumers opting for an experimental treatment through a clinical trial.[20] Newly available treatments in period $t$ are innovations around the previous-period centroid $\omega_{t-1}$. The experimental treatment available through clinical trial participation at $t$ is an innovation around the current-period centroid $\omega_t$, which includes newly available treatments. All treatments draw their characteristics $\theta_k$ from the distribution $g_\theta(\theta|w_k, s_{et-1})$ described by the following two-dimensional process:

$$\theta_k - w_k = \phi_0^v + \phi_1^v \cdot s_{et-1} + v_k, \quad E[v_k|s_{et-1}, w_k] = \mathbf{0} \tag{3}$$

where $w_k \in \{\omega_{t-1}, \omega_t\}$ depending on whether the process describes a new treatment or a current experimental treatment. The left-hand side of equation (3) denotes new product characteristics relative to the centroid, i.e., both a magnitude and a direction of innovation. Innovations are equal to a constant $\phi_0^v$ plus the previous share of individuals using the experimental treatment $s_{et-1}$ with coefficient $\phi_1^v$. This captures how consumption of experimental treatments can affect the quality of new treatments. Innovations also depend on an *iid* innovation vector $v_k$ drawn from the bivariate non parametric distribution $f_v(v)$.[21]

The number of new treatments $N_t$ entering the market changes each period. Similar to Acemoglu and Linn (2004), to capture this empirical pattern we allow the number to follow a binomial negative process $g_N(N_t|\kappa_{t-1}, s_{et-1})$ where:

$$E[N_t] = \exp(\phi_1^N \kappa_{t-1} + \phi_2^N s_{et-1}) \tag{4}$$

and $\kappa_{t-1}$, the magnitude of previous innovations, is defined as:

$$\kappa_{t-1} \equiv \sum_{r=1}^{2} \delta_r \left( \max_{\{k:\ k\in\mathbf{P}_{t-1}, k\notin\mathbf{P}_{t-2}\}} \left\{ \theta_k^r - \omega_{t-2}^r \right\} \right) \tag{5}$$

---

[20]Hence, the share of individuals not consuming a treatment at $t$ is $1 - s_{et} - \sum_{k\in\mathbf{P}_t} s_{kt}$.

[21]Consistent with this setup, we test and cannot reject the hypothesis that the coefficient on the centroid in equation (3) is equal to 1, i.e., that new product characteristics are drawn from a distribution centered on the centroid. We also find that conditioning on the centroid captures the relationship between experimental treatments at $t$ and the characteristics of new treatments entering the market at $t+1$.

given a vector $\{\delta_1, \delta_2\}$ of scaling weights ensuring comparability among the multiple characteristics of a treatment. This distribution captures two empirical patterns. First, more experimentation by product makers can be conducted if a larger proportion of the population consumes experimental treatments. Second, large breakthroughs tend to be followed by a relatively large number of new treatments; this allows for the possibility that breakthroughs spur innovative activity as firms attempt to capture market share.

Exit occurs following declining treatment usage. Decompose treatment $k$'s market share into new $\underline{s}_{kt}$ and repeat $\bar{s}_{kt}$ consumers $(\underline{s}_{kt} + \bar{s}_{kt} = s_{kt})$. The exit rule is described by the dyad $\{\underline{s}, \bar{s}\}$. When $\underline{s}_{kt}$ falls below the critical number $\underline{s}$ the treatment is no longer available for new consumers, remaining available only for repeat consumers. When $\bar{s}_{kt}$ falls below the critical number $\bar{s}$ the treatment is withdrawn altogether.

**Constructing a Tractable Choice Set.** The number of products on the market at any given time can make the choice problem computationally difficult. To address this problem we constrain the set of consumer alternatives while retaining the main features of the model. In any given period market treatments with similar characteristics are grouped into $J$ clusters using a rule denoted by $c(\mathbf{P}_t)$, which is common knowledge and uniquely assigns every treatment on the market at $t$ to a period-specific cluster.[22] Consumers choose a cluster and are randomly assigned to a treatment within the cluster. Assignment is modeled as a probability calculated using observed treatment probabilities conditional on choosing a treatment within the cluster. Let $\mathbf{P}_{jt}$ denote the treatments that rule $c$ assigns to cluster $j$ in period $t$ and let $q_{kjt}(k|\mathbf{P}_{jt})$ be the probability that treatment $k \in \mathbf{P}_{jt}$ is assigned when cluster $j$ is chosen at $t$. Hence, the distribution of characteristics induced onto the $j^{th}$ cluster at $t$ is:[23]

$$f_j(\theta|\mathbf{P}_{jt}) = \sum_{k \in \mathbf{P}_{jt}} q_{kjt}(k|\mathbf{P}_{jt}) I\{\theta_k = \theta\} \tag{6}$$

---

[22]The clustering algorithm is described in greater detail in Appendix B.1.

[23]We specify the clustering rule $c$ using an algorithm typically used in machine learning known as $k$-means and $q_{kjt}$ using a flexible polynomial based on the the characteristics of treatments in the cluster. For details see Appendix B.1.

Using notation to accommodate clusters, consumer options include not purchasing a treatment ($j = 0$), one of the clusters of market treatments available at $t$ ($j \in \{1, \ldots, J\}$), and the current experimental treatment ($j = J + 1$). Let $\theta_{jit}$ denote the treatment characteristics for individual $i$ implied by choice $j$ at $t$. If the consumer chooses one of the clusters or the experimental treatment he is in fact choosing to take a draw from one of the $J$ period-specific within-cluster distributions $f_j$ or from $g_\theta$, respectively. Finally, we allow the agent to choose the same market treatment as in the prior period ($j = J + 2$), supposing it has not exited the market. If so, the agent bypasses the probabilistic treatment assignment he faces when choosing a cluster. Let $r_{it}$ be an indicator for whether individual $i$ consumed a market treatment at $t - 1$ that is still available at $t$ and let $d_{jit}$, the $j^{th}$ component of $d_{it}$, be the indicator for individual $i$ choosing alternative $j$ at period $t$. Consumer choices satisfy $\sum_{j=0}^{J+1+r_{it}} d_{jit} = 1$.[24]

While the primary motivation to use the clustering algorithm is to reduce the size of the choice set, doing so while adding the possibility of repeat consumption has some attractive features for the context we study. First, clusters capture how patients generally face intermediaries (e.g., doctors) that recommend specific treatments based on patients' stated preferences. Yet, patients who have already used a specific treatment could presumably request to continue using it. Second, random assignment within a cluster introduces uncertainty for patients switching treatments, including the risk of assignment to a worse treatment. This uncertainty can help to explain reluctance to switch from treatments that are inside the technological frontier.

## 3.3 Health, Outcomes, Survival and Preferences

We now provide more detail on additional components of the value function, including stochastic processes for individual-specific state variables, outcomes such

---

[24]Using clusters instead of the full set of treatments yields a modified version of equation (1):

$$V(z_{it}) \equiv E \left\{ \sum_{\tau=t}^{\infty} \sum_{j=0}^{J+1+r_{i\tau}} \beta^{\tau-t} d_{ji\tau}^e b_{i\tau} \left[ u_j(z_{i\tau}) + \varepsilon_{ji\tau} \right] \Bigg| z_{it} \right\} \tag{7}$$

Further details on the modified value function are found in Appendix B.2.

as ailments and the utility function. We continue to use notation reflecting the modified choice set that includes clusters and repeated consumption. The production function for health $h_{it}$ (CD4 count) is:

$$h_{it+1} = \sum_{s=0}^{5} \gamma_s^h h_{it}^s + \sum_{j=0}^{J+1+r_{it}} d_{jit}\theta_{jit}^1 + \varepsilon_t^h. \tag{8}$$

The first term is a polynomial on prior-period health and captures non linearities in the persistence of health over time. $\theta_{jit}^1$ is the effectiveness characteristic of the treatment consumed and $\varepsilon_t^h$ is drawn from a nonparametric distribution $f_{\varepsilon^h}(\varepsilon^h)$ with $\mathbb{E}[\varepsilon_{it}^h|h_{it},\theta] = 0$.

Treatment consumption also affects a vector of additional outcomes and state variables collected into a vector $y_{it}$. The probability of not suffering physical ailments depends on previous-period health and the side-effects characteristic of the treatment consumed $\theta_{jit}^2$:

$$\Pr[y_{1it} = 0|h_{it},\theta] = \left(1 + \exp\left(\sum_{s=0}^{5} \gamma_s^x h_{it}^s + \sum_{j=0}^{J+1+r_{it}} d_{jit}\theta_{jit}^2\right)\right)^{-1} \tag{9}$$

Labor supply $y_{2it}$ is a state variable that individuals learn at the beginning of each period before making their treatment decision; its transition probability depends on the vector $x_{it}^l = [1, h_{it}, \ldots, h_{it}^4, a_{it}, y_{2it-1}]$ as follows:

$$\Pr[y_{2it} = 1|x_{it}^l] = \left(1 + \exp\left(x_{it}^l \gamma^l\right)\right)^{-1} \tag{10}$$

where the individual demographics vector $a_{it}$ contains age (in half year increments), race/ethnicity (black, Hispanic, white), and education level (high school, some college, college or more than college). Gross income $y_{3it}$ is governed by the process

$$y_{3it} = x_{it}^m \gamma^m + \eta_i + \varepsilon_{it}^m \tag{11}$$

where $x_{it}^m = [1, h_{it}, \ldots, h_{it}^7, a_{it}, y_{1it}, y_{2it}]$, $\eta_i$ captures individual-specific productivity and $\varepsilon_{it}^m$ are *iid* income shocks that the individual observes before making their treat-

ment choice.[25] Out-of-pocket expenditures for health care $y_{4it}$ are determined by

$$y_{4it} = x_{it}^o \gamma^o + \varepsilon_{it}^o \tag{12}$$

where $\varepsilon_{it}^m$ are iid $Normal(0, \sigma_o^2)$ and $x_{it}^o = [1, h_{it}, \dots, h_{it}^6, a_t, y_{1it}, y_{2it}, d_{it}]$.[26] Expenditures increase from purchasing a treatment but may also increase due to underlying health and physical ailments. Since we do not directly observe prices, equation (12) assumes two constant treatment costs, one for market treatments and one for experimental treatments.[27]Finally, the probability of being alive during the current period is determined by the vector $x_{it}^d = [1, h_{it}, \dots, h_{it}^5, a_{it}, y_{1it-1}]$ as follows:

$$\Pr(b_{it} = 1 | x_{it}^d) = \left(1 + \exp\left(x_{it}^d \gamma^d\right)\right)^{-1} \tag{13}$$

Turning to the utility function, the systematic part of flow utility is given by:

$$u_j(h_{it}, y_{it}) = \alpha_m(y_{3it} - y_{4it}) + \alpha_s y_{1it} d_{0it} + \alpha_{jh} h_{it} + \alpha_{ja}' a_{it} \tag{14}$$

Since $y_{3it} - y_{4it}$ is net income, $\alpha_m$ captures consumption utility. The utility cost of ailments while using a treatment is normalized to zero; hence, $\alpha_s$ captures how distaste for ailments ($y_{1it} = 1$) changes when individuals are not consuming a treatment ($d_{0it} = 1$). $\alpha_{jh}$ and $\alpha_{ja}'$ capture choice-specific utility associated with health, and demographics. Since the only relevant differences among clusters—both within a period and over time—are given by the distributions of within-cluster treatment characteristics $f_j(\theta | \mathbf{P}_{jt})$, we assume that $\alpha_{jh}$ and $\alpha_{ja}'$ do not vary across clusters. Nevertheless, we allow for individuals to derive different utility from consuming an experimental treatment or from repeating consumption of a market treatment. Besides affecting lifetime utility indirectly (through its impact on future health,

---

[25]We do not need to make parametric assumptions on these shocks because they enter linearly in the payoffs from choosing all alternatives and therefore do not affect choices.

[26]Out-of-pocket expenditures are censored at zero, which is why we model them separately from gross income.

[27]End-users customarily pay a standardized deductible that is a fraction of the brochure price of the drug paid by the insurance company. Median out-of-pocket drug costs are about $300 every six months for a regime of drugs that would cost the insurance company between $5,000 and $15,000.

survival, and outcomes) current health affects utility directly; in particular, $\alpha_{jh}$ captures differences in the time and psychic costs of accessing an experimental treatment by health (e.g., if doctors are more willing to suggest experimental treatments to the sickest), and it also captures how individuals may be more willing to try a new treatment from a cluster when in poor health.[28] Physical ailments affect utility directly and through earnings and survival.

## 4 Identification and Estimation

Our data span the period of time when the market began supplying treatments until the time at which the market matured; that is, $t \in \{0, \ldots, T\}$. Over this period the survey tracks a (replenished) panel of individuals, collecting their treatment consumption $d_{it}$ (including experimental treatments in clinical trials), health $h_{it}$, state and outcome variables $y_{it}$, demographic background $a_{it}$ and survival $b_{it}$. We also observe the life cycle of each market treatment $k$ including its within-sample market share $s_{kt}$, decomposed by new $\underline{s}_{kt}$ and repeat $\bar{s}_{kt}$ consumers.

### 4.1 Identification

Treatments' characteristics $\{\theta_k\}_{\forall k}$ are identified using observed patient health outcomes for given treatment choices. Specifically, treatment effectiveness $\theta^1$ and (lack of) side effects $\theta^2$ enter linearly in the processes for future health $h_{it+1}$ and physical ailments $y_{1it}$, respectively, We do not consider individual-specific treatment effects.[29] However, bias arising from selection into treatment is often due to unobserved or omitted time-varying factors such as underlying health. In estimating treatment effects, a benefit of the MACS data set is that we can include a polynomial in CD4 count, which is a continuous, objective measure of underlying health. The transition function for health along with other state transitions and outcomes and the survival probability are identified using analogous transitions and outcomes from the data set.

---

[28]We normalize $\alpha_{0h}$ and $\alpha'_{0a}$ to zero. Therefore, $\alpha_{jh}$ measures utility relative to those who do not use a treatment.

[29]Our sample is not large enough to back out treatment-specific distributions of treatment effects for the 80 plus treatments we observe.

The objects describing the entry and exit of available treatments include the distribution of characteristics of new and experimental treatments $g_\theta$, the distribution of the number of new market treatments $g_N$, the exit rule $\{\underline{s}, \overline{s}\}$, and the distribution of treatment characteristics of each period-specific market cluster $f_j\left(\theta \mid \mathbf{P}_{jt}\right)$. Both $g_\theta$ and $g_N$ are identified using the treatment panel. We observe the menu of treatments introduced from the discovery of latent demand when the pandemic starts (around 1984) to when the market has matured (around 2007). This provides us with 43 observations from the equilibrium distribution of the number of new treatments and 94 observations from the equilibrium distribution of treatment characteristics. The thresholds in the exit rule are identified as the minimum observed treatment shares for new and for all consumers prior to exit. Period-specific characteristics distributions $f_j\left(\theta \mid \mathbf{P}_{jt}\right)$ of all clusters are identified directly from observed within-cluster shares.

Finally, following much of the literature in discrete choice, we assume the discount factor $\beta$ and parameterize the choice disturbance density as Type 1 Extreme value.[30] We also assume that individuals have rational expectations about the aggregate processes generating treatments, treatment characteristics, clusters, and within-cluster distributions. Hence, identification of the preference parameters in $u_j$ follows from the general arguments of Magnac and Thesmar (2002), more specifically covered in the framework of Arcidiacono and Miller (Forthcoming).

## 4.2 Estimation

Our estimation procedure follows the steps below. More detailed explanations of some of the steps below are found in Appendix C.

1. *Treatments*. We define a single experimental treatment per period as the one used by those individuals joining a clinical trial. Treatment characteristics are estimated together with the health and ailments processes in equations (8) and (9).

---

[30]We estimated the model for values of $\beta \in \{0.8, 0.9, 0.95\}$ and found that 0.95 delivered the lowest value of the criterion function.
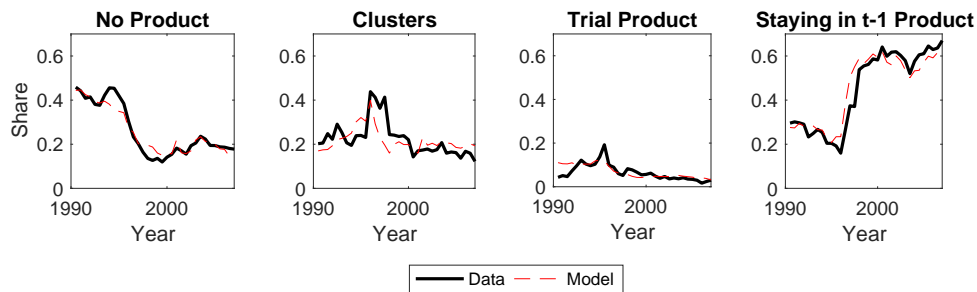
2. *Clusters*. Based on treatment characteristics from step 1, treatments are clustered at every period using the clustering rule $c$. Using the characteristics of the treatments in each cluster and treatment shares, we obtain the distribution of characteristics induced onto the $j^{th}$ cluster at $t$ given by equation (6).

3. *Innovation*. We calculate centroids for innovation and the magnitude of previous innovations for each $t$ using equations (2) and (5) and treatment characteristics from step 1. Then we estimate equation (3) that governs innovation. The residuals of this equation are used to non-parametrically estimate the two-dimensional distribution of innovation shocks $f_v$. Finally, we use the number of new treatments per period to estimate $g_N$.

4. *Transitions, outcomes and survival*. We estimate processes for labor supply transitions, income, out-of-pocket expenditures, and survival using equations (10), (11), (12), and (13).

5. *Utility function*. We estimate utility function parameters (equation (14)) using a GMM estimator and a forward simulation procedure whereby moment conditions equate the log odds ratio of current conditional choice probabilities with an expression involving utility parameters and simulated future CCPs, states and choices (Hotz et al., 1994; Altuğ and Miller, 1998).

   The forward simulation procedure we implement is modified to accommodate our context, in particular, consumer expectations over the evolving choice set. First, we estimate flexible parametric CCPs that control for the aggregate state as well as individual-specific state variables. (See Appendix C.5.) Then, for every observation $\{i,t\}$ we simulate a collection of aggregate paths describing treatment evolution using the stochastic entry and exit process; this entails simulating all individuals' behavior for each aggregate path because innovation is endogenous to aggregate patient choices. Then, for every observation $\{i,t\}$ we simulate individual choices and transitions taking as given a group of randomly selected aggregate simulated paths of innovations.

