

Managing Intelligence: Skilled Experts and AI in Markets for Complex Products*

Jonathan Gruber^{1,4}, Benjamin R. Handel^{2,4}, Samuel H. Kina³, and Jonathan T. Kolstad^{2,4}

¹MIT

²UC Berkeley

³Picwell, Inc.

⁴NBER

April 1, 2021

Abstract

New artificial intelligence (AI) technologies are increasingly being used to augment expert decisions. We study the implementation of an AI-based decision support tool in a private Medicare exchange where consumers are randomized to skilled agents over time. We find that agents are able to both use the new information and to effectively incorporate private information, leading to lower enrollee costs and higher ex-post satisfaction. While AI is a complement to skill on average, we find that it is a substitute across the skill distribution, but offer cautionary evidence that adverse selection worsens as a result.

*We thank Ned Augenblick and Filip Matejka for excellent comments and participants at MIT, Northwestern, University of Arizona, BU/Harvard/MIT health seminar, University of Rochester, National University of Singapore, Western Economic Assoc. Conference, University of North Carolina, Haas, BYU and the NBER workshops on Economics of AI and Machine Learning in Health Care Markets for feedback. All authors hold ownership stakes in Picwell. The findings represent our own views and all errors are our own.

1 Introduction

Skilled experts play a key role in assisting consumers in decision making in a variety of markets, from doctors advising patients to financial planners assisting investors, and beyond. The centrality of expertise in market function has long attracted the attention of economists. Since at least Arrow (1963), economists have focused on the potential role of skilled experts in improving choice quality and market function. The rise of artificial intelligence (AI) along with increasingly rich data offers an alternative to human expertise.¹ AI is distinguished by the potential to perform tasks typically reserved for human expertise as well as to exceed human performance on some domains involving complex computation.

Despite this promise, current incarnations of AI tools generally perform very well only on a subset of tasks required of a human expert. Thus, AI is frequently brought to market alongside human experts. Understanding AI, therefore, requires not simply asking the aggregate questions of substitution between humans and machines but also the detailed interaction between humans and AI tools where different aspects of the production function can be either augmented or replaced by technology. Whether and how technology complements or substitutes for skilled expertise will shape both product markets where expert recommenders play a role and the associated labor markets for that expertise (Acemoglu and Restrepo (2019); Athey et al. (2020)).

We study the role of AI in a specific setting — the market for health insurance — that lends itself to asking general questions about the interaction of AI and expertise in determining market outcomes. Many health care systems and policies rely on the private provision of insurance to cover consumer health risks, from Medicare Advantage, which offers private plans that seniors can purchase in lieu of enrolling in original Medicare, to the private insurance options offered on state exchanges under the Affordable Care Act.

There is now a substantial body of evidence showing that, in practice consumers face difficulties in making choices in insurance markets (see Chandra et al. (2018) and Handel and Kolstad (2015a) for overviews of this literature). Consumers leave large amounts of money on the table in their insurance, due to a combination of factors including search costs, switching costs, insurance literacy, inattention, and limited information. Structural models show systematic choice inconsistencies in consumer decisions over Medicare Part D prescription drug plans (Abaluck and Gruber (2011, 2016, 2017), Ketcham et al. (2012), Heiss et al. (2010), Polyakova (2016), Ericson (2014), Ho et al. (2017), Ketcham et al. (2016)), employer sponsored insurance (Bhargava et al. (2017), Handel (2013), Handel and Kolstad (2015b)), and Medigap (Fang et al. (2008)).

One possible solution to this issue could be to rely on skilled agents — insurance brokers and enrollment agents. Despite the potentially important role they play in market function, there is little empirical evidence on the role played by agents in choices in health insurance markets

¹We will use the term Artificial Intelligence broadly and interchangeably with machine learning. There are important differences but they depend on myriad definitions of each. We view the semantic distinction as beyond the scope of our paper and instead use the term AI which has been widely adopted by firms and policy makers considering the degree to which prediction can facilitate decision making (see e.g. Agrawal et al. (2018))

nor the degree to which reliance on human expertise addresses choice errors.² Evidence from other settings suggests that skilled agents do not ameliorate market failures stemming from a lack of consumer information. For example, financial advisers maximize their own fees not client results (Mullainathan et al. (2012), Egan et al. (2019), Egan (2019), Gambacorta et al. (2017)) and auctioneers have dramatic differences in the prices they obtain for homogeneous products in relatively simple auction formats (Lacetera et al. (2016)).

As technological capabilities improve and detailed data become available, decision aids — AI-based assessment of choices — offer an additional option to improve consumer decisions in insurance markets. Research in this area is, however, still nascent, focusing on either rudimentary forms of information provision (e.g. Kling et al. (2012)), lab experiments (e.g. Bhargava et al. (2017)), tools that are poorly designed (e.g. Abaluck and Gruber (2016)) or tools that focus only on known costs (e.g. current drugs or planned utilization) not predicted total spending and risk (e.g. Bundorf et al. (2019)). The limited evidence to date finds that consumers are unwilling to engage with decision support even when it is a default (Abaluck and Gruber (2016)) or when it is demonstrated to provide financial value for those that do adopt it (Bundorf et al. (2019)). This raises the question of whether decision support can be more effective when paired with skilled experts.

In this paper we study the roles of skilled expert intermediaries and sophisticated AI-based decision support technology in health insurance markets. We focus on a retiree health insurance exchange, one of the largest private Medicare exchanges in the United States (hereafter The Exchange). The Exchange offers products from across the Medicare universe including Medicare Advantage (MA), Medicare Part D (prescription drug coverage) and Medigap (supplemental financial coverage). Firms contract with the Exchange to offer insurance options to their retiree population. We focus our analysis on the MA market, which in 2019 enrolled 22 million seniors nationwide (34% of all seniors in Medicare). Over the time period 2015-2017 we study the behavior of the approximately 800 enrollment agents advising seniors on MA choices. Agents expend considerable effort in terms of time spent advising enrollees; in 2015, prior to the adoption of AI-based decision support, the mean agent-consumer call time in our study sample was 59.4 minutes.

During our study period, the Exchange partnered with Picwell, a technology firm, to implement AI-based decision support for its agents who were tasked with assisting consumers in their plan choices. Beginning in 2017, this decision support was fully integrated into the enrollment software used by agents, providing them by default with individual-specific information on the quality of different plan options. Throughout our study period, potential enrollees were randomly assigned to agents, allowing us to avoid issues related to endogenous matching between agents and consumers. Combining random assignment to agents with the integration of Picwell into the agent enrollment software allows us to investigate the effectiveness of skilled agents in aiding consumer insurance choices and how that effectiveness changes when their production function is disrupted by new technology.

We begin our analysis by assessing how AI changes expert behavior. We develop a simple model of the joint agent-consumer decision, with a focus on (i) how agent advice shapes consumer decisions

²See Karaca-Mandic et al. (2018) for a discussion of the role of brokers.

and (ii) how changing agent information influences consumer decisions. We pay particular attention to the fact that the AI tool provides sophisticated information on some dimensions (notably the expected plan financial impacts on spending and risk) but not on other dimensions (e.g. plan networks and consumer preferences for those networks). The model motivates our empirical work in which we model how agents weight (i) information incorporated into the algorithm and (ii) information that is excluded from the algorithm but is likely observed by the agent. We begin by assessing the positive question of whether and how experts use the AI-based information before turning to a discussion of the likely welfare implications.

To understand the factors impacting choices by experts and consumers with and without AI, we estimate a structural model of plan choice as a function of plan characteristics. We show that in 2015, before decision support was fully integrated, agent-consumer choices display a number of inconsistent choice weights; for example, choices are between 6 and 7 times more sensitive to plan premiums than to expected out-of-pocket spending. We find that the average plan enrollment decision in 2015 led to a \$1,260 financial loss for consumers, relative to the best financial option available to them. These foregone savings amount to 30% of total costs, which is comparable to findings from previous studies of consumer choice of Part D drug plans and consumer choice of employer-provided insurance.

In contrast to 2015, once experts have fully-integrated AI tools in 2017, the decisions they make with consumers place almost identical weights on premiums and expected out-of-pocket spending. We use our structural model to control for changes in the nature of choice set. We find that the 2017 choices made under 2015 decision-making weights are \$278 worse on average than 2017 choices made under 2017 decision-making. Notably, large financial losses are substantially more likely under 2015 decision-making than under 2017 decision-making. We also use aggregate data on plan market shares in these same markets to show that there was no general shift towards more appropriate plans over this time period.

Studying the financial aspects of plan choice alone limits our attention to the domain where the AI tool is plausibly able to affect choices but says little about the role of expertise in the other dimensions of plan choice, a central goal of the paper. To address this, we also study a key element of plan choice that is *not* included in the algorithm: provider network breadth. Network breadth and preferences for that breadth are a feature where human agents might have a comparative advantage in assessment. Since this information is an “unused observable” (Finkelstein and Poterba (2014)) in the algorithm, we can test for asymmetric information in recommendations and the degree to which following algorithmic recommendations crowds-out or distorts recommendation components observed by decision makers but not the algorithm.

Our estimates of the importance of network breadth on decision-making remain essentially unchanged from 2015 to 2017. Moreover, we find that incorporating AI leads to a dramatic reduction in the weighting of brand effects across brands which deliver comparable products along observed dimensions. There continues to be a large brand preference for one brand (Kaiser Permanente) which delivers substantially different care through a vertically integrated provider network – but decision support leads to fewer individuals incurring very large financial losses under Kaiser and

more incurring more modest losses. This suggests the skilled agents retain some value in eliciting private preferences that may not be included in, and could be distorted by, algorithmic recommendations. These results demonstrate agents' ability to synthesize information that trades-off marginal effects, rather than relying on more blunt heuristics such as blanket brand preferences or blindly following the AI recommendations.

Taken together, we show that agents rely on AI-based information but also integrate other information that is welfare-relevant to the consumer. This positive assessment of the role of AI and expertise is also suggestive of consumer welfare improvements, but does not allow us to precisely assess whether the combination improves consumer-plan matches. To address this issue, we assess the quality of agent/consumer choices using a simple reduced form measure of experienced utility: whether an enrollee re-enrolled in the plan chosen in a subsequent year after experiencing the plan.

We find that people who enrolled in plans that the algorithm did not recommend in 2017 were more than twice as likely to switch plans the following year, compared to people who enrolled in a recommended plan. We augment this analysis by instrumenting for plan choice using a judges design with agent fixed effects. We find a causal impact of higher AI-based recommendations/score on subsequent enrollee experience utility, measured by turnover. For an agent who is one standard deviation better (3 plan score points) in terms of predicted plan score, consumers are 7 percentage points less likely to switch plans the following year. This effect is equivalent to nearly the full propensity to switch plans in the overall sample.

While we conclude that AI is a complement to agent skill overall, we also demonstrate that it is a substitute for skill within the agent distribution. We estimate a model of plan choice quality in 2015 that includes agent fixed effects as a measure of average agent skill. Measured agent skill heterogeneity is large, with foregone savings that are twice as large in the worst quintile of agents as in the best quintile. Offering AI-based decision support compresses this distribution, with little impact on the top quintile and substantial improvements for the bottom quintile. We show that these results are very unlikely to result from statistical mean reversion but, instead, the result of the compression of skill towards the top once AI is implemented. We also find that the introduction of AI lowered call times by roughly 20% throughout the distribution of agent skill, so that the return to the tool is much larger for the least skilled agents.

Our final piece of empirical analysis concludes on a cautionary note, however. We consider the inherent link between choice adequacy and adverse selection, as discussed in, e.g., Handel (2013): it is possible that better individual consumer choices can facilitate more acute sorting of sicker consumers to generous plans (and vice-versa), leading to a greater degree of adverse selection. While a number of papers have discussed this concern, no paper that we are aware of shows how reduced choice errors actually impact adverse selection in a given empirical context. We demonstrate that improved choices do lead to more acute sorting, implying the potential for greater adverse selection if tools like the one we study are rolled out to all market participants.

The rest of the paper proceeds as follows. Section 2 presents the data and setting. Section 3 develops our model and empirical approach. Section 4 presents results and Section 5 concludes.

2 Data and Setting

2.1 Medicare Advantage

Medicare provides universal government-sponsored health insurance for the elderly and disabled in the U.S.. Enrollees can access coverage through various channels. Medicare-eligible individuals are automatically enrolled in the Medicare Part A program, which covers inpatient hospital expenses. Eligible individuals can elect to enroll in Medicare Part B to cover outpatient expenses and choose among privately provided Medicare Part D plans to cover prescription drug expenses. Beyond this, privately offered Medigap plans that cover out-of-pocket costs under Medicare Parts A and B (Original Medicare) are also available. A combination of these options constitutes an enrollment in Original Medicare.

Alternatively, an enrollee can opt out of Original Medicare by choosing among a set of competing private Medicare Advantage plan. MA plans can be offered as a stand-alone plan only covering medical care or a product that combines this coverage with prescription drug insurance: MA-PD.³ MA plans can offer additional benefits above and beyond those provided by Original Medicare and may charge additional premiums. MA plans are offered on a county-by-county basis, and their plans offer services through managed care networks. Nationwide, approximately 22 million people - or about 34% of those eligible for Medicare - enrolled in an MA plan in 2019. In addition to the patient premium, MA plans receive reimbursements from CMS based on bids that they submit, costs relative to local original Medicare costs, and risk-adjustments to account for differences in the enrolled population (Geruso and Layton (2015), Brown et al. (2014), Newhouse et al. (2012)).

We focus on MA-PD plans in our analysis. This allows us to study choices of insurance across the bulk of health care utilization (i.e. prescription drugs, inpatient and outpatient medical care). In this way, MA-PD plans more closely resemble the kinds of health benefits chosen by those outside of Medicare (e.g. employer-based) and in other settings (e.g. Medicaid and outside of the U.S.). The complexity of this more general setting may also make expert recommendations more beneficial than in simpler settings where consumer choice quality is studied (e.g. Medicare Part D coverage for prescription drugs).