   To understand how our procedure departs from standard forward simulation, note that in other contexts where forward simulation is used there is not an

aggregate process (Hotz et al., 1994), or the aggregate process is not endogenous (Altuğ and Miller, 1998) to agent decisions. In our context, the evolution of the aggregate state, including the choice set, is endogenous to consumer choices and consumers are neither fully aware nor unaware of future choice sets. Thus, we use flexible, parametric CCPs estimated from observed consumer behavior under various aggregate states to predict treatment choices given possible counterfactual future choice sets drawn from the endogenous processes of entry and exit.

Figure 6 plots observed treatment choices over time along with those generated by the model given the state at every point in time. The estimated model captures the key trends, including the rise in repeated usage as treatments improve over time and the decline in the share of individuals not consuming a treatment. The model also captures shifts over time in the share of individuals trying something new—either by consuming experimental treatments or by choosing a cluster that entails assignment to a new treatment.
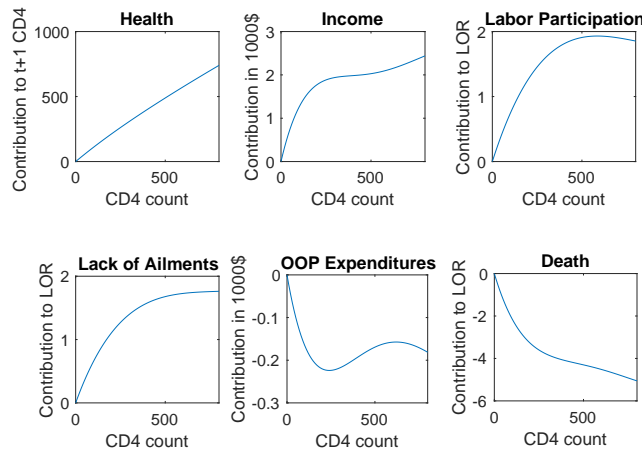
**FIGURE 6:** Goodness of Fit
Notes: Simulated and empirical choice shares over time.

## 5 Parameter Estimates

## 5.1 Individual Processes

27

**Transitions and outcomes.** Figure 7 depicts the estimated relationship between current-period health and other state variables and outcomes.[31] The relationship between current and one-period-ahead health is nearly linear, if slightly concave, reflecting that health deteriorates over time but is somewhat persistent. Persistence in health is consistent with key observed behavioral patterns, including reluctance to use medication or delays in switching to the most effective medications. Other relationships are highly non-linear, which is consistent with the well-known fact that declines in CD4 count caused by HIV infection have little-to-no impact on observed health until very low levels are reached (AIDS), at which point physical health rapidly deteriorates. In general, at CD4 counts below 250, changes in health begin to generate large shifts in ailments, income, labor supply, expenditures and survival.



**FIGURE 7:** Effect of Current Health on Future Health and Outcomes
Notes: CD4 Count measured in hundreds of cells per microliter. LOR stands for log odds ratio. OOP stands for out-of-pocket. Semestral income and expenditures measured in thousands of dollars of 2000.

Estimated equations for individual-level outcomes and transitions reveal additional patterns that accord with priors and are consistent with patterns found in other data sets (see point estimates in Tables S2 to S6 in Appendix D.1). Survival is higher for black men and for those not suffering physical ailments. Labor supply increases with age until age 40 and then flattens; it also increases with education

---

[31]We relegate the point estimates of all individual processes to Tables S2 to S6 in Appendix D.1; with very few exceptions, estimated coefficients are statistically significant at the 5% level.

and past participation in the job market. Gross income decreases with ailments, likely reflecting how poor physical health can reduce productivity, increases with employment and education, and is concave in age. Racial minorities earn less on average than white men, which is noteworthy since it means that racial inequality in labor outcomes, evident in many other data sets, extends to a population of HIV+ MSM. We also find that out-of-pocket expenditures increase with age, education, and ailments (controlling for treatment usage), perhaps due to expenditures on other health conditions. On average, minorities spend less out-of-pocket, even controlling for treatment usage. Employment increases expected expenditures, which may reflect different pricing schemes for public versus private insurance. Finally, estimates reveal that both treatments on the market and experimental treatments accessed via clinical trial entail some out-of-pocket costs.

**Preferences.** Table 2 shows that utility increases with net income, capturing utility from consumption. Utility decreases with physical ailments, reflecting consumer distaste for the symptoms of illness along with side effects of treatment (Chan and Hamilton, 2006; Papageorge, 2016).[32] The positive estimate of $\alpha_s$ implies that the cost of ailments is larger when individuals are not consuming a treatment. This suggests that the utility cost of ailments from treatment side effects is lower than the utility cost of symptoms of illness. We also find that although all individuals face utility costs from using a treatment, African Americans and Hispanics face higher costs. In particular, African Americans have the highest utility cost for consuming an experimental treatment, a finding that is consistent with a broad literature investigating historical reasons why African Americans are reluctant to participate in clinical trials (Harris et al., 1996; Alsan and Wanamaker, 2018). Age mitigates the utility costs of new treatment (commercial and experimental) which may be due to older agents becoming accustomed to trying new medications as they have more contact with the medical community. Consuming new treatments (especially experimental) is more costly for the healthy; this is consistent with less healthy indi-

---

[32]Table S7 in Appendix D.1 shows that taking the estimated ancillary parameters of the CCPs as given, both net income and physical ailments remain significant. Our final specification in equation (14) was determined by the statistical significance of results before the computationally intensive correction of standard errors in the last stage of estimation.

viduals having more frequent contact with doctors which may lower the utility cost of new treatment.[33] Finally, the utility of remaining on a treatment is positive, although not significantly different from the no-treatment base, further underscoring the individual's reluctance to try new treatments.

**TABLE 2:** Utility Parameters, $u_t$

$$u_j(h_{it}, y_{it}) = \alpha_m(y_{3it} - y_{4it}) + \alpha_s y_{1it} d_{0it} + \alpha_{jh} h_{it} + \alpha'_{ja} a_{it}$$

| coef. | variable | | est. | se |
|---|---|---|---|---|
| $\alpha_m$ | $NetIncome_t$ $(y_{3t} - y_{4t})$ | | 0.057 | (0.057) |
| $\alpha_s$ | $NoAilments_t \cdot NoTreatment_t$ $(y_{1t} d_{0t})$ | | 1.019 | (1.767) |

| | | *Cluster* $j = 1, \dots, J$ | | *Experimental* $j = J+1$ | | *Repeat* $j = J+2$ | |
|---|---|---|---|---|---|---|---|
| coef. | variable | est. | se | est. | se | est. | se |
| $\alpha_{ja1}$ | *White* | -3.546 | (0.744) | -1.468 | (0.280) | 0.502 | (0.567) |
| $\alpha_{ja2}$ | *Black* | -4.190 | (0.762) | -2.553 | (0.334) | 0.276 | (0.613) |
| $\alpha_{ja3}$ | *Hispanic* | -3.967 | (0.958) | -1.585 | (0.356) | 0.707 | (0.454) |
| $\alpha_{ja4}$ | $Age_t$ | 0.043 | (0.011) | 0.032 | (0.005) | 0.009 | (0.007) |
| $\alpha_{jh}$ | $h_t/10^3$ | -2.021 | (0.423) | -2.461 | (0.203) | | |

Notes: Estimation of (14). Discount factor $\beta = .95$. $J = 3$. $NoTreatment_{it}$ indicates whether he did not consume a treatment. $h_t$ is defined as the number of white blood cells per cubic millimeter of blood. In parentheses, standard errors computed using subsampling with 100 subsamples.

## 5.2 Product Entry and Exit

The first component of the innovation process is the distribution of characteristics of new and experimental treatments $g_\theta$.[34] Estimates of the systematic part of the innovations equation (3), which partly determines $g_\theta$, are presented in Table 3. Parameters $\phi_{11}^v$ and $\phi_{12}^v$ indicate that the previous share of the experimental treatment (i.e., participation in clinical trials) has a positive effect on average characteristics of new treatments. According to Table 3, expected effectiveness innovations are positive for lagged trial shares above 5.6 percent, and expected innovations on the ailments dimension are positive for lagged trial shares above 7.7 percent.[35]

Figure 8 depicts the estimated distribution of innovation shocks $f_v(v)$, which

---

[33]The model may also be capturing a stronger preference for sicker individuals in clinical trials.

[34]Estimates of treatment characteristics, which are obtained from the the health and ailments processes, are included in Table S8 in Appendix D.1.

[35]Mean trial participation is 7 percent in our sample (Table 1). Hence, new treatments are on average more effective than the prevalent technology but do not offer fewer side effects.

determines the stochastic portion of the innovations equation (3). Conditional on the previous trial share, the likelihood of innovation shocks decreases monotonically with their size, which is consistent with the fact that most innovations are small improvements over existing technology. Thus, an innovation like HAART is neither a fully anticipated shock nor a likely event, but instead a low probability innovation. The distribution exhibits positive correlation (0.24) between the two quality dimensions, suggesting that shocks that improve efficacy tend to be accompanied by fewer side effects. This modest correlation suggests that shocks can lead to new drugs that are improvements on both dimensions of quality. Yet, there remains a high likelihood that an innovation could mark an improvement on one dimension of quality, but a step backwards on the other dimension. Importantly, and consistent with the observed evolving set of treatments, this process can generate new products that are worse on both dimensions of quality, though endogenous market shares mean such products tend to exit the market quickly.

**TABLE 3:** Innovation Components

$$\theta_k - w_k = \phi_0^v + \phi_1^v \cdot s_{et-1} + v_k$$

| | | Health Innovation | | | Ailments Innovation | |
|---|---|---|---|---|---|---|
| | coef. | est. | se | coef. | est. | se |
| $s_{et-1}$ | $\phi_{11}^v$ | 433.11 | (19.95) | $\phi_{12}^v$ | 1.93 | (0.34) |
| Constant | $\phi_{01}^v$ | -24.14 | (1.47) | $\phi_{02}^v$ | -0.15 | (0.03) |

Notes: Estimates from (3). In parentheses, standard errors computed using subsampling with 100 subsamples.

The second component of the innovation process is the distribution of the number of new treatments $g_N$.[36] Table S9 in Appendix D.1 shows that both $\phi_1^N$ and $\phi_2^N$ in equation (4) are significant and greater than zero. Therefore, the expected number of new treatments increases with both the size of previous innovations and the previous trial share. This is consistent with firms vying for market share following breakthroughs by introducing similar treatments, and also with firms increasing their experimental activity as more consumers select experimental treatments, which may increase the quantity of viable new treatments that can be commercialized. Figure S3 in Appendix D.1 shows that the estimated distribution fits the data well.

---

[36]The third component of the law of motion of the set of available treatments are the exit rules. Our estimates of the exit thresholds $\underline{s}$ and $\bar{s}$ are 0.0047 and 0.0012, respectively.

**FIGURE 8:** The Distribution of Innovation Shocks, $f_v(v)$.
Notes: $f_v(v)$ is estimated non-parametrically off the residuals from (3).

Finally, as mentioned in Section 3.2, we ensure tractability using a clustering rule that reduces the number of available alternatives. Whenever individuals decide to try a different market treatment, they select one of $J$ period-specific clusters of treatments and are assigned a treatment from the chosen cluster according to the assignment probabilities $q_{kjt}\left(k\,|\,\mathbf{P}_{jt}\right)$ in equation (6). Point estimates of the assignment process, included in Table S10 in Appendix D.1, indicate that treatments with relatively harsher side effects within their cluster are less likely to be assigned. This makes sense as consumers dislike side effects and are thus less likely to ask for treatments with harsher side effects if similarly effective options are available that are less toxic, which would be the case within a cluster that groups together similar treatments.

## 6 The Likelihood of Observed Technological Progress

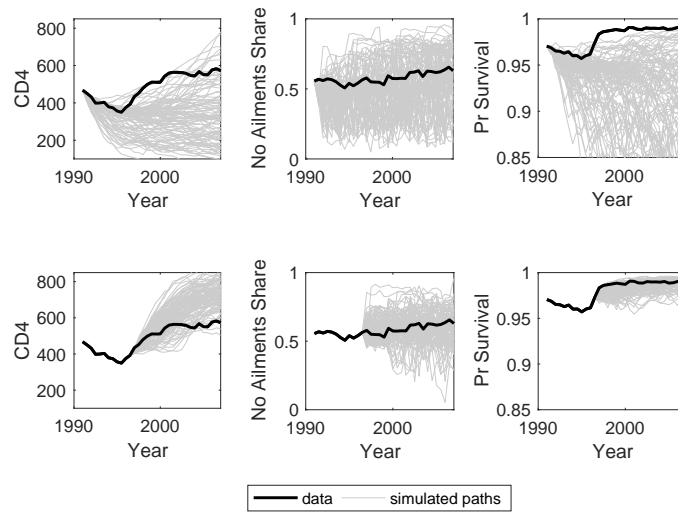Our model departs from standard models of dynamic demand in how we treat consumer expectations over innovation and the resulting evolution of the choice set. Consumers are modeled neither as fully informed (possessing perfect foresight) nor as fully uninformed (viewing each innovation as an unexpected shock). Rather, we specify a rich stochastic process governing the number of new products and their

32

qualities, which is estimated using observed innovations and is used to model consumer expectations. This is important as misspecification of consumer expectations would lead to biased estimates.

In this section we using our estimated model to assess how consumer expectations generated by the stochastic process describing innovation compare to the observed path of technology. For example, how likely or unlikely did consumers perceive a large innovation, such as HAART, to be? We simulate 100 paths of technological innovation spanning until the end of our sample (2008) and compare their distribution against the realized path observed in the data. We start our simulations at two initial states: 1991 (1st semester), prior to the introduction of breakthroughs when individuals' health was declining, and 1996 (2nd semester), shortly after the introduction of HAART when the trend in average health was reversed.

The top three panels in Figure 9 show average population CD4 count, ailments and survival as predicted by the process (grey lines) along with the realized path (black line). The first panel on the top is informative regarding the discrepancy between expectations and realized outcomes. While expectations are close to simulated paths until 1996, the introduction of HAART marks a clear, abrupt departure after which expectations over aggregate health are nearly always lower than what was observed. Two other patterns from the panel are noteworthy. First, while simulated paths tend to lie below the realized path, they also cease their descent, in part due to mortality below CD4 counts of 200 and also because the innovation process would on average generate modest innovations that would eventually lead to health improvements. Second, there are some paths that lie near or even above the realized path. Thus, an innovation like HAART was not seen as an impossible event, but instead as a low probability, large innovation. The second and third panels show ailments and survival. While survival depicts similar patterns to aggregate health, the process generates a disperse set of simulated paths that tend to lie below the realized path. This means that on average consumers expect more ailments than we observe, reflecting low expectations about aggregate health.

The bottom panels in Figure 9 repeat the exercise, but begin simulations in 1996. Shortly after HAART, consumers expected average health to continue to rise modestly, which is not what occurred. Instead, average health remained roughly

**FIGURE 9:** Distribution of Technology Paths: Individuals
Notes: 100 simulated paths conditional on the state of the world at 1991 and 1996.

constant, thus under-performing what the innovation process would have predicted. Average simulated paths of ailments are close to the realized path. Finally, simulated survival over time is concentrated near realized survival of nearly 100%. In other words, post-HAART consumers expected nearly all HIV+ patients to survive once HAART was invented, in part due to the view that further improvements would encourage universal use of HAART. In contrast, beliefs prior to HAART are much noisier and on average far below actual survival. These contrasts illustrate how rational expectations can both depart from or align with a realized path depending on the size of realized shocks. In our case, expectations, especially those prior to HAART, are not well-characterized by more standard assumptions, such as perfect foresight or that any innovation is fully unexpected.

## 7   The Externality and Policy

Since the stochastic process governing product entry and exit depends on aggregate market shares, an externality emerges because consumers do not take account of their impact on the path of innovation when making individual choices. A social planner's optimal policy incorporating the externality is a mapping from consumer

characteristics into treatment alternatives. However, solving for this mapping is intractable given the size of the state space. Hence, in this section we study temporary policy changes lasting only one period before reverting to the decentralized economy. This approach allows us to compute continuation values using the CCPs estimated in Section 4.2. As limited as these policies may appear, they have a long term impact because they affect the state variables (both aggregate and individual-specific) of the decentralized economy resuming next period.

Importantly, the validity of our counterfactual exercises requires the assumption that the three objects determining entry and exit ($g_\theta$, $g_N$, $\{\underline{s}, \bar{s}\}$) remain constant.[37] This is straightforward in a centralized economy as the planner can ensure this to be the case. However, in practice, these objects may not be immutable to counterfactual policy changes, which would expose us to the Lucas critique. Thus, we choose to investigate policies where this assumption is more likely to hold, relegating to Appendix D.3 counterfactual policies that may be more exposed to the critique.[38] Given that we approximate innovation in a reduced-form manner that does not take explicit account of firm behavior, it would be a mistake to use our model to investigate policies whose main impact would work through changes in firm behavior. Instead, we focus on one-period shifts to consumer behavior, whose main effects should work through demand, which we model in much greater detail, rather than through strategic interaction among firms.

## 7.1 Mandated Treatment

In our first policy counterfactual a planner assigns individuals to alternatives based on health and previous treatment. The population is split into four groups with high or low health and between those who are or are not potential repeat customers (i.e., those using a market treatment in the previous period or not). For each of the groups the planner either assigns one of the alternatives in the choice set to all members of the group, or assigns the individual-specific allocation from the decentralized

---

[37]Note that for distributions $g_\theta$ and $g_N$ what we require is that the functions remain constant, not the conditioning variables, which can respond to policy changes.