People who choose to enroll in MA-PD plans cannot also enroll in Medigap or Part D plans. That is, the Medicare enrollment decisions that people make can be described as a decision tree where, at the first level they choose whether to build their coverage around Original Medicare or MA. Those who choose Original Medicare must then select a Part D plan and they must decide whether they want additional Medigap insurance. Those who choose an MA-PD plan must choose from among the plans available, but they only need to make one plan choice. This simplifies our analysis and makes it likely that the choices we observe reflect the near entirety of consumers' health benefits.

³See Starc and Town (2015) for a discussion of the impacts of bundling coverage components together on benefit design in the MA market.

2.2 The Exchange

The Exchange is one of the largest private Medicare exchanges in the United States. It is offered as a service to employees of a set of large employers. Medicare-eligible workers, retirees and family members can select plans on the Exchange throughout the year as they turn 65 or have another qualifying event (e.g., losing employer-provided coverage) that triggers a special enrollment period. Approximately, 10% of enrollees are making first time enrollments. Individuals do not have to enroll in Medicare through this exchange, but if they use the exchange they can only enroll through an agent.

Most enrollment in the Exchange occurs between October 15 and December 31 each year during the Open Enrollment period. During the 2017 Open Enrollment period (which occurred between October 15, 2016 and December 31, 2016), 87,691 Medicare eligibles enrolled in plans through the Exchange. 41,563 of these people enrolled in a Medigap plan, 44,883 enrolled in a Part D plan, and 34,616 enrolled in a MA plan.

The Exchange employs enrollment agents who help customers understand their Medicare options and enroll in a plan. Each open enrollment period, the Exchange schedules appointments for customers to review their options and enroll in a plan with an agent. All agents are licensed to sell Medicare policies in multiple states, and agents are randomly assigned to appointments with customers.

To confirm that assignment is random we divide agents by baseline skill level. We return to our definition of skill below but demonstrate here that assignment is orthogonal to agent skill. This allows us to test for random assignment that is a primary threat to our analysis: that particular types of enrollees (e.g. more complex) might be assigned to particularly agents (e.g. more skilled at handling particular levels of complexity or specific conditions).

Table 1 presents gender, age and number of prescriptions by quintiles of agents skill. We can reject differences across agents in any of these key demographics that determine both the value of a specific product and the potential complexity of identifying the right plan.

Table 1: Summary of 2015 customer characteristics by agent skill quintile

Agent skill quintile	Female share of customers		Customer age		Rx per customer	
	Mean	SD	Mean	SD	Mean	SD
1	0.55	0.50	71.2	5.3	4.4	3.1
2	0.55	0.50	71.0	5.1	4.5	3.0
3	0.56	0.50	71.0	5.2	4.5	3.1
4	0.56	0.50	71.2	5.3	4.5	3.0
5	0.56	0.50	71.8	5.7	4.6	3.0

Agents were paid hourly and received additional bonuses when customers enrolled in a plan. Enrollment bonuses did not vary by plan within a market, but did vary conditional on market (e.g. MA-PD, vs. Medigap vs. Part D only). The bonus for customers enrolling in just a Part D plan was less than half of the bonus for enrolling in both medical and drug coverage, so, while agents

did not face financial incentives to direct customers towards one particular health plan in a given market, they did face incentives to enroll agents in a plan or set of plans that would cover both medical and prescription drug costs.⁴

In 2015 agents used web-based enrollment software that did not include decision support for MA plans. Agents could access a list of plans available and had access to a third party tool to look up whether physicians were included in specific plan networks. Decision support was available in 2015 but only a few weeks prior to open enrollment and to use the tool an agent had to leave their existing software platform. Accordingly, few agents consistently used Picwell decision support for a variety of potential reasons including because they were unaware of it, they were not familiar with how to use it and how to explain the resulting recommendations to customers, and/or they did not fully trust the recommendations.

The Exchange embedded decision support in the software used for all of their customers in 2017. This change was accompanied with training on how Picwell scores and cost estimates are generated and how to use the decision support tool to help Exchange customers choose Medicare plans.

To receive a Picwell recommendation, agents would walk customers through a user intake process that required entering in personal information across several steps. In the first page, agents would enter in the customer’s age, sex, zip code and county. Agents would then ask customers about any prescription drugs that they routinely take. After entering this information, agents would request recommendations, and Picwell would return sets of recommendations for all Medicare Advantage, Medigap and Part D plans sold on the Exchange in the customer’s county. Within each type of Medicare plan, recommendations would be returned, sorted by the Picwell Score. Picwell Scores were presented with one of three color tiers - Green, Yellow, and Red - that indicate levels of plan suitability for a customer, in descending order. Agents were instructed to interpret the color tiers as “Good”, “Fair” and “Poor” matches, respectively. In addition to the Picwell Score and color tier, recommendation results also included the expected OOP cost and a range of estimated costs that indicated the 20th and 80th percentiles from the distribution of predicted OOP costs for each person and plan. In addition to these features, more conventional information was offered including a description of plan characteristics such as the plan type (e.g., PPO or HMO), deductibles, OOP maximums, and, as in 2015, a third party tool allowed agents/consumers to look up the network status of their providers.

It is important to highlight here that our intervention captures the effect of providing decision support within the context of skilled agents. A separate and important question is whether providing similar decision support directly to individuals has similar effects. Bundorf et al. (2019) provides insight into this question through the offer of this same decision support tool to individuals purchasing Medicare Part D prescription drug coverage. They find that for those who use the tool, choices improve substantially along the lines that we show here. Improvement depends

⁴This level of alignment between agents and enrollees is not representative of many insurance transactions in the U.S. Frequently, insurance brokers are paid commissions that vary by carrier and need to be disclosed to the enrollee. Thus, we expect our results on agent skill to be an upper bound for the quality of agency, conditional on agent ability.

on providing the algorithm as an “expert recommendation” but, despite the improvements, only a very small minority of individuals use the tool.

2.3 Decision Support Background

Agents on the Exchange were able to use decision support technology to evaluate and compare Medicare plans available to Exchange customers. The decision support evaluates all options within a specific type of Medicare plan type (e.g. all MA-PD plans), but it does not make comparisons across plan types.⁵

The plan evaluations include (i) predicted “RealCost” which combines annual premiums with mean estimated OOP, (ii) a Picwell plan Score that rates plans on a 100 point scale, and (iii) a color tier that is simply a mapping of plan score to one of 3 color tiers with scores of 90 or greater assigned to the “Green” tier, scores of 75 to 89 assigned to the “Yellow” tier, and scores of 74 and lower assigned to the “Red” tier. The Picwell Score identifies the “utility maximizing” plan for a risk averse consumer within each choice set and assigns the highest score to this plan. Scores for all other plans identify how close the expected utility of each plan is to the plan with the highest expected utility. Agents were instructed to interpret the Picwell Score as an identifier for how well each plan matches a customer’s preferences and to treat any plan on the “Green” color tier as a good match.