[38]Policies in Appendix D.3 explore how the set of available treatments would evolve over longer periods of time if consumers had less influence over the process of innovation.

economy. We solve the problem in the first semester of 1991 by computing average simulated lifetime utility under all possible allocation rules.[39]

TABLE 4: Mandated Treatment

| | | % Gain/Loss over CE | | Groups, Group Shares and Assignment | | | |
| | | | | High H No Repeat | High H Repeat | Low H No Repeat | Low H Repeat |
| | Average Welfare ($1000) | High H | Low H | 0.50 | 0.26 | 0.05 | 0.19 |
|---|---|---|---|---|---|---|---|
| | 351.61 | 2.3 | -2.8 | 0 | 6 | 6 | 0 |
| | 351.28 | 2.3 | -3.2 | 0 | 6 | 0 | 0 |
| Top Rules | 350.84 | 2.1 | -3.1 | 0 | 5 | 0 | 0 |
| | 350.82 | 2.1 | -2.9 | 0 | 5 | 6 | 0 |
| | 350.63 | 1.0 | 1.2 | 0 | 6 | 0 | 6 |
| | ⋮ | | | | | | |
| Decentralized Allocation | 346.11 | - | - | 6 | 6 | 6 | 6 |
| | ⋮ | | | | | | |
| | 167.97 | -54.2 | -40.9 | 1 | 4 | 6 | 4 |
| | 167.96 | -53.3 | -44.6 | 1 | 4 | 1 | 4 |
| Bottom Rules | 167.33 | -53.5 | -44.9 | 1 | 4 | 3 | 4 |
| | 167.24 | -53.5 | -45.1 | 3 | 4 | 2 | 4 |
| | 165.90 | -54.4 | -43.1 | 1 | 4 | 4 | 4 |

Notes: Planner's problem solved at 1991. *High H* (*Low H*) individuals have $CD4 > (\leq)250$. *Repeat* (*No Repeat*) costumers can (cannot) repeat their prior period market treatment. Population shares shown below each group label. Numbers 1 to 3 correspond to clusters and numbers 0, 4, 5, and 6 stand for no treatment, experimental treatment, repeat consumption, and the decentralized allocation, respectively.

Table 4 presents the top and bottom five assignment rules ordered by average consumer welfare. In the worst rules the planner assigns the experimental treatment to healthy patients who dislike it most and often assigns individuals to low quality clusters, meaning they also incur switching costs. In the best rules, the planner improves technology by relying on healthy potential repeat customers because their previous choices incorporate treatment quality information. Since treatment quality is low in 1991, few others are assigned treatment. The top rule increases average welfare by 1.6% but decreases equality between health groups. Relative to the decentralized allocation, healthy individuals ($CD4 > 250$) gain 2.3% in average lifetime utility while the unhealthy lose 2.8%. However, the fifth top rule increases average welfare by almost as much as the top rule, but both health groups gain. Using this rule, a social planner aiming to reduce health inequality could do so largely without sacrificing average health improvements.

## 7.2 Optimal Consumption of Experimental Treatments

Consumers choosing trials draw a treatment from the stochastic innovation distribution that is not yet available on the market. The benefit to the individual is a potentially high draw, which could be particularly attractive to ill patients facing low survival probabilities. The risk is drawing a new treatment that is not much better (or even worse) than current market treatments. An externality arises because aggregate trial participation both accelerates and influences the direction of innovation, which benefits all patients, though individuals do not take these aggregate effects into account when making individual trial choices. Our second policy counterfactual investigates this externality. This is an especially interesting policy because, although consumption of experimental treatments benefits all patients, it is the sicker individuals who tend to enter trials. Thus, we ask whether a social planner can improve average welfare, but also increase equality, by re-distributing the burden of experimentation so that it is not concentrated on the sickest patients.

As opposed to our mandated treatment exercise above, in which the planner bases policy only on health and previous treatment, the planner here assigns alternatives based on all components of the individual state. However, the planner only has two options to assign: the experimental treatment or the decentralized allocation (excluding experimental treatment). Facing a tradeoff between innovation and individual utility costs from consuming experimental treatments, the planner chooses a cutoff for consumption of experimental treatments $s_{et}^*$ such that the gain in average welfare from allocating an additional individual to the experimental treatment is no longer positive. We solve the problem at the first semester of 1991 and again at the second semester of 1996.[40] Results are presented in Table 5.

In 1991 the planner's optimal share of consumers using the experimental treatment is approximately the same as the decentralized share ($s_{et}$). The utility costs of increased consumption of experimental treatments outweigh the benefits of new drugs in a time when individuals are very sick, no good treatments have been invented and previous innovations have been small. By 1996, large innovations had

---

[40]We discretize the trial in increments of 0.005 units and simulate aggregate lifetime utility 1000 times for each value.
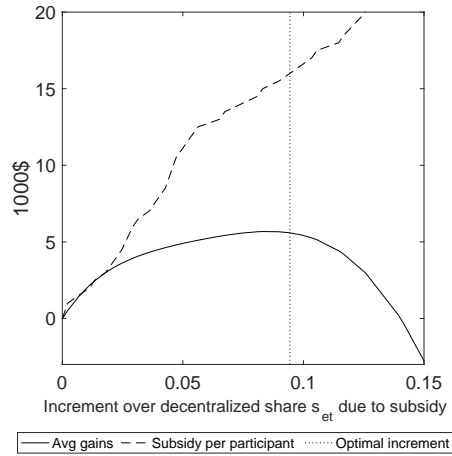
| | Year | |
|---|---|---|
| | 1991 | 1996 |
| Experimental treatment's share in planner's solution $s_{et}^*$ | 0.100 | 0.185 |
| Experimental treatment's share in decentralized economy $s_{et}$ | 0.102 | 0.092 |
| Average lifetime utility in planner's solution | 346 | 360 |
| Average lifetime utility in decentralized economy | 346 | 354 |
| Increment in share from marginal person assigned to experimental treatment at $s_{et}$ | 0.001 | 0.002 |
| Individual loss for marginal person assigned to experimental treatment at $s_{et}$ | -0.178 | -0.628 |
| Social gain from marginal person assigned to experimental treatment at $s_{et}$ | -1133 | 1051 |
| Flat subsidy (per consumer of experimental treatment) to attain $s_{et}^*$ | - | 16.0 |

Notes: Planner's problem solved at the first semester of 1991 and the second semester of 1996. Monetary values in $1,000s.

occurred, further innovations were therefore more probable, and consumers' health was improving fast. At this point in time the planner's optimal share of the experimental treatment is twice as large as the decentralized share and yields an average lifetime utility of $360,000, about 2% higher than in the decentralized economy. Figure 10 illustrates the planner's problem at 1996. The horizontal axis is trial participation above the decentralized proportion. The vertical axis is the gain in average welfare. According to the figure experimental treatment shares up to 9 percent points above the decentralized share generate welfare gains that outweigh individual losses due to consumption of experimental treatments (see the solid line).[41] Average lifetime utility drops precipitously for shares beyond the planner's optimal share $s_{et}^*$ because the new individuals assigned to the experimental treatment face larger losses relative to their optimal choice, and the innovation benefits to additional demand for experimental treatments are not large enough to offset these losses.

To measure the magnitude of the externality we also obtain the derivative of average lifetime utility with respect to the experimental treatment's share, evaluated at the decentralized share $s_{et}$. The benefit of measuring changes at the margin is that the Lucas critique is less of a concern. We assess the net social benefit of assigning the marginal consumer to the experimental treatment. Focusing on year 1996, we find that the marginal consumer loses roughly $600 (Table 5). However,

---

[41]The solid line in Figure 10 applies a fifth degree local smoothing polynomial over the original less smooth line shown in Figure S7 in Appendix D. We use this smoothed version to evaluate marginal gains.

**FIGURE 10:** Optimal Assignment to Experimental Treatments with a Flat Subsidy

Notes: On the x-axis are increments in experimental treatment's share over the decentralized share $s_{et}$. The solid line represents average gains in welfare over the decentralized allocation. The dashed line indicates the subsidy per participant necessary to decentralize a given increment. The dotted line indicates the planner's optimal increment over $s_{et}$. Year is 1996.

because demand for experimental treatments spurs innovation by raising the expected quality and the expected number of new treatments, the net social gain is over $2,000 per person. In our sample of 445 individuals in 1996, this means that a $600 loss from raising demand for experimental treatments by 1 person (about 0.22 percentage points) leads to a lifetime utility gain of roughly $1,000,000.

Although these results suggest a substantial externality associated with demand for experimental treatments, implementing this type of policy may not be feasible since it requires that the planner have a lot of information about consumer preferences. Furthermore, even with sufficient information governments with the authority to assign individuals to clinical experimentation may not do so according to who would suffer the least. For example, they might instead choose vulnerable populations to incur the costs of experimentation that benefit others. The infamous Tuskegee experiment is an example of this (Harris et al., 1996; Alsan and Wanamaker, 2018). Thus, we explore whether a flat Pigouvian subsidy could improve welfare. Figure 10 plots the flat Pigouvian subsidy necessary to attain a given share of the experimental treatment (see the dashed line) in 1996. The subsidy that attains the planner's optimal share $s_{et}^*$ is about $16,000 per participant. The size of the subsidy must be large enough to induce the marginal individual into consuming

an experimental treatment. Moreover, the subsidy represents a large reallocation of resources as it must be paid to *all* consumers of experimental treatments, including those who would have entered a trial absent subsidy. As we saw earlier, the planner's optimal share $s_{et}^*$ yields an average lifetime utility about 2% higher than in the decentralized economy. It also increases equality. As opposed to most of the mandated treatment policies described in Table 4, the subsidy decreases the gap in lifetime utility between the sickest individuals ($h_t < 200$) and everyone else by about ten percent. Additionally, equality does not increase at the expense of healthier individuals as the lump sum tax (about \$3,000) to pay for the subsidy is below their lifetime utility gains under $s_{et}^*$. Equity increases because the sickest individuals benefit the most from faster innovation and because they are more likely to consume experimental treatments. In other words, the subsidy reduces technological free-riding undertaken by healthy individuals.

## 8   Conclusion

We provide a framework to assess how consumer choices affect technological progress. In our case, aggregate consumer demand affects not only the speed of innovation, but also the direction of innovation in cases where product quality is multidimensional. We apply our framework to study consumer behavior and innovation in the market for HIV drugs. We capture several mechanisms through which consumer demand affects innovation, including experimentation with new drugs by participating in clinical trials, which accelerates the entry of new treatments. We show that individually optimal consumer behavior can slow the process of innovation due to a distaste for experimentation, and bend it towards less effective treatments that lower survival probabilities due to preferences for lower side effects. Moreover, individuals do not internalize the consequences of their treatment choices on other consumers' welfare, implying an externality that arises through the impact on technological progress. Our estimates show that a constrained planner can increase average welfare by at least two percent (approximately \$6,000 per individual), and that providing incentives for trial participation can improve social welfare.

Future research could incorporate a more structural model of supply into the type of framework we construct here. The aim would be to evaluate a richer set of counterfactual policies affecting both firm interaction and consumer behavior. For example, consider FDA policies legislating which types of medical treatments are allowed to be tested or approved. Such policies could affect innovation not only through their impact on consumer choices, but also through their impact on the behavior of strategically interacting firms responding to consumer preferences. However, incorporating supply is not a straightforward prospect and would likely require a simpler setting, e.g., one with fewer firms, less heterogeneity in consumer preferences, unidimensional product quality and less variation in the size of innovations.

## References

**Acemoglu, Daron and Joshua Linn**, "Market Size in Innovation: Theory and Evidence from the Pharmaceutical Industry," *Quarterly Journal of Economics*, 2004, *119* (23), 1049–90.

**Allcott, Hunt, Rebecca Diamond, and Jean-Pierre Dubé**, "The Geography of Poverty and Nutrition: Food Deserts and Food Choices Across the United States," NBER WP 24094 2017.

**Alsan, Marcella and Marianne Wanamaker**, "Tuskegee and the Health of Black Men," *Quarterly Journal of Economics*, 2018, *133* (1), 407–455.

**Altuğ, Sumru and Robert A. Miller**, "The Effect of Work Experience on Female Wages and Labour Supply," *Review of Economic Studies*, 1998, *65* (1), 45–85.

**Arcidiacono, Peter and Robert A. Miller**, "Identifying Dynamic Discrete Choice Models off Short Panels," *Journal of Econometrics*, Forthcoming.

**Bayer, Patrick and Robert McMillan**, "Tiebout Sorting and Neighborhood Stratification," *Journal of Public Economics*, 2012, *96* (11), 1129–1143.

**Bhaskaran, K., O. Hamouda, M. Sannes, F. Boufassa, A.M. Johnson, Lambert P.C., K. Porter, and CASCADE Collaboration**, "Changes in the Risk of Death after HIV Seroconversion Compared with Mortality in the General Population," *JAMA*, 2008, *300* (1), 51–59.

**Bhattacharya, Jay and Mikko Packalen**, "The Other Ex Ante Moral Hazard in Health," *Journal of Health Economics*, 2012, *31* (1), 135–146.

**Bolton, Patrick and Christopher Harris**, "Strategic Experimentation," *Econometrica*, 1999, *67* (2), 349–374.

**Chan, Tat Y. and Barton H. Hamilton**, "Learning, Private Information, and the Economic Evaluation of Randomized Experiments," *Journal of Political Economy*, 2006, *114* (6), 997–1040.

_ , _ , **and Nicholas W. Papageorge**, "Health, Risky Behaviour and the Value of Medical Innovation for Infectious Disease," *Review of Economic Studies*, 2016, *83* (4), 1465–1510.

**Crawford, Gregory S. and Matthew Shum**, "Uncertainty and Learning in Pharmaceutical Demand," *Econometrica*, 2005, *73* (1), 1137–1173.

**Darden, Michael**, "Smoking, Expectations, and Health: A Dynamic Stochastic Model of Lifetime Smoking Behavior," *Journal of Political Economy*, 2017, *125* (5), 1465–1522.

**Dickstein, Michael J**, "Efficient Provision of Experience Goods: Evidence from Antidepressant Choice," working paper 2018.

**Dranove, David, Craig Garthwaite, and Manuel Hermosilla**, "Pharmaceutical Profits and the Social Value of Innovation," NBER WP 20212 2014.

**Duda, Richard O. and Peter E. Hart**, *Pattern Classification and Scene Analysis*, Wiley, 1973.

**Fernandez, Jose M**, "An Empirical Model of Learning under Ambiguity: The Case of Clinical Trials," *International Economic Review*, 2013, *54* (2), 549–573.

**Finkelstein, Amy**, "Static and Dynamic Effects of Health Policy: Evidence from the Vaccine Industry," *Quarterly Journal of Economics*, 2004, *119* (2), 527–564.

**Goettler, Ronald L. and Brett Gordon**, "Does AMD Spur Intel to innovate More?," *Journal of Political Economy*, 2011, *119* (6), 1141–1200.

**Gowrisankaran, Gautam and Marc Rysman**, "Dynamics of Consumer Demand for New Durable Goods," *Journal of Political Economy*, 2012, *120* (6), 1173–1219.

**Harris, Yvonne, Philip B Gorelick, Patricia Samuels, and Isaac Bempong**, "Why African Americans May Not Be Participating in Clinical Trials.," *Journal of the National Medical Association*, 1996, *88* (10), 630.

**Hicks, John R**, *The Theory of Wages*, London: Macmillan, 1932.

**Hotz, V. Joseph and Robert A. Miller**, "Conditional Choice Probabilities and the Estimation of Dynamic Models," *Review of Economic Studies*, 1993, *60* (3), 497–529.

**_ , _ , Seth Sanders, and Jeffrey Smith**, "A Simulation Estimator for Dynamic Models of Discrete Choice," *Review of Economic Studies*, 1994, *61* (2), 265–289.

**Igami, Mitsuru**, "Estimating the Innovator's Dilemma: Structural Analysis of Creative Destruction in the Hard Disk Drive Industry, 1981-1998," *Journal of Political Economy*, 2017, *125* (3), 798–847.

**Jovanovic, Boyan and Glenn M. MacDonald**, "Competitive Diffusion," *Journal of Political Economy*, 1994, *102* (1).

**Lancaster, Kelvin J.**, "A New Approach to Consumer Theory," *Journal of Political Economy*, 1966, *74* (2), 132–157.

**Magnac, Thierry and David Thesmar**, "Identifying Dynamic Discrete Decision Processes," *Econometrica*, 2002, *70* (2), 801–816.

**Miller, Robert A.**, "Innovation and Reputation," *Journal of Political Economy*, 1988, *96* (4), 741–765.

**Papageorge, Nicholas W.**, "Why Medical Innovation is Valuable: Health, Human Capital, and the Labor Market," *Quantitative Economics*, 2016, *7* (3), 671–725.

**Petrin, Amil**, "Quantifying the Benefits of New Products: The Case of the Minivan," *Journal of Political Economy*, 2002, *110* (4), 705–729.

**Rosen, Sherwin**, "Hedonic Prices and Implicit Markets: Product Differentiation in Pure Competition," *Journal of Political Economy*, 1974, *82* (1), 34–55.

**Scherer, Frederic M**, "Demand-Pull and Technological Invention: Schmookler Revisted," *Journal of Industrial Economics*, 1982, pp. 225–237.

**Schmookler, Jacob**, *Invention and Economic Growth*, Harvard University Press, 1966.

**Stigler, George J.**, "The Cost of Subsistence," *Journal of Farm Economics*, 1945, *27* (2), 303–314.

**Waldfogel, Joel**, "Preference Externalities: An Empirical Study of Who Benefits Whom in Differentiated-Product Markets," *RAND Journal of Economics*, 2003, *34* (3), 557–568.

# For Online Publication Supplement to
# "Innovation and Diffusion of Medical Treatment"

## A   Data Appendix

Data collection for the Multi-Center AIDS Cohort Study started in 1984 with 4,954 men enrolled.[42] Two more enrollments have taken place: one in 1987-1991 (668 additional men) and another in 2001-2003 (1,350 additional men). We only use data from the first two enrollments. Since data is semi-annual each period $t$ corresponds to 6 months. Below we describe the main variables we use in our study:

*Health* ($h_{it}$): at every visit individuals undertake a physical examination that includes a blood sample which provides a measure of underlying health status: the individual's CD4 count. We denote as $h_{it}$ the CD4 count of the individual at the start of period $t$. According to the official U.S. government's website for HIV:[43]

> The CD4 count is [...] a snapshot of how well your immune system is functioning. CD4 cells (also known as CD4+ T cells) are white blood cells that fight infection. [...] These are the cells that the HIV virus kills. As HIV infection progresses, the number of these cells declines. When the CD4 count drops below 200 [cells per microliter] due to advanced HIV disease, a person is diagnosed with AIDS. A normal range for CD4 cells is about 500-1,500.

*Ailments* ($y_{1it}$): starting at visit 4, individuals are asked about physical symptoms. We focus on unusual bruises lasting at least two weeks, unintentional weight

---

[42]Data in this manuscript were collected by the Multi-Center AIDS Cohort Study with centers (Principal Investigators) at The Johns Hopkins Bloomberg School of Public Health (Joseph B. Margolick, Lisa P. Jacobson), Howard Brown Health Center, Feinberg School of Medicine, Northwestern University, and Cook County Bureau of Health Services (John P. Phair, Steven M. Wolinsky), University of California, Los Angeles (Roger Detels), and University of Pittsburgh (Charles R. Rinaldo). The MACS is funded by the National Institute of Allergy and Infectious Diseases, with additional supplemental funding from the National Cancer Institute. UO1-AI-35042, 5-MO1-RR-00052 (GCRC), UO1-AI-35043, UO1-AI-35039, UO1-AI-35040, UO1-AI-35041. Website located at http://www.statepi.jhsph.edu/macs/macs.html.

[43]See https://www.hiv.va.gov/patient/diagnosis/labs-CD4-count.asp

loss of at least 10 pounds, fatigue, diarrhea, fever, night sweats, and tender/enlarged glands. The last 5 ailments must be felt for at least 3 days during the period. Although individuals are asked explicitly about side effects starting at visit 13, we choose not to use this part of the data because it lacks consistency over time and more importantly, because individuals are most likely unable to correctly distinguish between side effects and symptoms. Thus, in our model $y_{1it}$ takes the value of 1 if an individual reports having any of the problems mentioned above.

*Labor supply* $(y_{2it})$: whether the individual worked full time (35 hours or more per week) during period $t$.

*Income* $(y_{3it})$: starting at visit 14, individuals answer the question "*Which of the following categories describes your annual individual gross income before taxes?*" For visit 14, categories are brackets that increase every $10,000, the last category being censored at "$70,000 or more." For visits 15 to 35 the brackets are censored at $50,000 and for visits 36 to 41 the brackets are censored at $60,000. We censor at $50,000 to obtain a uniform question over time. Then we assign the middle point to individuals in the bracket. For the highest bracket we assign the upper limit ($50,000). We divide gross income by two since our periods are half-years. Gross income as well as out-of-pocket expenditures (below) are in constant dollars of 2000.

*Out-of-pocket expenditures* $(y_{4it})$: starting at visit 14, individuals are asked a version of the following question "*Please, estimate the TOTAL out-of-pocket expenses that you or other personal sources (your lover, family or friends) paid for prescription medications since your last visit.*" This question is open so values are not categorized.

*Demographics* $(a_{it})$: individuals are either white, black or Hispanic, and their age increases by half a year every period.

## A.1  Products and Product Components

Starting at visit 6 individuals are asked about their medication. From visit 13 forward, as the number of treatments available increases, they answer separate survey modules for antiretroviral drugs (ARVs) and non antiretroviral drugs (NARVs). We

focus on ARVs since these are the drugs used to treat HIV infection. Below we provide the empirical definition of trial and market products that we use in the paper.

*Trial Products.* Individuals are asked to name specifically which drugs they took as well as whether or not they took the drug as part of a research study. In the original data, some of the reported drugs are themselves coded as trials. We regard these instances as individuals participating in trials. If an individual consumes one of his drugs as part of a trial we regard the individual as consuming a trial product in that period.

*Market Products.* We define a market product as a combination of components where no component is consumed in trial. This definition generates $1,835$ products. We reduce the number of market products using the following algorithm:

1. We start with the set of treatments that have more than 40 observations in the sample and denote this the set of "core market products."[44] Our core market products are listed in Table S1 which shows that there are 70 core market products overall with at most five components. Out of 20,142 subject-visit observations of individuals taking market products, 13,767 are covered by treatments classified as core market products.

2. We code the remaining 6,375 observations of non-core market products as core market products using the steps below. Each step sequentially assigns the remaining observations that were not assigned in previous steps.

   (a) Non-core market product $k$ is assigned to core market product $k'$ if $k'$ is the core market product with the highest number of components that is contained by $k$. Of the remaining 6,375 observations of non-core market products, this rule assigns 2,963 uniquely and leaves 3,412 with unassigned (1,647 that were assigned to multiple core market products plus 1,765 that were not assigned to any core market product).

   (b) If assigned to multiple core market products in step (a):

_____

[44]We tried different criteria for the minimum number of observations and product classification did not change substantially. Since our definition of core market products can miss treatments appearing near the end of the time period studied, we select the core products using all periods but exclude the last 4 periods from estimation.

i. First, we use the past history of the individual. If at period $t$ the individual is consuming non-core market product $k''$ that was assigned to both core market products $k$ and $k'$ in step (a), and he was observed consuming core market product $k$ in period $t-1$, then his treatment at $t$ is recoded as $k$. We repeat this procedure until no further gains are obtained. Out of the remaining 1,647 observations assigned to multiple core market products, 428 are assigned uniquely in this step.

ii. Second, we use the future history of the individual. If at period $t$ the individual is consuming non-core market product $k''$ that was assigned to both core market products $k$ and $k'$ in step (a), and he was observed consuming core market product $k'$ in period $t+1$, then his treatment at $t$ is recoded as $k'$. We repeat this procedure until no further gains are obtained. Out of the remaining 1,219 observations assigned to multiple core market products, 274 are assigned uniquely in this step.

iii. Third, we use the core market product with the highest share at $t$. If at period $t$ the individual is consuming non-core market product $k''$ that was assigned to both core market products $k$ and $k'$ in step (a), and $s_{kt} > s_{k't}$, then his treatment at $t$ is recoded as $k$. This final step assigns uniquely the remaining 945 observations assigned to multiple core market products.

(c) If not assigned to a core market product in step (a): we regard all 1,765 observations as "fringe treatments" since they do not contain any core market product. We aggregate all fringe treatments that appear at period $t$ into one single "fringe mix," and assign to it all users consuming this product over time. We only consider fringe mixes that have at least 40 users. This reduces the number of observations by 345 (which represents 1.6% of the number of observations of individuals using a treatment). This aggregation leads to 16 fringe mixes that we pool with the set of core market products, which amounts to a total of 86 market products overall. (See Table S1.)

A-4

3. In the paper we specified that a treatment gets withdrawn from the market altogether when its share falls below $\bar{s}$ for 2 consecutive periods. However, in the data, a treatment may have a share below $\bar{s}$ for more than 2 consecutive periods and then reappear again. 78 out of 86 core market products have unique spells without "reappearance." We regard the remaining treatments with multiple spells as measurement error and follow the next procedure to ensure that treatments have unique spells without reappearance. For every core market product $k$ with reappearance:

   (a) We identify all spells that treatment $k$ has in the data. This is, we identify the first spell and all reappearances.

   (b) From those spells we select the one that contains the period $t'$ in which $s_{kt}$ was the highest. We drop all observations of individuals consuming market product $k$ in other spells.

   Out of 19,797 (20,142 minus 345 from step 2(c)) observations of individuals taking market products, this smoothing procedure drops 42 observations leaving 19,755 observations of individuals taking market products. Supporting the importance of the spells selected by this procedure, the maximum share in the selected spell is on average about 24 times larger than the maximum share in other spells of the same market product.[45] Table S1 includes entry and exit dates implied by this spell smoothing procedure.

---

[45]In addition to this procedure we tried (i) selecting the spell with the highest average share and (ii) selecting the spell with the highest sum of shares. All criteria result in very similar entry and exit dates.
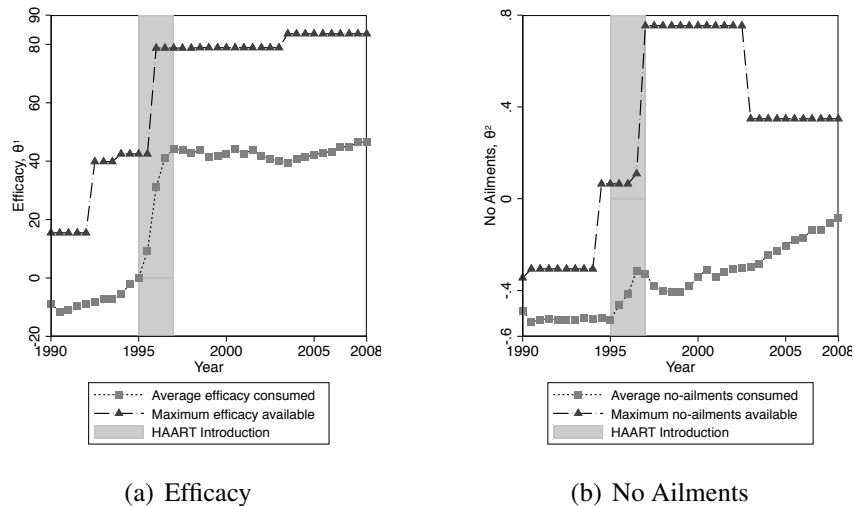
## APPENDIX TABLE S1: Market Products

| Market Product | Entry | Exit | Market Product | Entry | Exit |
|---|---|---|---|---|---|
| AZT | 1987 S1 | - | ddI , d4T, Nevirapine | 1997 S2 | - |
| Interferons ($\alpha$ and/or $\beta$), AZT | 1987 S2 | 1995 S2 | ddI , 3TC, Nelfinavir | 1997 S2 | - |
| AL-721 egg lecithin | 1987 S2 | 1991 S2 | ddI , d4T, Efavirenz | 1998 S2 | 2008 S1 |
| AZT, Acyclovir | 1989 S2 | 2000 S1 | 3TC, Abacavir, Efavirenz | 1998 S2 | - |
| Acyclovir | 1989 S2 | 2000 S1 | AZT, Nevirapine, 3TC, Abacavir | 1999 S1 | - |
| AZT, Acyclovir, ddI | 1990 S1 | 1997 S1 | AZT, 3TC, Abacavir, Efavirenz | 1999 S1 | - |
| Acyclovir, ddI | 1990 S1 | 2000 S1 | AZT, 3TC, Efavirenz | 1999 S1 | - |
| AZT, ddC | 1990 S1 | 2001 S2 | AZT, 3TC, Abacavir | 1999 S1 | - |
| AZT, ddI | 1990 S1 | 2004 S2 | d4T, 3TC, Efavirenz | 1999 S1 | 2006 S1 |
| ddI | 1990 S1 | - | Nevirapine, 3TC, Abacavir | 1999 S2 | - |
| AZT, ddC, Acyclovir, ddI | 1991 S1 | 1997 S1 | d4T, 3TC, Kaletra | 2001 S1 | 2006 S1 |
| AZT, ddC, Acyclovir | 1991 S1 | 1999 S2 | 3TC, Kaletra, Abacavir | 2001 S2 | - |
| AZT, ddC, ddI | 1991 S1 | 1995 S2 | AZT, 3TC, Kaletra | 2001 S2 | - |
| ddC, Acyclovir | 1991 S1 | 1997 S2 | AZT, 3TC, Kaletra, Abacavir | 2002 S1 | - |
| ddC | 1991 S1 | 1999 S1 | 3TC, Abacavir, Efavirenz, Tenofovir | 2002 S1 | - |
| d4T | 1993 S1 | - | AZT, 3TC, Abacavir, Tenofovir | 2002 S1 | - |
| AZT, Acyclovir, 3TC | 1994 S2 | 2000 S1 | AZT, 3TC, Kaletra, Tenofovir | 2002 S1 | - |
| AZT, 3TC | 1995 S1 | - | Nevirapine, 3TC, Tenofovir | 2002 S1 | 2007 S1 |
| Acyclovir, d4T, 3TC | 1995 S2 | 2000 S1 | 3TC, Kaletra, Tenofovir | 2002 S1 | - |
| AZT, 3TC, Saquinavir | 1996 S1 | 2005 S1 | Kaletra, Efavirenz, Tenofovir | 2002 S1 | - |
| d4T, 3TC | 1996 S1 | - | 3TC, Efavirenz, Tenofovir | 2002 S1 | - |
| AZT, 3TC, Saquinavir, Ritonavir | 1996 S2 | - | AZT, 3TC, Kaletra, Abacavir, Tenofovir | 2002 S2 | - |
| AZT, Acyclovir, 3TC, Indinavir | 1996 S2 | 2000 S1 | ddI , Kaletra, Tenofovir | 2002 S2 | - |
| Acyclovir, d4T, 3TC, Indinavir | 1996 S2 | 2000 S1 | ddI , Efavirenz, Tenofovir | 2002 S2 | - |
| AZT, 3TC, Ritonavir, Indinavir | 1996 S2 | 2006 S2 | Abacavir, Efavirenz, Tenofovir | 2002 S2 | - |
| d4T, 3TC, Ritonavir, Indinavir | 1996 S2 | 2006 S2 | Kaletra, Abacavir, Tenofovir | 2002 S2 | - |
| d4T, 3TC, Saquinavir, Ritonavir | 1996 S2 | 2004 S2 | 3TC, Ritonavir, Abacavir, Atazanavir | 2003 S2 | - |
| ddI , d4T, Indinavir | 1996 S2 | 2004 S2 | Efavirenz, Tenofovir, Emtricitabine | 2003 S2 | - |
| d4T, 3TC, Indinavir | 1996 S2 | 2008 S1 | Ritonavir, Efavirenz, Tenofovir, Emtricitabine, Atazanavir | 2004 S1 | - |
| AZT, 3TC, Indinavir | 1996 S2 | - | 3TC, Ritonavir, Abacavir, Tenofovir, Atazanavir | 2004 S1 | - |
| d4T, Nevirapine, 3TC | 1997 S1 | - | ddI , Ritonavir, Tenofovir, Atazanavir | 2004 S1 | - |
| AZT, Nevirapine, 3TC | 1997 S1 | - | Ritonavir, Tenofovir, Emtricitabine, Atazanavir | 2004 S1 | - |
| AZT, 3TC, Nelfinavir | 1997 S1 | - | Nevirapine, Tenofovir, Emtricitabine | 2004 S1 | - |
| ddI , d4T, Nelfinavir | 1997 S1 | 2005 S2 | Kaletra, Tenofovir, Emtricitabine | 2004 S2 | - |
| d4T, 3TC, Nelfinavir | 1997 S2 | - | Ritonavir, Tenofovir, Emtricitabine, Lexiva | 2005 S1 | - |

*Fringe Mixes*

| Market Product | Entry | Exit | Market Product | Entry | Exit |
|---|---|---|---|---|---|
| Isoprinosine, Ribavirin, Interferons ($\alpha$ and/or $\beta$) | 1987 S1 | 1992 S1 | Nevirapine, 3TC, Ritonavir, Kaletra, Tenofovir | 2003 S1 | - |
| Interferons ($\alpha$ and/or $\beta$), 3TC, Saquinavir, Indinavir, Efavirenz | 1997 S1 | 2007 S1 | 3TC, Ritonavir, Kaletra, Abacavir, Tenofovir, Atazanavir | 2004 S1 | - |
| Nevirapine, 3TC, Saquinavir, Ritonavir, Indinavir | 1997 S2 | 2006 S2 | Ritonavir, Tenofovir, Emtricitabine, Atazanavir, Lexiva | 2004 S2 | - |
| Nevirapine, 3TC, Saquinavir, Ritonavir, Nelfinavir | 1998 S1 | 2006 S2 | Saquinavir, Ritonavir, Tenofovir, Emtricitabine, Atazanavir | 2005 S1 | - |
| Nevirapine, Saquinavir, Ritonavir, Abacavir, Efavirenz | 1999 S1 | 2005 S2 | 3TC, Ritonavir, Abacavir, Tenofovir, Atazanavir, Lexiva | 2005 S2 | - |
| Nevirapine, Ritonavir, Nelfinavir, Abacavir, Efavirenz | 1999 S2 | - | Saquinavir, Ritonavir, Abacavir, Tenofovir, Emtricitabine | 2007 S1 | - |
| Nevirapine, Ritonavir, Kaletra, Abacavir, Efavirenz | 2001 S2 | 2008 S2 | 3TC, Ritonavir, Tenofovir, Emtricitabine, Raltegravir | 2008 S1 | - |
| Nevirapine, 3TC, Nelfinavir, Abacavir, Tenofovir | 2002 S2 | - | Ritonavir, Tenofovir, Emtricitabine, Darunavir, Raltegravir | 2008 S2 | - |

Notes: Entry and exit dates implied by the smoothing of spells in Step 3 of the algorithm used to reduce market products in Section A.1. S1 and S2 indicate the semester within a year. Many products had not exited by the end of the sample. For "fringe mixes" we only include the 5 or 6 most used products in the mix.

## A.2   Available and Prevalent Technology

Our animated appendix (alternatively, see Figure 5) shows that the path of technology advances mostly on the efficacy dimension first and then on the side effects dimension; it also shows that new products tend to appear around the centroid. Figure S1 summarizes some of the information contained in our animated appendix. It depicts the average quality demanded and the maximum quality available for each quality dimension (effectiveness and side effects). Before the mid 1990s, when people were concerned about survival, average efficacy consumed increased while average lack of side effects remained flat, even though less toxic products were available. After the mid 1990s, when people had attained higher immune system health, average lack of side effects consumed increased while average effectiveness remained flat, even though more effective products were available. Both figures reveal the supply side technological capability to develop and introduce products in both dimensions of quality over time. Thus, we argue that firms' profit motives likely led innovation to occur around what consumers were buying, generating a demand externality.



(a) Efficacy     (b) No Ailments

**APPENDIX FIGURE S1:** Available and Prevalent Technology

Notes: Left panel shows the maximum efficacy available and the average efficacy consumed by individuals. Right panel shows the same statistics for the propensity not to cause ailments.

# B Model Appendix

## B.1 Evolution of the Choice Set

In this section we provide further details of the law of motion of the set of available treatments as well as its empirical implementation.