The process of generating a set of plan recommendations can be divided into three distinct steps. In the first step of this process, a machine learning model predicts annual medical spending for individual k . The machine learning model is trained on an extract from the IBM Watson MarketScan claims database that includes 2 years of continuous claims and enrollment for approximately 1.2 million MA enrollees.⁶ Model features are defined based on observable characteristics in the first year (t) of the 2 year claim period, and the prediction target is the total allowed costs incurred in year $t + 1$. This generates a mapping from individual characteristics μ_k (including age, sex, and a list of prescription drugs) to a “risk group” of individuals K in the claims data with a distribution of allowed costs $f(ALLOWED|\mu_K)$. The performance of the machine learning model compares favorably to other oft-used risk rating models in terms of out-of-sample prediction (see Appendix A for further detail).

In the second step, the decision support applies benefit calculators for each plan j to year $t + 1$ claims for each of the individual in risk group K to generate a distribution of OOP cost, including all drug and medical costs, for each plan $f(OOP|\mu_K, \psi_j)$. The benefit calculators account for detailed plan information (ψ_j) including deductibles, OOP maximums, formularies and coverage and cost sharing rules for every benefit category in each plan.⁷ The decision support calculates

⁵This technology is available. However, Medicare marketing rules limit the ability to make recommendations across bundles of plans.

⁶Each year, the model uses the most recent two year window of claims available, which covers the calendar years 2 and 3 years prior to the forecast date. For example, in the 2017 open enrollment period, the decision support tool generated cost predictions for 2018 that relied on claims from 2015 and 2016. The model uses the CPI-medical to adjust for inflation between 2016 and 2018, and it uses the 3 year average inflation rate to forecasts inflation for 2018.

⁷The benefit calculators are based on benefit design information contained in the Plan Benefit Package (PBP)

$E(OOP_{kj})$ based on $f(OOP|\mu_K, \psi_j)$, and uses this to return the RealCost, or expected total cost, for each person-plan pair, where $RealCost_{kj} = Premium_{kj} + E(OOP_{kj})$.

In the third step, the decision support calculates utility and a score from $f(OOP)$. Utility is calculated as a function of personal and plan attributes using a constant absolute risk aversion model.⁸ U_{kj} is translated to dollars by calculating a certainty equivalent CEQ_{kj} . This allows us to estimate a risk penalty, r_{kj} , where

$$r_{kj} = CEQ_{kj} - (P_{kj} + E(OOP_{kj})) \quad (1)$$

The risk penalty can be interpreted as the additional annual premium that an individual with risk aversion γ and individual characteristics μ_k would be willing to pay to eliminate all variance around $E(OOP_{jk})$. In other words, it represents a marginal willingness to pay to eliminate risk that represents both individual preferences toward risk and exposure to risk. Finally, the decision support applies a score function that translates each CEQ_{kj} in individual k 's choice set into a score between 0 and 100.

In 2017, the third year that decision support was available, agents were required to generate recommendation requests for all customers. Furthermore, just prior to the 2017 Open Enrollment period agents received comprehensive training in the use of the decision support technology. Combined, this led to near universal adoption of the tool and, potentially, enhanced trust in the recommendations.

2.4 Data

We study health plan choices for individuals enrolling in MA-PD plans through a private health insurance exchange (“the Exchange”) for the 2015 and 2017 Open Enrollment Periods. Unfortunately, we are unable to include the interim year of 2016 as the number of observations drops dramatically.⁹ We observe detailed information on approximately 59,000 MA-PD enrollees, their agents, and their enrollment options in both 2015 and 2017. At the enrollee level, we observe age, sex, zip code, county of residence and a list of prescription drugs that they take. At the agent level, we observe the identity of each customer, the number and duration of calls for that customer, the plan that customer enrolls in, and the number of years of experience each agent has working on the Exchange. For each choice made we observe the set of Medicare plans available to each enrollee and detailed information about those plans including premiums, deductibles, out-of-pocket maximums, coinsurance rates for various categories of coverage, actuarial value, plan type (e.g. PPO,

files that CMS publishes.

⁸In typical implementations, the decision support technology assigns risk aversion parameters to customers based on responses to survey questions about attitudes towards risk, but in this particular application, such survey questions were not permitted, so all customers were assigned a risk aversion parameter of $\gamma = 4.0 * 10^{-4}$, which is similar to estimates from Handel (2013).

⁹The low volume of agents and the low volume of customers per agent in 2016 do not allow us to draw any meaningful conclusions about the relationship between decision support and plan selection. The low number of observations is due to idiosyncratic market factors – the number of enrollments is a function of the customer population of the exchange and 2016 appears to be an unusually small year.

HMO, POS), brand name and network breadth. We also observe the predicted spending in each plan based on the Picwell algorithm and a plan score that enrollment agents could use to compare plans.

Table 2 presents summary statistics for the enrollee population in 2015 and 2017. We restrict our study population to people between the ages of 64 and 90 at the time of enrollment and only consider enrollment decisions made by people who ultimately enrolled in a MA plan that also covers prescription drugs (MA-PD). After these restrictions, the remaining populations for 2015 and 2017 were 31,090 and 27,739, respectively. In both years, approximately 55% of enrollees were female. The study population was slightly older in 2017, with a mean age of 72.7 compared to 71.2, but 2017 enrollees took slightly fewer prescriptions, with a mean number of 3.4 prescriptions per enrollee in 2017 compared to 3.7 in 2015. Overall, the population in the MA-PD market looks similar based on observables between 2015 and 2017. This suggests that the introduction of AI-based decision support did not systematically alter the extensive margin decision to select MA versus original Medicare.¹⁰

We observe enrollment appointment information for 835 agents in 2015 and 732 agents in 2017. Table 2 presents summary statistics. Average appointment duration was 59 minutes in 2015 and average appointment duration was 48 minutes in 2017. Conditional on ultimately purchasing a plan, 78% of consumers needed only one call to finalize the purchase while the remaining 22% needed more than one call to do so.¹¹ Of the 732 agents, 305 worked for the Exchange in both years. There was a shift in average experience from 2015 to 2017, reflecting new or seasonal hires. The average years of experience for agents in 2017 was 2.98 compared to 4.02 in 2015.

Table 2 presents summary statistics for the MA-PD plans available to enrollees in 2015 and 2017. In both 2015 and 2017, average premiums were similar, with a mean monthly premium of \$49.21 in 2015 and \$55.77 in 2017 (in addition to the base Part B premium), and the range in premiums available in a choice set was similar, with a mean difference between the highest and lowest monthly premium available of \$156.47 in 2015 and \$154.25 in 2017. In both years, HMO plans made up a similar share of plans offered. In 2015, slightly more regional PPOs and slightly fewer local PPOs were available. In both years, cost plans and private Fee for Service plans were rare.

In both periods there were a large number of plan choices available. The mean number of options is 12.5 with a 95th percentile of 23 options in 2015 and remains very similar in 2017. The one notable difference in Table 1 is a reduction in the typical experience of agents over time. While we do not know what is driving this change, we note that, to the extent that experience improves performance, we would expect this to bias against observing improved decision making. Moreover,

¹⁰These empirical results are also consistent with our understanding of the marketplace where the process for selecting a particular type of coverage did not change following the introduction of Picwell. AI was available once an enrollee selects a particular type of coverage to select among options but not in choosing MA versus product combinations under Original Medicare. AI was available in all product markets including MA, MA-PD, Medigap and Part D, suggesting that there was no reason to differentially steer an enrollee between different market options after the introduction of Picwell.

¹¹We do not observe calls that were made by consumers who did not end up purchasing a plan. Our understanding is that a vast majority of consumers purchase a plan once they engage with an agent in this market.

Table 2: Summary Statistics

		Enrollee		Agent		MA Plans	
		2015	2017	2015	2017	2015	2017
Enrollees		31,090	27,739	Agents	835	732	
Age				Years of experience			# Plans in Choice Set
Mean	71.15	72.73	Mean	4.02	2.98	Mean	12.43
p25	67	67	p25	3	1	p25	7
p50	70	71	p50	4	1	p50	13
p75	74	77	p75	5	6	p75	17
p95	82	86				p95	23
Drugs				Total customers			Monthly Premium
Mean	3.68	3.37	Mean	37.32	37.10	Mean	49.21
p25	1	1	p25	15	77	p25	0.00
p50	3	3	p50	29	33	p50	26.00
p75	5	5	p75	55	54	p75	80.00
p95	9	9				p95	188.00
Female	55.3%	55.5%	Call time			Plan Type	
			Mean	59.4	47.8	Cost	0.7
			p25	29	11	HMO	53.2
			p50	50	40	HMO-POS	8.1
			p75	79	69	PFFS	3.0
			p75	79	69	PPO	21.8
						Regional PPO	13.2
						Premium Range	156.47
							154.25

when we estimate our choice model just on the subset of agents who were present in both periods, our results do not materially change.

3 Choice Model and Empirical Specification

3.1 Insurance Demand

We begin with a baseline model of an expected utility maximizing enrollee choosing among insurance products. Each individual, indexed by $k \in K$ faces a distribution of cost outcomes for Medicare Advantage plan options $j \in J$ $F_{k,j}()$. Enrollees get latent utility from choosing a plan j according to the following von Neumann-Morgenstern (vNM) expected utility:

$$U_{kj} = \int_0^\infty f_{kj}(s) u_k(W_k, x_{kj}(P_{kj}, s)) ds \quad (2)$$

Here u_k is a vNM utility index and s is a realization of out-of-pocket cost from the distribution $F_{k,j}()$. Individual specific wealth is captured by W_k and P_j is the premium for plan j . x_{kj} is an individual's level of consumption based on their realized health shock, s . We model this as:

$$x_{kj} = W_k - P_j - s + \epsilon_{kj} \quad (3)$$

where ϵ_{kj} is a mean zero individual specific shock. We follow the literature and assume that families have constant absolute risk aversion (CARA) preferences implying that, for a given *ex-post* consumption level x :¹²

$$u_k(x) = \frac{1}{\gamma_k} e^{-\gamma_k x} \quad (4)$$

This specification constitutes the baseline choice model; reflecting the basic role that insurance plays as a product to mitigate financial risk. However, we extend the model to allow for three specific additional features that are likely to impact choices: (i) heuristic/behavioral decision making (ii) the role health insurance plays in allowing access to different health care providers and (iii) brand preferences for insurers. For (i), we include a set of salient plan features λ_j including the deductible and out-of-pocket maximum for plan j . The elements of λ_j are observable but, as equation 2 shows, do not affect realized utility conditional on the realized spending level s . Even though these financial characteristics of a plan only affect utility through their impact on the distribution of risk, consumers may still heuristically place weight on them when making choices (see, e.g., Abaluck and Gruber (2011)). For (ii) we also allow the network of available providers in plan j to enter utility captured by n_j .¹³ For (iii), we incorporate brand effects where a given insurer's brand κ is constant across all plans j that insurer sells from their set of plans J . Preferences for the insurer's

¹²As shown in the online appendix in Abaluck and Gruber (2011), the distinction between CRRA and CARA risk preferences in our context is very unlikely to matter materially for our empirical results.

¹³The Exchange used a third party vendor for the provider search tool that was included in the decision support experience. We use data that captures the percent of providers requested that were in network for each state and plan in 2017 as a measure of network breadth.

brand could capture pure non-welfare-relevant brand effects or welfare-relevant effects such as, e.g., differences in administrative support or online tools. Incorporating these features we express the utility of each state s as:

$$x_{kj} = W_k - P_j - s + \lambda_j + n_j + \kappa_j + \epsilon_{kj} \quad (5)$$

To implement the model empirically we parameterize utility as:

$$u_{kj} = \rho P_j - \delta E(OOP_{kj} | \mu_k, \psi_j) + \phi R_j(\gamma, s) + \omega \lambda_j + \beta n_j + \alpha_{j \in J} \kappa_j + \epsilon_{kj} \quad (6)$$

where P_j is the premium for plan j , $E(OOP_{kj} | \mu_k, \psi_j)$ is a measure of expected out-of-pocket spend for individual k in plan j based on the predictive model and R_j is a function reflecting the risk protective value of plan j . We assume a uniform risk aversion value in the population ($\gamma = 4.0 * 10^{-4}$) and compute R for each customer and plan following equation 1. Finally, ϵ_{kj} is a type-I extreme value error term.

3.2 Agency and AI

To this point our model has simply specified an enrollee utility from insurance. The central question of this paper, however, is the role played by both agents and AI in choices. We model enrollees who rely on agents to help them choose an insurance plan that maximizes their utility. Assume that true enrollee utility can be expressed as:

$$E(u_{kj}) = \Lambda X_{kj} + \epsilon_{kj} \quad (7)$$

where Λ captures a vector of weights on plan attributes that map to realized plan utility.¹⁴ Utility, with full information, is maximized over an expectation based on unbiased observations of plan attributes available to the enrollee and accords with a full information version of equation 7.

We assume agents have an (unobserved) information set/signal Ω^b that consists of information from enrollee k and knowledge of plan attributes for plan j . Latent agent skill as well as effort affect the realization of this signal. Agents then develop a set of weights for plan attributes that scale true enrollee utility captured by the vector Θ . Agents and enrollees then make enrollment decisions by maximizing expected utility according to:

$$E(u_{kj}) = (\Theta | \Omega^b) \Lambda X_{kj} + \epsilon_{kj} \quad (8)$$

Agents' weights (Θ) reflect the degree to which they are able to capture the true, latent preferences of enrollees (Λ). For example, elements of the vector Θ equal to 1 represent perfect agency; the agent puts the same weight on an attribute that an enrollee would in the standard model.¹⁵

¹⁴Without loss of generality we express utility here as a linear function of plan attributes. One could alternatively specify utility as $E(u_{kj}) = \Lambda f(X_{kj}) + \epsilon_{kj}$ to capture a flexible function of observables. We follow our empirical implementation in equation 6 in which we express CARA risk preferences as a linear term because we implement the model of agency and AI empirically using the same specification.

¹⁵We assume that enrollees have preferences following the benchmark model developed in equation 4. An alternative

AI-based decision support enters as an additional source of information that an agent includes in their information set, captured in the model as a new, post-AI, information set:

$$\Omega' = \alpha\Omega^{AI} + (1 - \alpha)\Omega^b \quad (9)$$

where α captures the relative weight on AI-based information versus prior beliefs used to form the new information set.¹⁶

We assume that the AI-based signal provides information on only a subset of the attributes that enter choice. This accords with our specific setting, where the AI-based tool focused primarily on financial/cost components of the choice. It also captures a general phenomenon in which AI-based tools are particularly good at accounting for quantifiable aspects of decisions but rarely account for the universe of welfare-relevant aspects of a decision.