**The distribution of the number of new treatments.** $g_N\left(N_t \mid \kappa_{t-1}, s_{et-1}\right)$ is a negative binomial that permits dispersion in the mean:

$$N_t \sim Poisson\left(\mu_{t-1}^*\right); \qquad \mu_{t-1}^* \sim Gamma\left(1/\alpha_{t-1}^N, \alpha_{t-1}^N \mu_{t-1}\right)$$

$$\mu_{t-1} = \exp(\phi_1^N \kappa_{t-1} + \phi_2^N s_{et-1}); \qquad \alpha_{t-1}^N = \exp(\phi_3^N + \phi_4^N \kappa_{t-1}) \qquad \text{(S1)}$$

where the magnitude of previous innovations $\kappa_{t-1}$ is defined in (5) and the scaling weights, which account for the fact that different characteristics may be measured in different scales, are given by the maximum innovations observed in the data:

$$\delta_r^{-1} \equiv \max_{k:\, k \in \mathbf{P}_{\tau-1},\, k \notin \mathbf{P}_{\tau-2},\, \forall \tau > 1} \left\{\theta_k^r - \omega_{\tau-2}^r\right\}, \text{ for } r \in \{1,2\} \qquad \text{(S2)}$$

**The end of a treatment's life cycle.** In the empirical implementation we relax the exit rule $\{\underline{s}, \bar{s}\}$ defined in Section 3.2 as follows. Recall that the market share of treatment $k$ can be decomposed by new $\underline{s}_{kt}$ and repeat $\bar{s}_{kt}$ consumers ($\underline{s}_{kt} + \bar{s}_{kt} = s_{kt}$). Define the conditional share for new consumers as:

$$\underline{\tilde{s}}_{kt-1} \equiv \frac{\underline{s}_{kt-1}}{\sum_{k' \in \mathbf{P}_{t-1}} \underline{s}_{k't-1}} \qquad \text{(S3)}$$

No new consumers can access treatment $k$ if $\underline{\tilde{s}}_{kt-1}$ falls below the critical number $\underline{s}$ for three consecutive periods. Treatment $k$ reaches the end of its life cycle when $\tilde{s}_{kt-1}$, defined in (2), falls below the critical number $\bar{s}$ for two consecutive periods. The number of consecutive periods for each exit rule are chosen to match the data, where a single period of low demand does not always signal the end of a treatment's life cycle. This relaxation of the exit rule adds two state variables to the aggregate state of the problem, $\mathcal{E}_{t-1}^1$ and $\mathcal{E}_{t-1}^2$, which are indicators of to what extent the

conditions for exit are binding:

$$\mathscr{E}_{kt}^1 = \mathbf{I}\left\{\tilde{\underline{s}}_{kt-1} < \underline{s}\right\}\left(\mathscr{E}_{kt-1}^1 + \mathbf{I}\left\{\tilde{\underline{s}}_{kt-1} < \underline{s}\right\}\right) \qquad \text{(S4)}$$

$$\mathscr{E}_{kt}^2 = \mathbf{I}\left\{\tilde{s}_{kt-1} < \bar{s}\right\}\left(\mathscr{E}_{kt-1}^2 + \mathbf{I}\left\{\tilde{s}_{kt-1} < \bar{s}\right\}\right) \qquad \text{(S5)}$$

where $\mathscr{E}_{kt_k}^1 = \mathscr{E}_{kt_k}^2 \equiv 0$. Exit for new consumers binds when $\mathscr{E}_{kt}^1 = 3$ and exit for all consumers binds when $\mathscr{E}_{kt}^2 = 2$.

**A Tractable Choice Set**   The clustering rule $c\left(\mathbf{P}_t\right)$, which allows us to reduce the size of the choice set, is characterized as the solution to a $k$-means clustering algorithm. At every period $t$ the clusters $j = 1, \dots, J$ are chosen to minimize:[46]

$$c\left(\mathbf{P}_t\right) = \sum_{j=1}^J \sum_{k \in \mathbf{P}_t} \mathbf{I}\left\{k \in j\right\} \left\|\theta_k - \theta_j^c\right\|^2, \qquad \theta_j^c \equiv \frac{\sum_{k \in \mathbf{P}_t} \mathbf{I}\left\{k \in j\right\} \theta_k}{\sum_{k \in \mathbf{P}_t} \mathbf{I}\left\{k \in j\right\}} \qquad \text{(S6)}$$

where $\sum_{j=1}^J \mathbf{I}\left\{k \in j\right\} = 1$ for all $k \in \mathbf{P}_t$.

The within-cluster assignment probability is given by:

$$q_{kjt}\left(k \mid \mathbf{P}_{jt}\right) = \frac{\exp\left(x_{kt}^w \gamma^w\right)}{\sum_{k \in j} \exp\left(x_{kt}^w \gamma^w\right)} \qquad \text{(S7)}$$

where $x_{kt}^w$ includes a constant term, the ranking (within its cluster) of the characteristics of the treatment, the number of members in the cluster, whether the treatment is new, and several interactions. The vector of parameters $\gamma^w$ is obtained from a nonlinear regression of within cluster shares $s_{kt|j}$ such that:

$$\mathbb{E}\left[s_{kt|j} \mid x_{kt}^w\right] = \exp\left(x_{kt}^w \gamma^w\right), \qquad s_{kt|j} \equiv \frac{s_{kt}}{\sum_{k' \in \mathbf{P}_{jt}} s_{k't}} \qquad \text{(S8)}$$

---

[46]See Duda and Hart (1973) and Andrew W. Moore's *K-means and Hierarchical Clustering* tutorial at http://www.cs.cmu.edu/~awm/tutorials.html. (See Appendix C.3 for implementation details.)

## B.2   The Modified Value Function

At the beginning of $t$ the realization of treatment assignment for those who selected a cluster in the previous period is drawn using the within-clusters probabilities $q_{kjt}$, the realization of last period's experimental treatment characteristics is drawn from $g_\theta$; health $h_{it}$, ailments $y_{1it-1}$, income $y_{3it-1}$, out-of-pocket payments $y_{4it-1}$, survival $b_{it}$, and current labor supply $y_{2it}$ are realized. The number of new treatments is drawn from $g_N$ and their characteristics drawn independently from $g_\theta$; market treatments finish their life cycle following the $\{\underline{s}, \bar{s}\}$ rule. Clusters of treatments are formed according to the clustering rule $c$. Under the law of motion specified in Section 3.2 the aggregate state $z_t$ includes the treatments remaining on the market $\mathbf{P}_t$, the centroid for innovation $\omega_t$, the magnitude of previous innovations $\kappa_t$, the previous share of the experimental treatment $s_{et-1}$, and the joint distribution of consumer demographics (including previous consumption) $\mathcal{F}_t$. The individual state is formed by idiosyncratic preference-shocks $\varepsilon_{it}$, and $z_{it}$, which includes the aggregate state $z_t$ together with a collection of individual-specific variables: health $h_{it}$, labor supply $y_{2it}$, recent usage $\theta_{J+2,it-1}$, demographics $a_{it}$ and productivity $\eta_i$. Individuals have rational expectations and zero measure in the population. They observe their current state and choose $j \in \{0, 1, \dots, J+1+r_{it}\}$. Aggregate choices at $t$ determine market shares. The individual's ex-ante value function in the modified decentralized problem is:

$$V(z_{it}) \equiv E\left\{ \sum_{\tau=t}^{\infty} \sum_{j=0}^{J+1+r_{it}} \beta^{\tau-t} d_{ji\tau}^e b_{i\tau} \left[ u_j(h_{i\tau}, y_{i\tau}) + \varepsilon_{ji\tau} \right] \,\middle|\, z_{it} \right\} \qquad (S9)$$

Because individuals in the decentralized economy do not take into account the consequences of their actions (e.g., their consumption of experimental treatments or their adoption of treatments with certain characteristics) on treatment development and hence on other individuals' future payoffs, the aggregate process generates an externality.

## C  Estimation Appendix

### C.1  Treatment Characteristics

We estimate treatment characteristics using the larger sample (visits 6 to 49) thereby using all data available on previous health, individual treatment usage, and subsequent health and ailments. Estimation equations follow from (8) and (9):

$$h_{t+1} = \sum_{s=0}^{5} \gamma_s^h h_t^s + \sum_{k \in \mathbf{P}_t} \tilde{d}_{kt} \theta_k^1 + d_{J+1,t} \theta_{et}^1 + \varepsilon_t^h \tag{S10}$$

$$\Pr[y_{1t} = 0 | h_t, \theta] = \left( 1 + \exp \left( \sum_{s=0}^{5} \gamma_s^x h_t^s + + \sum_{k \in \mathbf{P}_t} \tilde{d}_{kt} \theta_k^2 + d_{J+1,t} \theta_{et}^2 \right) \right)^{-1} \tag{S11}$$

Along with estimates of treatment characteristics, (S10) and (S11) provide parameter vectors $\gamma^h$ and $\gamma^x$ that describe the health transition in (8) and the process for physical ailments in (9).

### C.2  Exit Rule

Recall the definitions of $\underline{\tilde{s}}_{kt}$ in (S3) and $\tilde{s}_{kt}$ in (2). We set the values of the exit rule using the aggregate data on new $\underline{s}_{kt}$ and repeat $\bar{s}_{kt}$ consumers for each treatment:

$$\underline{s} = \min_{k,t} \{ \underline{\tilde{s}}_{kt} \} \quad \text{and} \quad \bar{s} = \min_{k,t} \{ \tilde{s}_{kt} \} \tag{S12}$$

### C.3  Clusters

In our empirical implementation we assume there are $J$ clusters every period. We implement the following version of the $k$-means algorithm. At every period $t$:

1. Select the treatments for which the $\underline{s}$ rule has not been applied. In other words, select treatments that are still available for new consumers at $t$. Denote this set of treatments $\mathbf{A}_t$.

2. In order to keep comparability, re-scale the characteristics of all treatments

available for clustering at $t$ by computing:

$$\tilde{\theta}_k^r = \frac{\theta_k^r}{\max_{k \in \mathbf{A}_t} |\theta_k^r|}, \text{ for } r = 1, 2 \tag{S13}$$

Thus, by construction $\tilde{\theta} \in [-1, 1] \times [-1, 1]$.

3. Select the first $J$ centroids using the scaled characteristics $\tilde{\theta}$ of $J$ randomly selected treatments from $\mathbf{A}_t$.

4. Allocate all remaining treatments $k \in \mathbf{A}_t$ to clusters sequentially. At each step select for allocation the treatment whose scaled characteristics $\tilde{\theta}_k$ are closest to one of the existing clusters. Assign treatment $k$ to the closest cluster and update the centroid of the cluster. Repeat this process until all treatments in $\mathbf{A}_t$ are assigned to a cluster.

5. Taken the centroids as given, reallocate all treatments to their closest centroid.

6. Calculate the value of the clustering rule $c(\mathbf{P}_t)$ in (S6) for the current allocation.

7. Repeat 200 times steps 3 to 6 using the scaled characteristics $\tilde{\theta}$ of different groups of $J$ randomly selected treatments in $\mathbf{A}_t$ as initial centroids. The allocation with the lowest value of $c(\mathbf{P}_t)$ is chosen.[47]

## C.4  Innovation

According to (3), the characteristics of new and experimental treatments are displaced innovations about the centroid (current or previous), and depend on previous trial participation and a draw from the distribution of innovation shocks $f_v(v)$. To estimate (3) and $f_v(v)$ we use all periods in the MACS data with relevant information on treatment consumed, health and ailments (1986 to 2008). Over the time span in our data, and given our definition of treatments, we observe 86 realized innovations from newly introduced market treatments and 22 realized innovations

---

[47]In estimation, whenever we simulate clusters we only repeat the process 50 times.

from experimental treatments. Consistent with our definition of market treatments, we only consider experimental treatments that have at least 40 users. We do not impose that innovations vectors cannot be strictly negative. In other words, relative to the centroid, inferior treatments with lower quality in both dimensions (health and ailments) may be introduced.[48]

## C.5   Utility Parameters

We estimate the utility parameters in (14) using a GMM estimator and moment conditions that equate the log odds ratio of current conditional choice probabilities with a representation of the differences in conditional value functions in terms of utility parameters and future CCPs, states and choices (Hotz et al., 1994; Altuǧ and Miller, 1998). Below we explain this step of the estimation process in more detail.

### C.5.1   Moment Condition

Our moment conditions appeal to well-known results following from our assumption that the taste shocks $\varepsilon_{jit}$ are iid Extreme Value Type I distributed (Hotz and Miller, 1993). They rely on differences between the log odds ratio and an alternative representation of differences in conditional value functions $(v_j(z_{it}) - v_0(z_{it}))$ in terms of future conditional choice probabilities, choices, states and utility parameters. Recalling the definition of $V(z_{it})$ in (S9), the conditional value function of choosing alternative $j$ at period $t$ is:

$$v_j(z_{it}) = E\left\{u_j(h_{it}, y_{it}) + \beta V(z_{it+1}) \,\middle|\, z_{it}, d_{jit} = 1\right\} \tag{S14}$$

Let $p_{jit}(z_{it})$ be the probability that individual $i$ chooses option $j$ at time $t$ conditional on his state $z_{it}$. Let $\psi_{jit}(z_{it})$ be the expected value of the $j^{th}$ taste shock conditional on alternative $j$ being optimal, and let $\gamma$ be the Euler constant. Since the joint distribution of $\varepsilon_t$ is Extreme Value Type-I:

$$\psi_j(z_{it}) \equiv E_\varepsilon\left[\varepsilon_{jit} | z_{it}, d_{jit}^e = 1\right] = \gamma - \ln\left(p_{jit}(z_{it})\right) \tag{S15}$$

---

[48]This is consistent with what we observe in the data, and theoretical reasons why this may happen have been provided in the literature (Miller, 1988).

Define $E_j\{\cdot\}$ as the expectation conditional on $d_{jit} = 1$. Dropping the individual subindex $i$ for simplicity, using (S15), we can write the conditional value function in (S14) in terms of future utility flows induced by all available alternatives, weighted by the future probabilities of those alternatives being chosen and corrected by the fact that the alternative may not be optimal. Notably, the weighted average of corrected flow payoffs of a given period must be discounted by the probability of survival up to that period conditional on today's state and choice. Letting $T^*$ be an arbitrary period with $t < T^* \leq T$, the alternative representation of the conditional value function is given by:

$$
\begin{aligned}
v_{jt}(z_t) \;=\;& E_j\left\{u_j(h_t,y_t)\,|z_t\right\} + \beta E_j\left\{V(z_{t+1},\varepsilon_{t+1})\,|z_t\right\}\\[4pt]
\;=\;& E_j\left\{u_j(h_t,y_t)\,|z_t\right\} + \beta E_j\left\{ b_{t+1} E_\varepsilon \left\{ \sum_{j'=0}^{J+1+r_{t+1}} d^e_{j't+1}\left[ u_{j'}(h_{t+1},y_{t+1}) + \psi_{j'}(z_{t+1}) \right] \right\} \Big| z_t \right\}\\
& + \beta^2 E_j\left\{ b_{t+2} V(z_{t+2},\varepsilon_{t+2})\,|z_t \right\}\\[4pt]
\;=\;& E_j\left\{u_j(h_t,y_t)\,|z_t\right\} + \beta E_j\left\{ b_{t+1} \sum_{j'=0}^{J+1+r_{t+1}} p_{j't+1}(z_{t+1})\left[ u_{j'}(h_{t+1},y_{t+1}) + \psi_{j'}(z_{t+1}) \right] \Big| z_t \right\}\\
& + \beta^2 E_j\left\{ b_{t+2} V(z_{t+2},\varepsilon_{t+2})\,|z_t \right\}\\[4pt]
\;=\;& E_j\left\{u_j(h_t,y_t)\,|z_t\right\} + \beta E_j\left\{ b_{t+1} \sum_{j'=0}^{J+1+r_{t+1}} p_{j't+1}(z_{t+1})\left[ u_{j'}(h_{t+1},y_{t+1}) + \psi_{j'}(z_{t+1}) \right] \Big| z_t \right\}\\
& + \beta^2 E_j\left\{ b_{t+1} b_{t+2} \sum_{j'=0}^{J+1+r_{t+2}} p_{j't+2}(z_{t+2})\left[ u_{j'}(h_{t+2},y_{t+2}) + \psi_{j'}(z_{t+2}) \right] \Big| z_t \right\}\\
& + \beta^3 E_j\left\{ b_{t+1} b_{t+2} V(z_{t+3},\varepsilon_{t+3})\,|z_t \right\}\\[4pt]
\;=\;& E_j\left\{u_j(h_t,y_t)\,|z_t\right\} + \sum_{\tau=1}^{T^*} \beta^\tau E_j\left\{ \left( \prod_{r=1}^{\tau} f_b(h_{it+r}) \right) \sum_{j'=0}^{J+1+r_{t+\tau}} p_{j't+\tau}(z_{t+\tau})\left[ u_{j'}(h_{t+\tau},y_{t+\tau}) + \psi_{j'}(z_{t+\tau}) \right] \Big| z_t \right\}\\
& + \beta^{T^*+1} E_j\left\{ \left( \prod_{r=1}^{T^*+1} f_b(h_{it+r}) \right) V(z_{t+T^*+1},\varepsilon_{t+T^*+1}) \Big| z_t \right\}
\end{aligned}
$$
(S16)

Let $w(z_{it})$ be a vector of instruments orthogonal to the difference between the log odds ratio and the alternative representation. Hence, we can form the following moment conditions:

$$
\mathbb{E}\left\{ w(z_{it}) \otimes \begin{bmatrix} \ln\left( \dfrac{p_{0it}(z_{it})}{p_{1it}(z_{it})} \right) + v_{1it}(z_{it}) - v_{0it}(z_{it}) \\ \vdots \\ \ln\left( \dfrac{p_{0it}(z_{it})}{p_{J+1+r_{it},it}(z_{it})} \right) + v_{J+1+r_{it},it}(z_{it}) - v_{0it}(z_{it}) \end{bmatrix} \right\} = 0.
\tag{S17}
$$