Attributes are partitioned into X_{kj}^i and X_{kj}^g where i indexes attributes observed by agents and included in AI and g indexes attributes observed by agents but not included in AI. After AI-based decision support is introduced enrollment decisions are made to maximize expected utility according to:

$$E(u_{kj}) = (\Theta'_i|\Omega')\Lambda X_{kj}^i + (\Theta'_g|\Omega')\Lambda X_{kj}^g + \varepsilon_{kj} \quad (10)$$

Instead of specifying the micro-foundations either for information acquisition in equation 9 or the resulting weights in recommendations in equation 10, we allow for a flexible model in which AI can concurrently change the information available and weights on different attributes — those included in AI and those excluded. This allows us to flexibly capture a variety of ways in which agents might incorporate information. It nests models in which agents efficiently (in a Bayesian sense) integrate information (e.g. Diamond (1971), Dranove and Satterthwaite (1992)) as well as a broader class of models in which agents are rationally inattentive in selecting attributes, the associated weights and in making recommendations with and without AI (see Mackowiak et al. (2018) for a review). The structure also allows for more behavioral models in which attention is heuristically allocated in ways that need not be optimal and may reflect a variety of biases (see e.g. Handel and Schwartzstein (2018)).¹⁷

Following equation 10 we expect the introduction of AI-based decision support to alter the weight placed on attributes included in the AI-based tool when agents i) incorporate the AI-based

model would be that agents and information can not only alter weights on attributes but also the weights themselves in the utility function. Bundorf et al. (2019) study this distinction explicitly in an experimental design, though they do not observe experts/agents themselves. This distinction does not affect our empirical implementation as, in practice, we study the joint realization of weights between agents and enrollees with and without AI and focus only on how they change.

¹⁶We express information acquisition as a linear combination of beliefs. One could specify the updating process more generally. Since we do not empirically model learning itself we simplify exposition with linear weights.

¹⁷With sufficient assumptions on both the model of learning and the covariance of the elements of X_{kj} we could recover estimates for α . Rather than undertake this structural approach (for which one might want to implement a more flexible parameterization of this model) we focus on a measure of how attributes are updated that accord with measures of enrollee outcome utility to evaluate how agents update and what the positive and normative impacts of that updating are.

recommendations and ii) AI provides new information that changes beliefs. The aggregate change in attribute weights therefore depends on the combination of agent updating (captured by Ω') as well as the covariance of the attributes across choice set options.

For example, when an agent is able to access a predicted measure of out-of-pocket cost, this may make out-of-pocket costs more salient and provide more precise individual-plan-specific information on those costs. Both of these factors could contribute to an increased weight on expected out-of-pocket spending in observed plan choices.

Equation 10 also conditions the weights on X_{kj}^g — attributes not directly changed by AI — on the new information set Ω' . Weights on attributes not included in the AI still change because we expect relative attribute weights to change as the overall information set changes. For example, were agents to rely on heuristics to deal with the complex prediction problem of estimating cost, they might have put a high weight on plan premium or simple measures of generosity such as the level of deductible or out-of-pocket maximum prior to AI. If they gain new information on total cost we expect the weights on those measures to decline relative to the out-of-pocket cost prediction that can be generated with AI.

Our model does not assume that AI improves choices, in the sense of better approximating true preferences. AI-based decision support might make some attributes of a plan more salient at the expense of harder to observe but nevertheless valuable components of a plan. For example, particular insurance brands like Kaiser Permanente that enrollees prefer do not have their non-financial features accounted for in the AI (e.g. the breadth of physician network or specific health care delivery models). As a result, weights on these excluded attributes may decline relative to the weights on included attributes. Alternatively, if AI provides valuable new information and agents are rationally inattentive (i.e. they efficiently integrate information to optimize recommendations) we expect recommendations to (weakly) better approximate enrollee preferences. Which model better reflects behavior is an empirical question.

3.3 Identification

Our empirical analysis relies on two key sources of identifying variation in our data. First, we use the fact that customers are randomly assigned to agents to deal with any issues related to agent-customer matching. Second, we leverage the change from 2015 to 2017 when the marketplace moved from little/no AI-based decision support to integrating the AI-based recommendation into the enrollment software used by all agents for all enrollments.

Relying on intertemporal variation alone presents an obvious challenge: what other features changed over time that might affect plan choice quality? The summary statistics demonstrate consistency over time in key statistics for both enrollees and agents, suggesting that we do not need to be concerned about large-scale changes to those participating in the exchange. We are concerned, however, about changes in the set of plans available to enrollees.

We address this issue in two ways. First, we control for the features of plans offered to the enrollees in our sample. In particular, we estimate a variant of equation 6 separately for 2015 and

2017. The associated model parameters capture the weights placed on plan attributes in each of those years. Based on estimated choice model parameters, we simulate choices in 2017 *holding fixed* the set of plans available. We compute:

1. 2017 choices based on 2017 demand parameter estimates
2. 2017 choices based on 2015 demand parameter estimates

We use this approach to compare how choices change moving from 2015 choice parameters to 2017 choice parameters for the set of plans available in 2017. Put differently, using our structural plan model we estimate the choices that would have been made in 2017 by the same agents had they behaved as they did in 2015.

Second, we estimate comparable choice models both in our sample and in national data on MA-PD plan enrollment. In particular, we use data from the Centers for Medicare Medicaid Services (CMS) on the aggregate market share of each MA plan in every county served by our exchange. Our sample represents only .19 percent of the total enrollees that are used to measure market shares, so that this aggregate provides an essentially independent means of assessing time trends in choice. To match the sample as closely as possible we including counties from which our sample is drawn (approximately 54% of total counties, though about 90% of counties on a MA population-weighted basis).¹⁸ We then match the remaining counties to that sample using propensity score methods based on county-level demographics for the Medicare-eligible population from the American Community Survey.

To provide a comparison of choices quality, we estimate models of plan choice as a function of measures of plan premiums and plan generosity that can be estimated comparably both in our micro-data and in the aggregate market share data. We find that the impact of premiums on plan choices evolves comparably in both our data and the aggregate market share data, but that while in our data there is a notable shift towards plans with more generous out of pocket coverage, in the aggregate data choices evolve in the opposite direction over this time period. This strongly suggests that general changes in tastes are not driving our results. These results are described in Appendix F.

4 Results

Before moving to our primary results, using our choice model estimates, we discuss observed money left on the table in 2015. Figure 1 plots the money left on the table for a consumer’s actual choice in 2015, relative to the best possible financial choice in their choice set.

Clearly, there are substantial sums of money left on the table in 2015 when looking just at the financial dimension of consumers’ choices. More than half of consumers leave \$1,000 on the table while a meaningful proportion leave over \$2,000 on the table. The fact that consumers leave this much money on the table descriptively is in line with the basic facts put forth in other studies in

¹⁸In counties in which there are any exchange enrollments the average county-level share is .3%.

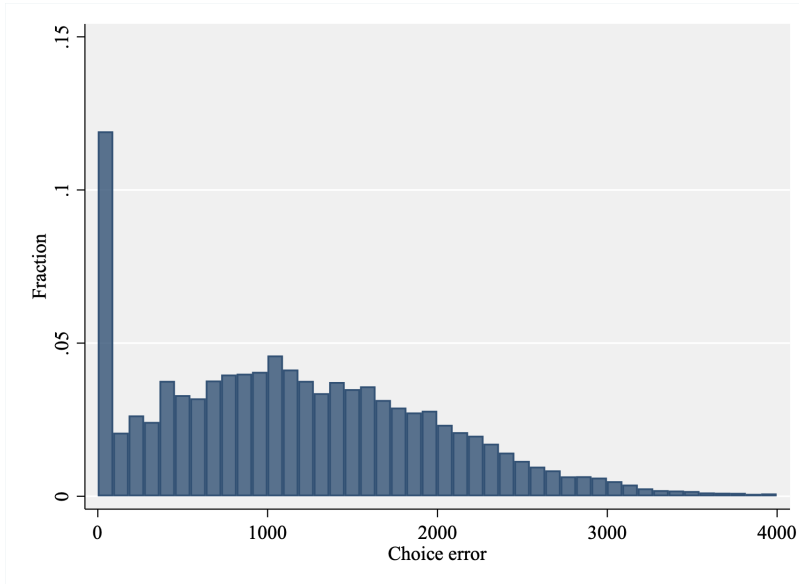


Figure 1: Histogram of observed 2015 money left on the table (financial choice error).

the insurance choice literature including Abaluck and Gruber (2011), Handel (2013), Handel and Kolstad (2015b), Bhargava et al. (2017) and others, as summarized in Chandra et al. (2018).

Though this evidence is suggestive of meaningful choice errors, there are quite a few reasons why this kind of histogram could be consistent with good choices. First, plans could be differentiated by networks of providers and brand effects (some portion of which may reflect substantive differentiation). Second, this is only an expected value for consumers. Consumers who value risk reduction (one of the primary purposes of insurance) will want to give up some expected value in return for lower cost variance. We now turn to our primary choice model estimates in order to (i) assess the weights placed on different financial dimensions of choice and (ii) account for potentially important non-financial plan dimensions.

4.1 Impact of Decision Support: Financial Value

Table 3 presents the estimates from our choice model, described in equation 4. The first two columns of the table present estimates for 2015, prior to the widespread use of algorithmic decision support. The first specification is a simplified version that includes both key inputs into the algorithm (annual premium and predicted OOP) and potentially important choice factors excluded from the algorithm (network coverage, plan type dummies, brand dummies). The second specification, our primary specification, also includes risk aversion and plan financial characteristics whose value should be fully subsumed by the predicted OOP variable but may not be due to customer / agent use of heuristics.

For our primary specification, in 2015, joint agent/consumer decisions place substantially more weight on plan premium than they do on expected plan out-of-pocket. For a rational, informed consumer — ‘homo economicus’ — these attributes should be valued identically. In practice,

consumers choosing plans in 2015 weight premiums 6.5 times more than expected plan out-of-pocket spending.

A number of other results in the 2015 specification are at odds with standard economic models of choice. Even holding constant the individual’s own out of pocket risk, individuals have a strong distaste for higher deductibles and higher maximum out-of-pocket spending levels. Once those distastes are factored in, consumers then (i) have a preference for plans with lower actuarial values and (ii) are willing to pay more for plans with higher risk premia, both inconsistent with typical ‘homo economicus’ choice models.

How do consumer choices and associated preference estimates change when algorithmic decision support is introduced in 2017? The third column of Table 3 presents estimates for our primary choice model specification estimated for 2017.

A number of important results emerge. First, the implementation of algorithmic decision support entirely removes the large bias weighting premiums more heavily than out-of-pocket spending: the ratio of these coefficients in 2017 is approximately 1 to 1, as opposed to 6.5 to 1 in 2015. This is a substantial change, especially given that this ratio is shown in the literature to be robustly different than one across many choice settings in prior empirical work.

It is important to note that this change is not ‘mechanical,’ in the sense that it has to follow from the integration of decision support. As discussed in the model in Section 3.2, agents were not required to accept recommendations. They can consider both the algorithm’s recommendation, and whether or not to take it, and the non-financial plan dimensions not included in the algorithm. In our coming analysis, we show that agents/consumers continue to value non-financial dimensions of plans, even while making choices that are more consistent with a ‘homo economicus’ model in terms of financial plan dimensions.

The results in 2017 are also much more consistent with the standard economic model in a variety of ways. The coefficients on the deductible and maximum OOP are greatly reduced, as should be the case given the inclusion of individual out of pocket spending in the recommendation algorithm. The coefficient on actuarial value is right-signed and remains small. The risk penalty coefficient becomes negative, consistent with individuals preferring plans that are less risky all else equal. Taken together, these results show clearly that decisions improve greatly when focusing on the financial dimensions that the recommendation algorithm incorporates.

Figure 2 illustrates the magnitude of the monetary improvement in choices from 2015 to 2017. It plots three distributions. First, it plots the actual distribution of money left on the table due to 2017 plan choices. Next, it plots the predicted distribution of money left on the table in 2017 choices using estimates from the choice model estimated on 2017 choices. These two lines are essentially on top of each other, showing strong model fit. In addition, both lines show that there are still large sums of money left on the table in 2017 choices, though this could be because of preferences for non-financial plan dimensions, e.g., preferences for the Kaiser delivery model (discussed momentarily).

The third line plots the distribution of money left on the table for predicted 2017 choices using estimates from the 2015 choice model. This line reflects how agents/consumers choose in 2017 if they act like they did in 2015. The figure clearly shows that 2015 choice model parameters lead

	2015		2017
	(1)	(2)	(2)
Annual Premium (\$100)	-0.0746*** (0.00117)	-0.0984*** (0.00132)	-0.0633*** (0.00212)
Predicted OOP (\$100)	-0.0110*** (0.001)	-0.0151*** (0.00146)	-0.0721*** (0.00251)
Deductible (\$100)		-0.347*** (0.00704)	-0.0428*** (0.00606)
Max OOP (\$100)		-0.0384*** (0.000702)	-0.0129*** (0.00114)
Risk Penalty (\$100)		0.204*** (0.00497)	-0.0579*** (0.00246)
Actuarial Value		-0.0118*** (0.00252)	0.0130*** (0.00296)
Network Coverage	0.0133*** (0.00072)	0.0192*** (0.000722)	0.0301*** (0.000754)
Plan Type			
HMO	-	-	-
PPO	0.974*** (0.0206)	1.181*** (0.0225)	1.200*** (0.0211)
Other	-1.672*** (0.114)	-3.070*** (0.0758)	-0.742*** (0.0957)
Brand			
Regional carrier	-	-	-
Aetna	0.792*** (0.031)	0.343*** (0.033)	0.240*** (0.0257)
Blue	1.129*** (0.0245)	0.980*** (0.0251)	0.0609* (0.0276)
Humana	0.708*** (0.0271)	0.985*** (0.0285)	-0.317*** (0.0321)
Kaiser Permanente	3.184*** (0.0327)	3.170*** (0.0388)	2.058*** (0.0474)
United	0.551*** (0.0305)	1.037*** (0.0325)	-0.362*** (0.0263)
Pseudo R-squared	0.135	0.171	0.13
Observations	385,883	385,883	337,198

Standard errors in parentheses, * p<0.05 ** p<0.01 *** p<0.001

Table 3: This table presents the estimates from our main structural choice models.