## C.5.2 Conditional Choice Probabilities

The individual's choice set $\{0, 1, \ldots, J+1+r_{it}\}$ includes the following alternatives: no treatment, one of $J$ clusters, an experimental treatment, and last-period's product (if $r_{it} = 1$). The probability that an individual chooses one of the alternatives depends on the individual and aggregate elements of his state, where the aggregate state is given by $z_t = \{\{\theta_k\}_{k \in \mathbf{P}_t}, \omega_t, \kappa_t, s_{et-1}, \mathscr{F}_t\}$. In estimation we include $\omega_t$, $\kappa_t$ and $s_{et-1}$ directly in the CCPs and characterize other components of $z_t$ as follows. The set of treatments available $\{\theta_k\}_{k \in \mathbf{P}_t}$ is characterized by the distribution of treatment characteristics of all clusters. We use the first two moments of these distributions in estimation. The distribution of consumer characteristics $\mathscr{F}_t$ is controlled for using a set of non parametric moments denoted $\tilde{\mathscr{F}}_t$.[49] Let $m_{jit}$ be the moments describing the distribution of characteristics induced by alternative $j$ for individual $i$ at period $t$, including the mean vector and the variance matrix. Effectively, $m_{jit}$ is heterogeneous across individuals only when $j = J+2$, i.e., when the individual decides to purchase the same treatment he consumed last period. Let $m_{jit} m_{jit}$ denote a vector of interactions between the elements of $m_{jit}$. Let $\tilde{x}_{it}$ and $\tilde{z}_{it}$ be subsets of the individual-specific components of the state.[50] Let $\omega_t m_{jit}$ denote a vector of interactions between the centroid and the elements of $m_{jit}$. Similarly, let $m_{jit} \tilde{z}_{it}$ be a vector of interactions between the components of $m_{jit}$ and individual-specific state components and let $\omega_t m_{jit} \tilde{z}_{it}$ be defined in a similar fashion. Our flexible CCPs are given by:

$$p_{jit} = \frac{\exp(I_{jit})}{\sum_{j'=0}^{J+1+r_{it}} \exp(I_{j'it})} \tag{S18}$$

where

$$I_{0it} \equiv 0 \tag{S19}$$

$$I_{jit} \equiv \gamma_J \tilde{x}_{it} + \beta_0 m_{jt} + \beta_1 m_{jt} m_{jt} + \beta_2 \omega_t m_{jt} + \beta_3 m_{jt} \tilde{z}_{it} + \beta_4 \omega_t m_{jt} \tilde{z}_{it} + \beta_5 m_{jt} \tilde{\mathscr{F}}_t + \beta_6 \kappa_t + \beta_7 s_{et-1}, \quad j = 1, \ldots, J \tag{S20}$$

$$I_{J+1,it} \equiv \gamma_{J+1} \tilde{x}_{it} + \beta_0 m_{J+1,t} + \beta_1 m_{J+1,t} m_{J+1,t} + \beta_3 m_{J+1,t} \tilde{z}_{it} + \beta_5 m_{J+1,t} \tilde{\mathscr{F}}_t + \beta_6 \kappa_t + \beta_7 s_{et-1} \tag{S21}$$

$$I_{J+2,it} \equiv \gamma_{J+2} \tilde{x}_{it} + \beta_0 m_{J+2,it} + \beta_1 m_{J+2,it} m_{J+2,it} + \beta_2 \omega_t m_{J+2,it} + \beta_3 m_{J+2,it} \tilde{z}_{it} + \beta_4 \omega_t m_{J+2,it} \tilde{z}_{it} + \beta_5 m_{J+2,it} \tilde{\mathscr{F}}_t + \beta_6 \kappa_t + \beta_7 s_{et-1} \tag{S22}$$

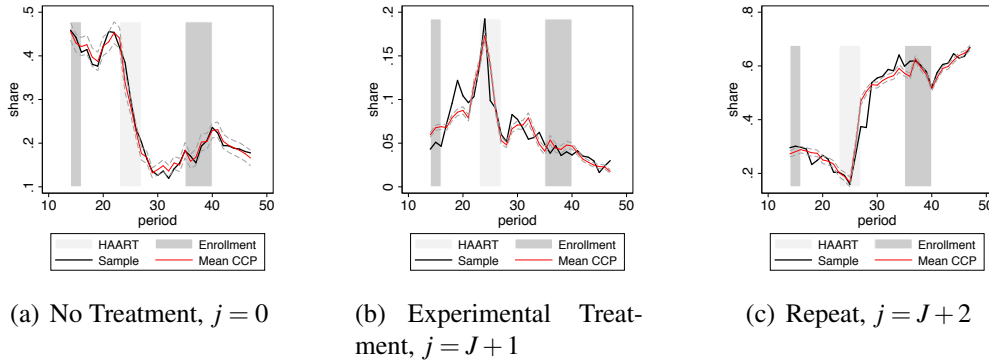Although the characteristics of the choice sets are non stationary due to treatment entry and exit, by interacting our time-varying regressors $\tilde{z}_{it}$ with the characteristics

---

[49]We specify these moments as shares of people with different sets of characteristics.

[50]$\tilde{z}_{it}$ includes $h_{it-1}$, $a_{it-1}$, $b_i$, $y_{2it}$ while $\tilde{x}_{it}$ includes a constant, $a_{it-1}$, $b_i$.

of the choice for individual $i$, $m_{jit}$, we are able to control for the state of the world inside the CCPs.[51] This procedure gives us CCPs for any simulated world as long as our observed worlds cover the space of possible worlds. Additionally, we include in the CCPs ancillary coefficients that are unrelated to the state of technology, denoted $\gamma$ in (S20) to (S22), which capture stationary taste differences between alternatives. Because, conditional on cluster characteristics, all clusters are equivalent to "trying a new market treatment," we impose $\gamma_j = \gamma_J =$ for any $j = 1, \ldots, J$.

Figure S2 displays the mean predicted conditional choice probability using (S18) over time against the corresponding share of the population who chose the alternative.[52]



(a) No Treatment, $j = 0$

(b) Experimental Treatment, $j = J + 1$

(c) Repeat, $j = J + 2$

**APPENDIX FIGURE S2:** Average CCPs

Notes: The figure shows the average estimated conditional choice probability against the share of people choosing the alternative. Dashed lines represent 95% confidence intervals around the predicted CCPs. Three periods of special relevance are highlighted in the Figure: two periods during which enrollment into the sample was undertaken and the period in which treatments belonging to the HAART class were introduced.

### C.5.3 Simulation

In order to form the sample analog of the moment condition in (S17) we obtain a simulated version of the conditional value function in (S16) truncated at $T^*$ for

---

[51]Because some of the components of $m_{J+1t}$ are linear functions of $\omega_{t-1}$ (see (3)) we avoid collinearity by not including terms $\omega_t m_{J+1,t}$ and $\omega_t m_{J+1,t} \tilde{z}_{it}$ in (S21).

[52]We also explore the fit of our CCP estimates by comparing the relative shares that clusters received in reality against the predictions from our estimated CCPs. We ranked the three clusters at every period by the share they received and compare this ranking against the ranking obtained from our estimated CCPs. Predicted ranks match observed ranks in about 80% of the periods.

every observation $\{i,t\}$ and alternative $j \in \{0,1,\ldots,J+1+r_{it}\}$. We select $T^* = 10$ so that the treatment $\beta^{T^*+1} \prod_{r=1}^{T^*+1} f_b(h_{it+r})$ approaches zero, eliminating further differences in conditional value functions beyond $T^*$. Let $S$ denote the number of simulated paths for each $\{j,i,t\}$ and let the superscript $s$ indicate that a quantity is simulated. For individual $i$ and alternative $j$ at period $t$ we write the simulated counterpart of the truncated value function as

$$\bar{v}_{jit}(z_{it}) \equiv \frac{1}{S} \sum_{s=1}^{S} \left\{ u_j(h_{it}, y_{it}^s) + \sum_{\tau=1}^{T^*} \beta^\tau \left( \prod_{r=1}^{\tau} f_b(h_{it+r}^s) \right) \sum_{j'=0}^{J+1+r_{t+\tau}} d_{j'it+\tau}^s \left[ u_{j'}(h_{it+\tau}^s, y_{it+\tau}^s) + \psi_{j'}(z_{it+\tau}^s) \right] \right\}$$
(S23)

Each future path depends on the current individual state $z_{it}$, and hence on the current aggregate state $z_t$, and the current choice $j$. We first simulate as many aggregate paths at $t$ as there are individuals at period $t$. Overall this yields $IT$ paths of technological innovation. Then, because individuals are atomistic, for each observation $\{i,t\}$ and alternative $j$ we generate sequences of future choices and payoffs taking as given $S = 20$ artificial technological paths chosen at random from the set of $I$ simulated technological paths that start at date $t$.[53] This serves two purposes. It maintains the assumption, needed for consistency of the estimator, that the sample draws from the moment conditions—the contribution from each observation—are independent from each other, and it prevents simulation errors in technology paths from propagating across all observations.

**Simulation of Aggregate State.** Taking as given the current aggregate state $z_t$ we create as many simulated aggregate state paths $\{z_{t+\tau}^s\}_{\tau=1}^{T^*}$ as there are individuals at every $t$. In other words, we repeat the algorithm below to generate $I$ simulated aggregate paths for every period $t$:

1. Let $\tau = 1$.

2. *Entry and Exit of Treatments*. Simulate a number of new treatments at $t + \tau$, $New_{t+\tau}^s$, using the entry process in (S1). If $New_{t+\tau}^s > 0$, for each simulated new treatment draw simulated characteristics using (3). Simulate the charac-

---

[53]Notice that we could rely on Hotz et al. (1994) and set $S = 1$ and obtain consistency of our estimator. However, we choose $S = 20$ after trying different values for robustness.

A-17

teristics of the experimental treatment using (3). Obtain $\kappa_{t+\tau}^s$ using (5) and (S2). For all incumbent treatments, apply the exit rule $\{\underline{s}, \overline{s}\}$ taking into account the extent to which it binds according to (S5). From the simulated set of treatments in $\mathbf{P}_{t+\tau}^s$ that have not yet satisfied the $\underline{s}$ exit rule, form clusters following the clustering rule in (S6). Obtain the distribution of characteristics of each cluster using (S7) and (S8). For $\tau > 1$ compute the simulated centroid $\omega_{t+\tau}^s$ using (2).

3. *Demand.* For all individuals $i'$ at $t$: If $\tau = 1$, define $h_{i't+1}^s \equiv h_{i't+1}$ and $d_{i't}^s \equiv d_{i't}$, otherwise, simulate $h_{i't+\tau}^s$ using (8). Draw a simulated labor state $y_{2i't+\tau}^s$ using (10). Compute deterministic transitions (e.g., age). Using $z_{i't+\tau}^s$, and hence $z_{t+\tau}^s$, and (S18) to (S22) compute simulated CCPs $p_{ji't+\tau}^s \left( z_{i't+\tau}^s \right)$ for every alternative $j \in \{0, 1, \ldots, J+1+r_{it+\tau}^s\}$ and draw a decision $d_{i't+\tau}^s$. Obtain the simulated share of trial participation $s_{e,t+\tau}^s$ and the nonparametric representation of the simulated distribution of consumer characteristics $\tilde{\mathscr{F}}_{t+\tau}^s$.

4. *Cycle back.* If $\tau = T^*$ end the loop. Otherwise, let $\tau = \tau + 1$ and go back to step 2.

**Simulation of Individual Paths.** For every observation $\{i, t\}$ and every alternative $j \in \{0, 1, \ldots, J+1+r_{it}\}$ we generate $S$ sequences of future states, choices and outcomes $\{z_{it+\tau}^s, d_{it+\tau}^s, y_{it+\tau}^s\}_{\tau=1}^{T^*}$ taking as given a subset of $S$ simulated aggregate paths—that start at $t$—chosen at random without replacement. We follow the steps below:

1. Let $\tau = 1$.

2. *Demand.* Same as above but only for individual $i$. When $j$ is not equal to the observed choice for $\{i, t\}$, we also simulate health at the beginning of period $t+1$. For this we back out the realized health residual using (S10) and use (8) to simulate health $h_{it+1}^s$ under counterfactual choice $j$. Additionally, we compute the simulated one-period-ahead survival probability $f_b \left( h_{it+\tau}^s \right)$.

3. *Outcomes.* Only for individual $i$: Simulate (lack of) ailments using (9) and the relevant distribution of treatment characteristics implied by the simulated

A-18

choice $d^s_{it+\tau}$. Simulate income using (11) and out-of-pocket expenditures using (12).[54]

4. *Cycle back.* If $\tau = T^*$ end the loop. Otherwise, let $\tau = \tau + 1$ and go back to step 2.

When simulating a path following an alternative $j$ that is not the observed choice for $\{i,t\}$, we obtain current-period simulated payoffs $u_j(h^s_{it}, y^s_{it})$ by simulating current income, out-of-pocket expenditures and ailments conditional on the counterfactual choice $j$ at $t$.

### C.5.4 Estimator

Let $j = 0$ be the base alternative, and let $\delta_{it}$ be an indicator of whether individual $i$ is in the data at period $t$. The simulated sample analog of the moment condition in (S17) is

$$
\frac{1}{\sum_i \sum_t \delta_{it}} \sum_{i=1}^{I} \sum_{t=1}^{T} \delta_{it} w(z_{it}) \otimes
\begin{bmatrix}
\ln\left(\frac{p_{0it}(z_{it})}{p_{1it}(z_{it})}\right) + \bar{v}_{1it}(z_{it}) - \bar{v}_{0it}(z_{it}) \\
\vdots \\
\ln\left(\frac{p_{0it}(z_{it})}{p_{J+1+r_{it},it}(z_{it})}\right) + \bar{v}_{J+1+r_{it},it}(z_{it}) - \bar{v}_{0it}(z_{it})
\end{bmatrix}
= 0
$$

(S24)

Denote $\Lambda$ as the $M-$dimensional vector of parameters of the utility function. Following Hotz et al. (1994) we estimate $\Lambda$ as the vector that minimizes the following objective function:

$$
\left( (IT)^{-1} \sum_{i=1}^{I} \sum_{t=1}^{T} \delta_{it} w(z_{it}) \otimes A_{it}(z_{it}, \Lambda) \right)' W_n \left( (IT)^{-1} \sum_{i=1}^{I} \sum_{t=1}^{T} \delta_{it} w(z_{it}) \otimes A_{it}(z_{it}, \Lambda) \right)
$$

(S25)

---

[54]Even though individuals know their idiosyncratic income shocks $\varepsilon^m_{it}$ we do not need to simulate these shocks as they are iid, have mean zero, and enter linearly in the flow utility, which results in them averaging out to zero in the moment condition.

$$A_{it}(z_{it},\Lambda) \equiv \begin{bmatrix} \ln\left(\frac{p_{0it}(z_{it})}{p_{1it}(z_{it})}\right) + \bar{v}_{1it}(z_{it}) - \bar{v}_{0it}(z_{it}) \\ \vdots \\ \ln\left(\frac{p_{0it}(z_{it})}{p_{J+2it}(z_{it})}\right) + \bar{v}_{J+2it}(z_{it}) - \bar{v}_{0it}(z_{it}) \end{bmatrix} \tag{S26}$$

where $W_n$ is a square weighting matrix. Using the linear structure of the utility function in (14) we collect and factor terms in order to write the $j$th component of the vector $A_{it}(z_{it},\Lambda)$ as the linear form

$$\tilde{y}_{jit} - \tilde{x}'_{jit}\Lambda \tag{S27}$$

Define $Y$ as a vector with $(J+2)IT$ rows that stacks all $\tilde{y}_{jit}$, and $X$ as a $(J+2)IT \times M$ matrix that stacks all $\tilde{x}_{jit}$. Define $Z$ as the $IT \times R$ matrix whose columns contain the $R$ instruments orthogonal to the difference between the log odds ratio of current conditional choice probabilities and the alternative representation of the differences in conditional value functions.[55] Thus

$$Y = \begin{bmatrix} \tilde{y}_{1,1,1} \\ \tilde{y}_{1,1,2} \\ \vdots \\ \tilde{y}_{1,I,T-1} \\ \tilde{y}_{1,I,T} \\ \vdots \\ \tilde{y}_{J+2,1,1} \\ \tilde{y}_{J+2,1,2} \\ \vdots \\ \tilde{y}_{J+2,I,T-1} \\ \tilde{y}_{J+2,I,T} \end{bmatrix}, \quad X = \begin{bmatrix} \tilde{x}_{1,1,1,1} & \cdots & \tilde{x}_{1,1,1,M} \\ \tilde{x}_{1,1,2,1} & \cdots & \tilde{x}_{1,1,2,M} \\ \vdots & & \vdots \\ \tilde{x}_{1,I,T-1,1} & \cdots & \tilde{x}_{1,I,T-1,M} \\ \tilde{x}_{1,I,T,1} & \cdots & \tilde{x}_{1,I,T,M} \\ \vdots & & \vdots \\ \tilde{x}_{J+2,1,1,1} & \cdots & \tilde{x}_{J+2,1,1,M} \\ \tilde{x}_{J+2,1,2,1} & \cdots & \tilde{x}_{J+2,1,2,M} \\ \vdots & & \vdots \\ \tilde{x}_{J+2,I,T-1,1} & \cdots & \tilde{x}_{J+2,I,T-1,M} \\ \tilde{x}_{J+2,I,T,1} & \cdots & \tilde{x}_{J+2,I,T,M} \end{bmatrix}, \quad Z = \begin{bmatrix} w(z_{11})_1 & \cdots & w(z_{11})_R \\ w(z_{12})_1 & \cdots & w(z_{12})_R \\ \vdots & & \vdots \\ w(z_{IT})_1 & \cdots & w(z_{IT})_R \end{bmatrix}$$

$$\tag{S28}$$

Finally, let $\mathbf{I}_{[J+2]}$ be a $(J+2)$-dimensional identity matrix and define $\tilde{Z} \equiv \mathbf{I}_{[J+2]} \otimes Z$. Then we can write the objective function in (S25) as

$$\left((IT)^{-1}\tilde{Z}'(Y - X\Lambda)\right)' W_n \left((IT)^{-1}\tilde{Z}'(Y - X\Lambda)\right) \tag{S29}$$

---

[55]Hence $W_n$ is a $(J+2)R$-dimensional square matrix.

Equation (S29) is a linear arrangement so we can obtain a closed form solution for $\hat{\Lambda}$ as the optimal GMM estimator. It entails first and second stage estimators given by

$$\hat{\Lambda}^{1S} = \left(X'\tilde{Z}\tilde{Z}'X\right)^{-1}\left(X'\tilde{Z}\tilde{Z}'Y\right), \qquad \hat{\Lambda}^{2S} = \left(X'\tilde{Z}\hat{S}^{-1}\tilde{Z}'X\right)^{-1}\left(X'\tilde{Z}\hat{S}^{-1}\tilde{Z}'Y\right) \quad \text{(S30)}$$

where

$$\hat{S} = \frac{1}{I^*}\tilde{Z}'D\tilde{Z}, \qquad I^* = IT(J+1) + \sum_{i=1}^{I}\sum_{t=1}^{T} r_{it} \qquad \text{(S31)}$$

accounts for the fact that some individuals cannot repeat their previous consumption (for instance, if the treatment was withdrawn), and $D$ is the $I(J+2)$ square diagonal matrix with diagonal elements $\hat{u}_{jit}^2 = \left(y_{jit} - x'_{jit}\hat{\Lambda}^{1S}\right)^2$. As instruments we use initial health $h_{it}$, lagged labor state $y_{2it-1}$, income fixed effect $\eta_i$, race, education indicators, and age $a_{it}$, the centroid $\omega_t$ and the lagged share of trial participation $s_{et-1}$, as well as interactions between these variables. The variance-covariance matrix of the second stage estimator is

$$\hat{V}^{2S} = I^*\left(X'\tilde{Z}\hat{S}^{-1}\tilde{Z}'X\right)^{-1} \qquad \text{(S32)}$$

## C.6 Standard Errors

The uncorrected standard errors for our utility parameters yield from the variance-covariance matrix in (S32). In order to obtain corrected standard errors we undertake subsampling taking as given the following objects obtained from the full sample: the definition of treatments (i.e., what their components are, for instance, AZT or AZT + ddI), their corresponding entry and exit dates, and the exit thresholds $\{\underline{s}, \bar{s}\}$ specified in Section 3.2. We draw $R = 100$ subsamples containing a proportion $\tilde{p} = 0.9$ of the individuals in the sample drawn without replacement, and estimate all parameters in the model using each subsample. This includes estimating treatment characteristics, parameters governing transition and outcome processes, and simulating forward paths of technology to obtain utility parameters. For any parameter $\gamma$ with estimated value $\hat{\gamma}_r$ from the $r^{th}$ subsample, the subsampling standard

A-21

errors are obtained as

$$se(\hat{\gamma}) \approx se(\hat{\gamma}_r) \cdot \sqrt{\tilde{p}} \qquad \text{(S33)}$$

where $se(\hat{\gamma}_r)$ is estimated as the standard deviation of the R quantities $\hat{\gamma}_r$.

# D Results Appendix

## D.1 Estimates

**APPENDIX TABLE S2:** Health Effects on Future Health and Ailments

| | Ailments, $\gamma^x$ | | Health, $\gamma^h$ | |
|---|---|---|---|---|
| Variables | coef. | se | coef. | se |
| $h_t$ | 0.008 | (0.0004) | 1.152 | (0.013) |
| $h_t^2/10^3$ | -0.013 | (0.001) | -0.519 | (0.043) |
| $h_t^3/10^7$ | 0.109 | (0.017) | 4.375 | (0.546) |
| $h_t^4/10^{10}$ | -0.040 | (0.010) | -2.016 | (0.298) |
| $h_t^5/10^{14}$ | 0.054 | (0.021) | 2.803 | (0.546) |
| Constant | -0.929 | (0.038) | -5.874 | (1.350) |

Notes: Parameters estimated using (S10) and (S11). In parentheses, standard errors computed using subsampling with 100 subsamples.

**APPENDIX TABLE S3:** Labor Supply, $y_{2t}$

| variable | coef. ($\gamma^l$) | se |
|---|---|---|
| $h_t$ | 0.009 | (0.0003) |
| $h_t^2/10^3$ | -0.013 | (0.001) |
| $h_t^3/10^7$ | 0.075 | (0.005) |
| $h_t^4/10^{10}$ | -0.013 | (0.002) |
| $age_t$ | 0.102 | (0.009) |
| $age_t^2$ | -0.001 | (0.0001) |
| black | -0.168 | (0.025) |
| hispanic | -0.040 | (0.044) |
| some college | 0.312 | (0.031) |
| college | 0.537 | (0.029) |
| more than college | 0.613 | (0.033) |
| labor participation$_{t-1}$, $y_{2t-1}$ | 4.458 | (0.028) |
| constant | -5.914 | (0.190) |

Notes: Estimates of the Logit model in (10). Health is given by the CD4 count measured in hundreds of cells per microliter. In parentheses, standard errors computed using subsampling with 100 subsamples.

**APPENDIX TABLE S4:** Gross Income, $y_{3t}$

| variable | coef. ($\gamma^m$) | se |
|---|---|---|
| $h_t$ | 0.018 | (0.001) |
| $h_t^2/10^3$ | -0.064 | (0.007) |
| $h_t^3/10^7$ | 1.138 | (0.171) |
| $h_t^4/10^{10}$ | -1.030 | (0.213) |
| $h_t^5/10^{14}$ | 4.854 | (1.414) |
| $h_t^6/10^{18}$ | -11.270 | (4.712) |
| $h_t^7/10^{20}$ | 0.101 | (0.062) |
| $age_t$ | 0.482 | (0.034) |
| $age_t^2$ | -0.006 | (0.0004) |
| $black$ | -5.534 | (0.115) |
| $hispanic$ | -4.167 | (0.222) |
| $some\ college$ | 2.497 | (0.141) |
| $college$ | 5.812 | (0.157) |
| $more\ than\ college$ | 8.203 | (0.151) |
| $labor\ participation_t,\ y_{2t}$ | 5.738 | (0.074) |
| $lack\ of\ ailments_t,\ y_{1t}$ | 0.207 | (0.024) |
| $constant$ | -2.095 | (0.801) |

Notes: Estimates of (11). Random effects regression of gross income on covariates. $y_{3t}$ is measured in thousands of real dollars of 2000. Health is given by the CD4 count measured in hundreds of cells per microliter. In parentheses, standard errors computed using subsampling with 100 subsamples.

**APPENDIX TABLE S5:** Out-of-pocket Expenditures, $y_{4t}$

| variable | coef. ($\gamma^o$) | se |
|---|---|---|
| $h_t$ | -0.002 | (0.0004) |
| $h_t^2/10^3$ | 0.009 | (0.002) |
| $h_t^3/10^7$ | -0.133 | (0.032) |
| $h_t^4/10^{10}$ | 0.090 | (0.029) |
| $h_t^5/10^{14}$ | -0.266 | (0.118) |
| $h_t^6/10^{18}$ | 0.279 | (0.181) |
| $age_t$ | 0.037 | (0.004) |
| $age_t^2$ | -0.0002 | (0.0001) |
| $black$ | -0.240 | (0.014) |
| $hispanic$ | -0.119 | (0.016) |
| $some\ college$ | 0.169 | (0.016) |
| $college$ | 0.318 | (0.018) |
| $more\ than\ college$ | 0.336 | (0.018) |
| $market\ product_t$ | 0.429 | (0.016) |
| $trial\ product_t$ | 0.313 | (0.021) |
| $labor\ participation_t,\ y_{2t}$ | 0.105 | (0.009) |
| $lack\ of\ ailments_t,\ y_{1t}$ | -0.122 | (0.008) |
| $constant$ | -1.459 | (0.099) |
| | | |
| $\sigma^o$ | 0.862 | (0.027) |

Notes: Estimates of (12) using a Tobit Model for data censored at 0. $market\ treatment_t = d_{J+2,t} + \sum_{k=1}^{J} d_{kt}$. Out-of-pocket expenditures $y_{4t}$ are measured in thousands of real dollars of 2000. Health is given by the CD4 count measured in hundreds of cells per microliter. In parentheses, standard errors computed using subsampling with 100 subsamples.

**APPENDIX TABLE S6:** Death, $1 - b_t$

| variable | coef. ($\gamma^d$) | se |
|---|---|---|
| $h_t$ | -0.028 | (0.001) |
| $h_t^2/10^3$ | 0.079 | (0.005) |
| $h_t^3/10^7$ | -1.104 | (0.102) |
| $h_t^4/10^{10}$ | 0.704 | (0.088) |
| $h_t^5/10^{14}$ | -1.610 | (0.285) |
| $age_t$ | -0.116 | (0.021) |
| $age_t^2$ | 0.002 | (0.0002) |
| black | -0.509 | (0.069) |
| hispanic | 0.034 | (0.076) |
| some college | 0.060 | (0.057) |
| college | -0.353 | (0.053) |
| more than college | -0.512 | (0.060) |
| lack of ailments$_{t-1}$, $y_{1t-1}$ | -1.140 | (0.050) |
| constant | 1.682 | (0.474) |

Notes: Estimates of the Logit model in (13). Health is given by the CD4 count measured in hundreds of cells per microliter. In parentheses, standard errors computed using subsampling with 100 subsamples.

**APPENDIX TABLE S7:** Utility Parameters, $y_{it}$

| coef. | variable | est. | se | unc. se |
|---|---|---|---|---|
| $\alpha_m$ | $NetIncome_t$ $(y_{3t} - y_{4t})$ | 0.057 | (0.057) | (0.010) |
| $\alpha_s$ | $NoAilments_t \cdot NoTreatment_t$ $(y_{1t}d_{0t})$ | 1.019 | (1.767) | (0.260) |

| | | Cluster $j=1,\ldots,J$ | | | Experimental $j=J+1$ | | | Repeat $j=J+2$ | | |
|---|---|---|---|---|---|---|---|---|---|---|
| coef. | variable | est. | se | unc. se | est. | se | unc. se | est. | se | unc. se |
| $\alpha_{jw}$ | White | -3.546 | (0.744) | (0.179) | -1.468 | (0.280) | (0.136) | 0.502 | (0.567) | (0.130) |
| $\alpha_{jb}$ | Black | -4.190 | (0.762) | (0.190) | -2.553 | (0.334) | (0.142) | 0.276 | (0.613) | (0.145) |
| $\alpha_{jl}$ | Hispanic | -3.967 | (0.958) | (0.647) | -1.585 | (0.356) | (0.300) | 0.707 | (0.454) | (0.354) |
| $\alpha_{ja}$ | $Age_t$ | 0.043 | (0.011) | (0.004) | 0.032 | (0.005) | (0.003) | 0.009 | (0.007) | (0.002) |
| $\alpha_{jh}$ | $h_t/10^3$ | -2.021 | (0.423) | (0.104) | -2.461 | (0.203) | (0.078) | | | |

Notes: Estimates of (14). Discount factor $\beta = .95$. $J = 3$. $NoTreatment_{it}$ indicates whether he did not consume a treatment. $h_t$ is defined as the number of white blood cells per cubic millimeter of blood. In parentheses, uncorrected standard errors (unc. se) computed using (S32), and corrected standard errors (se) computed using subsampling with 100 subsamples.

# APPENDIX TABLE S8: Treatment Characteristics

| Market Product | Ailments, $\theta^2$ coeff | se | Health, $\theta^1$ coeff | se |
|---|---|---|---|---|
| AZT | -0.500 | (0.020) | -12.004 | (0.736) |
| Interferons ($\alpha$ and/or $\beta$), AZT | -0.600 | (0.061) | -55.796 | (3.102) |
| AL-721 egg lecithin | -0.433 | (0.087) | -19.655 | (3.917) |
| AZT, Acyclovir | -0.539 | (0.050) | -12.752 | (1.670) |
| Acyclovir | -0.783 | (0.047) | -0.017 | (2.678) |
| AZT, Acyclovir, ddI | -0.851 | (0.037) | -16.474 | (1.497) |
| Acyclovir, ddI | -0.348 | (0.043) | -4.159 | (2.479) |
| AZT, ddC | -0.439 | (0.029) | -5.155 | (1.309) |
| AZT, ddI | -0.571 | (0.061) | -16.615 | (2.488) |
| ddI | -0.375 | (0.071) | 15.263 | (2.587) |
| AZT, ddC, Acyclovir, ddI | -0.789 | (0.115) | -13.351 | (7.73) |
| AZT, ddC, Acyclovir | -0.514 | (0.086) | -13.186 | (2.168) |
| AZT, ddC, ddI | -1.440 | (0.047) | -32.700 | (1.801) |
| ddC, Acyclovir | -0.310 | (0.093) | 2.415 | (4.370) |
| ddC | -0.358 | (0.084) | -18.630 | (3.389) |
| d4T | -0.717 | (0.054) | 39.776 | (2.210) |
| AZT, Acyclovir, 3TC | -0.527 | (0.096) | 42.267 | (3.394) |
| AZT, 3TC | 0.064 | (0.051) | 34.398 | (1.875) |
| Acyclovir, d4T, 3TC | -0.509 | (0.100) | 33.792 | (4.664) |
| AZT, 3TC, Saquinavir | -0.271 | (0.052) | 38.283 | (1.992) |
| d4T, 3TC | -0.104 | (0.112) | 37.173 | (4.070) |
| AZT, 3TC, Saquinavir, Ritonavir | -0.591 | (0.085) | 57.776 | (10.571) |
| AZT, Acyclovir, 3TC, Indinavir | -0.479 | (0.056) | 63.734 | (2.201) |
| Acyclovir, d4T, 3TC, Indinavir | -0.295 | (0.108) | 78.559 | (3.665) |
| AZT, 3TC, Ritonavir, Indinavir | -0.567 | (0.102) | 35.032 | (6.629) |
| d4T, 3TC, Ritonavir, Indinavir | -0.767 | (0.049) | 33.510 | (3.321) |
| d4T, 3TC, Saquinavir, Ritonavir | -0.444 | (0.085) | 42.631 | (5.409) |
| ddI , d4T, Indinavir | -0.048 | (0.137) | 32.286 | (3.981) |
| d4T, 3TC, Indinavir | -0.395 | (0.096) | 53.128 | (4.546) |
| AZT, 3TC, Indinavir | -0.075 | (0.066) | 65.041 | (2.809) |
| d4T, Nevirapine, 3TC | -0.386 | (0.052) | 46.846 | (2.962) |
| AZT, Nevirapine, 3TC | 0.109 | (0.087) | 46.275 | (4.061) |
| AZT, 3TC, Nelfinavir | -0.432 | (0.072) | 50.776 | (3.924) |
| ddI , d4T, Nelfinavir | -1.049 | (0.060) | 57.227 | (3.672) |
| d4T, 3TC, Nelfinavir | -0.881 | (0.134) | 48.018 | (9.588) |

| Market Product | Ailments, $\theta^2$ coeff | se | Health, $\theta^1$ coeff | se |
|---|---|---|---|---|
| ddI , d4T, Nevirapine | 0.753 | (0.175) | 44.240 | (3.781) |
| ddI , 3TC, Nelfinavir | -0.810 | (0.083) | 47.816 | (6.848) |
| ddI , d4T, Efavirenz | -0.626 | (0.078) | 41.280 | (2.772) |
| 3TC, Abacavir, Efavirenz | 0.108 | (0.047) | 53.341 | (1.501) |
| AZT, Nevirapine, 3TC, Abacavir | 0.038 | (0.131) | 39.379 | (3.369) |
| AZT, 3TC, Abacavir, Efavirenz | 0.348 | (0.080) | 78.914 | (3.549) |
| AZT, 3TC, Efavirenz | 0.342 | (0.079) | 43.526 | (3.073) |
| AZT, 3TC, Abacavir | -0.442 | (0.078) | 54.824 | (3.175) |
| d4T, 3TC, Efavirenz | -0.346 | (0.069) | 47.978 | (3.876) |
| Nevirapine, 3TC, Abacavir | -0.470 | (0.099) | 17.866 | (12.148) |
| d4T, 3TC, Kaletra | -0.310 | (0.123) | 35.611 | (5.199) |
| 3TC, Kaletra, Abacavir | -0.934 | (0.124) | 51.570 | (5.325) |
| AZT, 3TC, Kaletra | -0.655 | (0.140) | 49.838 | (3.967) |
| AZT, 3TC, Kaletra, Abacavir | 0.298 | (0.234) | 9.855 | (9.404) |
| 3TC, Abacavir, Efavirenz, Tenofovir | -0.308 | (0.070) | 31.845 | (3.848) |
| AZT, 3TC, Abacavir, Tenofovir | -0.652 | (0.074) | 19.273 | (5.651) |
| AZT, 3TC, Kaletra, Tenofovir | -0.552 | (0.067) | 32.227 | (2.681) |
| Nevirapine, 3TC, Tenofovir | -0.258 | (0.163) | 27.246 | (4.619) |
| 3TC, Kaletra, Tenofovir | -0.092 | (0.082) | 51.672 | (2.709) |
| Kaletra, Efavirenz, Tenofovir | -0.966 | (0.100) | 47.617 | (2.684) |
| 3TC, Efavirenz, Tenofovir | -0.011 | (0.108) | 47.790 | (5.468) |
| AZT, 3TC, Kaletra, Abacavir, Tenofovir | -0.738 | (0.141) | 19.980 | (4.226) |
| ddI , Kaletra, Tenofovir | -0.276 | (0.112) | 18.396 | (4.015) |
| ddI , Efavirenz, Tenofovir | -0.420 | (0.117) | 2.381 | (2.505) |
| Abacavir, Efavirenz, Tenofovir | -0.762 | (0.140) | 39.457 | (3.150) |
| Kaletra, Abacavir, Tenofovir | -0.820 | (0.198) | 14.891 | (2.601) |
| 3TC, Ritonavir, Abacavir, Atazanavir | -0.061 | (0.039) | 26.850 | (1.181) |
| Efavirenz, Tenofovir, Emtricitabine | 0.118 | (0.082) | 54.798 | (2.464) |
| Ritonavir, Efavirenz, Tenofovir, Emtricitabine, Atazanavir | 0.306 | (0.053) | 83.823 | (1.706) |
| 3TC, Ritonavir, Abacavir, Tenofovir, Atazanavir | -0.403 | (0.163) | 38.313 | (10.521) |
| ddI , Ritonavir, Tenofovir, Atazanavir | 0.049 | (0.108) | 47.800 | (2.837) |
| Ritonavir, Tenofovir, Emtricitabine, Atazanavir | 0.138 | (0.104) | 53.028 | (3.940) |
| Nevirapine, Tenofovir, Emtricitabine | -0.205 | (0.079) | 37.227 | (2.303) |
| Kaletra, Tenofovir, Emtricitabine | -0.183 | (0.093) | 46.723 | (5.990) |
| Ritonavir, Tenofovir, Emtricitabine, Lexiva | -0.372 | (0.116) | 30.226 | (3.328) |

*Fringe Mixes*

| Market Product | Ailments, $\theta^2$ coeff | se | Health, $\theta^1$ coeff | se |
|---|---|---|---|---|
| Isoprinosine, Ribavirin, Interferons ($\alpha$ and/or $\beta$) | -1.017 | (0.110) | -21.950 | (6.644) |
| Interferons ($\alpha$ and/or $\beta$), 3TC, Saquinavir, Indinavir, Efavirenz | -0.054 | (0.243) | 65.353 | (5.179) |
| Nevirapine, 3TC, Saquinavir, Ritonavir, Indinavir | 0.068 | (0.134) | 6.457 | (7.335) |
| Nevirapine, 3TC, Saquinavir, Ritonavir, Nelfinavir | -0.689 | (0.156) | 30.293 | (7.841) |
| Nevirapine, Saquinavir, Ritonavir, Abacavir, Efavirenz | -1.121 | (0.161) | 19.278 | (4.112) |
| Nevirapine, Ritonavir, Nelfinavir, Abacavir, Efavirenz | -0.697 | (0.099) | 31.044 | (4.027) |
| Nevirapine, Ritonavir, Kaletra, Abacavir, Efavirenz | -0.410 | (0.174) | 43.495 | (5.757) |
| Nevirapine, 3TC, Nelfinavir, Abacavir, Tenofovir | -0.467 | (0.109) | 27.893 | (3.250) |

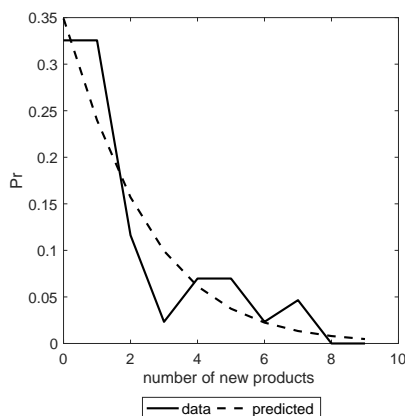| Market Product | Ailments, $\theta^2$ coeff | se | Health, $\theta^1$ coeff | se |
|---|---|---|---|---|
| Nevirapine, 3TC, Ritonavir, Kaletra, Tenofovir | -1.265 | (0.113) | 45.683 | (4.934) |
| 3TC, Ritonavir, Kaletra, Abacavir, Tenofovir, Atazanavir | -0.465 | (0.077) | 28.440 | (2.687) |
| Ritonavir, Tenofovir, Emtricitabine, Atazanavir, Lexiva | -0.612 | (0.142) | 42.050 | (3.579) |
| Saquinavir, Ritonavir, Tenofovir, Emtricitabine, Atazanavir | -0.665 | (0.120) | 31.824 | (3.879) |
| 3TC, Ritonavir, Abacavir, Tenofovir, Atazanavir, Lexiva | -0.210 | (0.078) | 26.678 | (5.890) |
| Saquinavir, Ritonavir, Abacavir, Tenofovir, Emtricitabine | 0.072 | (0.142) | 32.865 | (4.856) |
| 3TC, Ritonavir, Tenofovir, Emtricitabine, Raltegravir | 0.032 | (0.094) | 33.352 | (2.728) |
| Ritonavir, Tenofovir, Emtricitabine, Darunavir, Raltegravir | -0.221 | (0.067) | 47.736 | (2.929) |

Notes: Treatment characteristics are estimated as indicators for treatment usage in (S10) and (S11). In parentheses, standard errors computed using subsampling with 100 subsamples. For "fringe Mixes" we only include the 5 or 6 most used treatments in the mix.

**APPENDIX TABLE S9:** Distribution of Number of New Treatments, $F_N$

$$E[N_t] = \mu_{t-1} \equiv \exp(\phi_1^N \kappa_{t-1} + \phi_2^N s_{et-1})$$

| | $\ln \mu$ | | | | $\ln \alpha$ | | |
|---|---|---|---|---|---|---|---|
| variable | coef. | est. | se | variable | coef. | est. | se |
| $\kappa_{t-1}$ | $\phi_1^N$ | 0.432 | (0.246) | *Constant* | $\phi_3^N$ | -0.206 | (0.451) |
| $s_{et-1}$ | $\phi_2^N$ | 6.177 | (2.462) | $\kappa_{t-1}$ | $\phi_4^N$ | -1.019 | (0.626) |

Notes: Model is specified in (S1). $\kappa_{t-1}$ measures the magnitude of previous innovations. $E[N_t] = \mu_{t-1}$ and $Var[N_t] = \mu_{t-1}(1 + \alpha_{t-1}^N \mu_{t-1})$. In parentheses, standard errors computed using subsampling with 100 subsamples.



**APPENDIX FIGURE S3:** Distribution of Number of New Treatments

Notes: Model is specified in (S1). Figure shows the empirical distribution of the number of new treatments and the average over time of the predicted probabilities using the estimated parameters in Table S9.
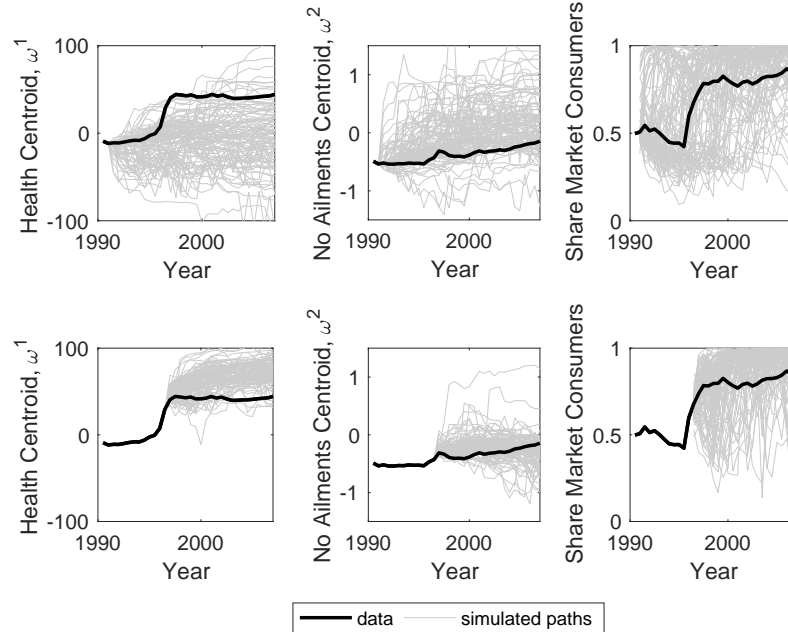
**APPENDIX TABLE S10:** Within Cluster Share Function

| variable | coef. ($\gamma^w$) | se |
|---|---|---|
| *Ailments Rk* | -0.427 | (0.124) |
| *Ailments Rk $\times$ Health Rk* | 0.074 | (0.020) |
| *Health Rk²* | -0.029 | (0.008) |
| *Ailments Rk²* | -0.019 | (0.006) |
| *NP* | -0.509 | (0.048) |
| *Health Rk $\times$ NP* | 0.046 | (0.009) |
| *Ailments Rk $\times$ NP* | 0.063 | (0.010) |
| *Ailments Rk $\times$ Health Rk $\times$ NP* | -0.007 | (0.002) |
| *New* | -0.352 | (0.508) |
| *New $\times$ NP* | 0.027 | (0.404) |
| *Constant* | 0.786 | (0.121) |

Notes: Parameters estimates from (S7) and (S8). *Rk* stands for the rank of the characteristic compared to other treatments within a cluster. *NP* is the cluster size. *New* indicates whether the treatment just entered the market. In parentheses, standard errors computed using subsampling with 100 subsamples.

## D.2 The Likelihood of Observed Technological Progress



**APPENDIX FIGURE S4:** Distribution of Technology Paths: Technology and Consumption
Notes: 100 simulated paths conditional on the state of the world in 1991 and 1996.
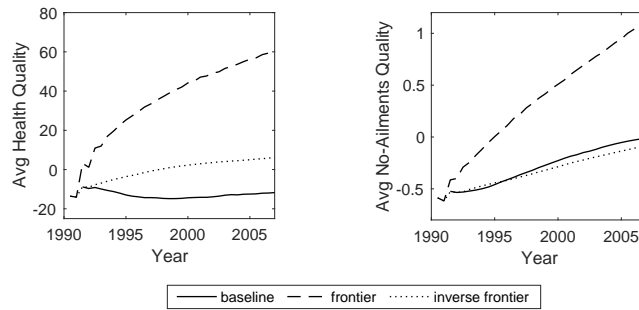
## D.3 Demand Pull: How Consumer Choices Affect Technology

The evolution of technology, and ultimately consumer welfare, is affected by demand externalities arising in the innovation process. We quantify the importance of these externalities by describing how the set of available treatments would evolve if consumers had less influence over the process of innovation, restricting the role of demand pull. However, since we do not explicitly model the primitives of the supply process (i.e., how product makers make their decisions), our first counterfactual here imposes the strong assumption that the objects describing the evolution of the set of available treatments ($g_\theta$, $g_N$, $\{\underline{s}, \bar{s}\}$) remain unchanged. We investigate two ways in which the process of innovation is detached from demand. The first counterfactual assumes a scientific body determines exit exclusively on the basis of treatment quality, and at the entry margin innovation is based equally on the char-

acteristics of all available treatments, with no special weight on the characteristics of popular treatments. The second counterfactual eliminates the effect of repeat purchase on innovation. For each experiment we present results averaging over 500 simulated paths starting at the first semester of 1991.

**Exogenous scientific intervention.**   In the first counterfactual regime innovation is independent of consumer demand. Thus, new treatments characteristics are no longer dependent on demand. At the entry margin we transform the centroid to be simply the average of the characteristics of treatments currently available on the market, as opposed to the share weighted average in the baseline model in Section 3, and take the estimates from the law of motion of the set of available treatments in Section 5.2 as given. Since $g_N(N_t | \kappa_{t-1}, s_{et-1})$ depends on $s_{et-1}$, we use the path for the experimental treatment share resulting from averaging the simulated experimental treatment share paths from the baseline model. By following this approach we keep that part of the comparison constant relative to the baseline. At the exit margin we also separate treatment exit from demand by adopting two alternative exogenous exit rules designed to resemble the actions of scientific authorities tasked with keeping only the best treatments on the market:

↪ *Frontier.* Any treatments that is not in the technological frontier is dropped from the market. This rule provides an upper bound for how fast innovation can move.

↪ *Inverse frontier.* We use the exit rate path resulting from averaging the simulated exit rate paths from the baseline model. This exit rate determines the number of treatments $n_t$ to be dropped. We define the inverse qualities of treatment $k$ as $-\theta_k$ and the inverse frontier as the technological frontier constructed using the inverse qualities. Then we drop $n_t$ treatments at random from the inverse frontier. If $n_t$ is larger than the amount of treatments in the inverse frontier, we construct the new inverse frontier and repeat the process until $n_t$ treatments are dropped from the market. This exit rule captures expert intervention in a less draconian way.

**APPENDIX FIGURE S5:** Alternative Regimes: Exogenous Scientific Intervention
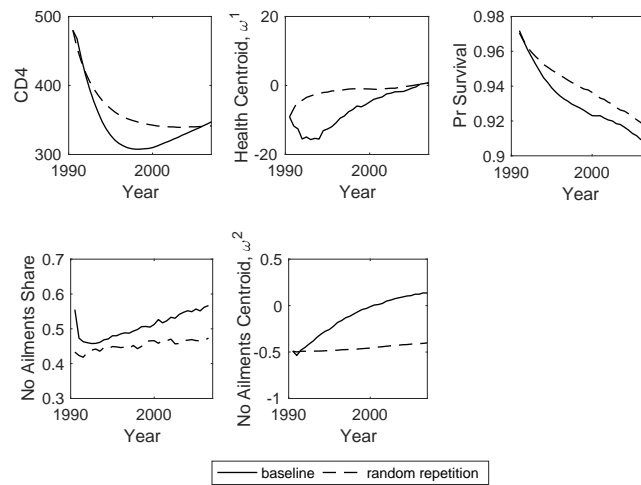
Notes: Evolution of the average treatment quality in the market under alternative regimes (500 simulations per regime) conditional on the state of the world at the first semester of 1991. The *baseline* is the estimated model of demand-pull innovation in Section 3. In both the *frontier* and *inverse frontier* regimes the centroid is not driven by demand as it is a simple average of the characteristics of treatments available on the market. In the *frontier* regime all treatments inside the quality frontier are exogenously withdrawn. In the *inverse frontier* regime the exit rate of treatments equals the average exit rate in the baseline simulations but the treatments that are withdrawn are the worst, independent of their demand.

Figure S5 shows that under the *frontier* regime innovation is more rapid, leading to much better treatments on both dimensions of quality. In contrast, the path of treatment quality is not as different from the baseline under the *inverse frontier* regime. Since individuals already avoid using the very worst treatments in the decentralized economy, an intervention that removes these treatments exogenously has little impact on the centroid, and hence on subsequent innovations. Nevertheless, our estimates imply that the inverse frontier regime does lead to somewhat higher average health quality.

**Eliminating the effect of repeat purchase.** Since consumers dislike changing treatment, they face a tradeoff between old and new technologies, and are more likely to repeat purchase if prior treatment offers better bundles of qualities than current clusters. In this regime we study the evolution of treatment quality when the process of innovation remains responsive to demand but demand by repeat consumers is not guided by their preferences and individual characteristics. Concretely, we assign individuals to alternatives in the choice set in the same proportions as the baseline (including the experimental treatment and no treatment), but make repeat consumption of old technologies random. By setting the unconditional shares of this counterfactual regime to match the unconditional shares in the baseline we

A-30

avoid spurious effects on the process of innovation yielding from arbitrary aggregate shares (e.g., $1/G$ for a choice set of size $G$). This regime neutralizes the dependence of the technological path on the preferences and characteristics of repeat consumers without changing the nature of the law of motion of available treatments.

Figure S6 shows that in the counterfactual regime the path of innovation is tilted towards more effective treatments with greater side effects. In other words, eliminating the effects of repeat consumption improves health and survival, but leads to more physical ailments. The reason is that individuals prefer medical treatments with fewer side effects despite the detrimental impact on their survival.



**APPENDIX FIGURE S6:** Alternative Regimes: Eliminating the Effect of Repeat Purchase.
Notes: Average paths computed over 500 simulations that are conditional on the state of the world at 1991. The *baseline* is the estimated model of demand-pull innovation in Section 3. The baseline solid lines in Figure S6 are the averages of the grey lines in Figure 9 and Figure S4 in Appendix D.2. Individuals in the alternative regime are assigned alternatives using the unconditional shares from the baseline model as assignment probabilities.

### D.4   Further Details of Policy Counterfactuals

**Mandated treatment.**   The first planner can only assign alternatives based on whether a person's health is high or low and whether the person decided to consume a market treatment last period (either by repeating his previous market treatment or by choosing a cluster). Hence, the planner's policy rules can be based only on four different categories. The planner can send all individuals in each of the four groups

to any of the $J + 2 + r_{it}$ alternatives available. We nest the baseline decentralized allocation by adding this allocation as an alternative in the planner's action set. Hence, there are $J + 3 + r_{it}$ alternatives in the planners action set and he can base his assignment on 4 categories. Since only two of the four categories can repeat their previous market treatment (when $r_{it} = 1$), this amounts to $7^2 \times 6^2 = 1,764$ policy rules. An example of a policy rule is presented in Table S11. We precompute a set of continuation values and match them to allocation rules to avoid forward simulation for each rule. We further explain these procedures below.

**APPENDIX TABLE S11:** Example of an Action-Constrained Planner's Policy Rule

| Category | | Alternatives | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Health status | Treatment $t-1$ | Cluster 1 | Cluster 2 | Cluster 3 | Trial | Repeat | No treatment | DA |
| high | yes | | | | | | | x |
| high | no | | x | | | | | |
| low | yes | | | | x | | | |
| low | no | | | | | | x | |

Notes: *Treatment $t-1$* column indicates whether individuals in this category consumed a market treatment in $t-1$. *DA* column indicates that the planner assigns the decentralized allocation.

**Optimal Consumption of Experimental Treatments.** The second planner we consider can base his policy on the entirety of the individual state but his action set has only two elements: he can give the person the experimental treatment or he can allocate the decentralized allocation (excluding the experimental treatment). Policy rules for this planner are experimental treatment shares and his problem also nests the decentralized allocation. For policy rules below the decentralized trial share $s_{et}$ the planner incurs a welfare cost by preventing people from rationally consuming an experimental treatment. For policy rules above $s_{et}$ he incurs a welfare cost by assigning people to consume an experimental treatment against their rational preferences. Welfare gains, if any, come from the externality via demand for experimental treatments, which pushes innovation. We evaluate policies in increments of 0.5 percent points, which amounts to 202 policy rules. Here we also use the set of precomputed continuation values and match them to allocation rules to avoid forward simulation for each rule.

### D.4.1 Continuation Values and Smoothing

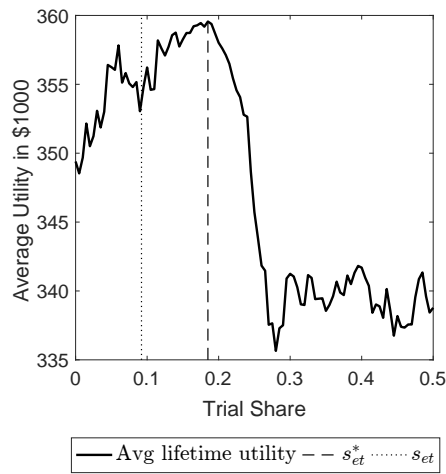We obtain continuation values for every planner rule by implementing the following algorithm:

1. Create a collection, denoted $\mathscr{A}$, of 500 continuation value vectors computed for all $t+1$ states. Each row in a value vector is an individual. Each value vector $v \in \mathscr{A}$ corresponds to a $t+1$ aggregate state $z_{t+1}^{v}$.

2. For each rule $n$ in a given planner problem, we compute each individual's current payoff and their future state, as well as the implied $t+1$ aggregate state $z_{t+1}^{n}$.

3. We match rule $n$ to the continuation value vector $v^* \in \mathscr{A}$ corresponding to the $t+1$ aggregate state that is closest to the aggregate state induced by rule $n$. In other words, we match rule $n$ to the continuation value vector $v^*$ that solves:

$$v^* = \arg\min_{v \in \mathscr{A}} ||z_{t+1}^{n} - z_{t+1}^{v}|| \tag{S34}$$

We use a measure of Euclidean distance that yields from discretizing the aggregate states $z_{t+1}^{n}$ and $z_{t+1}^{v}$ into vectors with 196 components. We scale each component of the discretized aggregate state vectors to be between zero and one by dividing over its largest value.

4. We repeat steps 2 and 3 one thousand times for every rule $n$ and average over repetitions.

As Figure S7 shows, our method of matching continuation values generates noise around the mapping from planner rules into average consumer lifetime utility for the planner who chooses the optimal experimental treatment share $s_t^*$. Hence, we use a local polynomial to smooth the mapping in an interval starting at the decentralized share $s_{et}$ and going 15 percent points above it (from 0.09 to 0.24). This produces Figure 10 and the results associated with it in Table 5.

**APPENDIX FIGURE S7:** Optimal Assignment to Experimental Treatments

Notes: The solid line represents average lifetime utility. The dashed line indicates the planner's optimal share $s_{et}^*$. The dotted line represents the decentralized share $s_{et}$. Year is 1996.