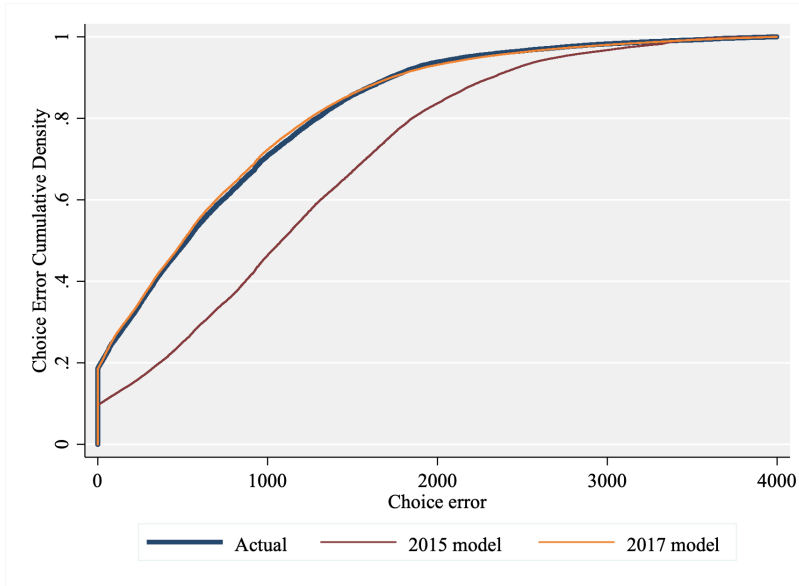


Figure 2: Observed 2017 choice error compared to simulated choice error

to substantively different outcomes with greater sums of money left on the table. The share of enrollments in plans with zero or near zero money left on the table is substantially lower. In 2015, 9.8% of people enrolled in the plan with the lowest expected cost compared to 18.0% in 2017, and only 24.2% of people enrolled in a plan that was within \$500 of the lowest expected cost available in 2015 compared to 47.4% in 2017.

Table 4 presents key related statistics. First, the table shows that the average actual money left on the table in 2015 was \$1,261, as compared to \$895 in 2017. The money left on the table thus went down by \$365 after decision support was in widespread use. To make this comparison ‘apples to apples’ we compare mean 2017 money left on the table under actual choices to counterfactual choices based on 2015 choice model parameter estimates. We find that now the average gain in terms of money left on the table from AI-based decision support is \$278. Table 4 also presents distributional statistics that for these key quantities showing meaningful variation in the differential money left of on the table pre and post decision support.

4.1.1 Robustness: Measurement Error

One potential concern with our analysis in this section is that our 2015 estimates are the result of agents measuring consumer plan-specific out-of-pocket spending with error, rather than jointly inconsistent decisions (on financial dimensions) between agents and customers. It seems clear that our results are not caused solely by econometric measurement error in plan-specific out-of-pocket spending. If that kind of measurement error were important, we would not find a major improvement in the weighting of premiums and out-of-pocket spending after the widespread decision support was introduced in 2017. Instead, we would find the same over-weighting of premiums in 2017.

Table 4: Actual and counterfactual choice error, full sample

	Mean	Percentile				
		10	25	50	75	90
Full Sample						
2015 Actual	\$1,261	\$66	\$550	\$1,124	\$1,762	\$2,342
2017 Actual	\$895	\$0	\$101	\$549	\$1,190	\$1,878
2017 Error w/ 2015 Sim. demand	\$1,173	\$15	\$500	\$1,083	\$1,712	\$2,331
2017 Error w/ 2017 Sim. demand	\$895	\$0	\$90	\$528	\$1,150	\$1,937
2015 Act. - 2017 Act.	\$365	\$66	\$449	\$575	\$571	\$464
2015 Sim. - 2017 Act.	\$278	-\$1,009	-\$68	\$389	\$1,104	\$1,738
2017 Sim. - 2017 Act.	-\$1	-\$1,134	-\$434	\$0	\$416	\$1,099

* Error differences >0 indicate reductions in cost error.

However, our estimated coefficients could reflect a combination of (i) agent/consumer measurement error in estimating out-of-pocket spending and (ii) jointly inconsistent decisions in valuing a dollar of premiums versus out-of-pocket costs. In this case, the removal of the premium versus out-of-pocket coefficient wedge between 2015 and 2017 cannot be interpreted solely as improving choices through reducing inconsistency – it may instead simply reflect a more mechanical reduction in measurement error in out-of-pocket estimation. To assess the importance of these mechanisms, we perform a series of exercises.

We run a series of simulations that assume that agents have the normatively correct weights for premiums, out-of-pocket spending, and plan financial characteristics. We assume that preferences for things besides these financial components are as estimated in our primary specification. We then simulate 2015 choices under the following scenarios for agent beliefs about out-of-pocket spending:

1. **Baseline:** algorithm-predicted individual-plan specific out-of-pocket used in primary model
2. **Coarse rounding:** assume that agents round consumer-plan-specific out-of-pocket to the nearest \$500 increment. We also implement this for the nearest \$1,000 increment.
3. **Normal Noise:** assume that agents have normally distributed mean 0 noise around the algorithmic projection. We use individual-plan-specification normal distributions with standard deviations equal to 200, 500, and 1,000, 2,000, and 3,000 in five different implementations.

After we simulate choices in these scenarios, we estimate our primary structural model. Table 13 in the appendix reports the estimates for these specifications. We find that estimates based on the simulation with baseline out-of-pocket predictions yield estimates where agents value premiums and out-of-pocket spending similarly and do not place any additional weight on financial characteristics, both as expected. When we move to the simulations with coarse rounding for predictions of out-of-pocket spending, we find that measurement error from those predictions do not meaningfully alter

the estimates from the baseline scenario (i.e. premium and out-of-pocket equally weighted and no emphasis on additional financial characteristics).

For the specifications which add normally distributed noise with σ of 200, 500 or 1,000 (truncated to 0 from below) the coefficients remain similar to the ‘homo economicus’ parameters that the underlying simulations are based on. We also consider extreme noise increases to 2,000 and 3,000. Even in these extreme cases, while premiums are weighted more heavily than OOP costs, the ratio never rises above 2 to 1, well below our 2015 estimate. Moreover, the coefficients on other plan characteristics never rise to more than a small fraction of their 2015 values shown in Table 3. These results confirm that the initial wedge in the respective weights for premiums vs. expected out-of-pocket spending in 2015 is due to behavioral foundations that are different than agent measurement error of individual-plan-specific out-of-pocket spending.

4.2 Non-Financial Dimensions

Overall, the results of Section 4.1 suggest that consumers and their agents make much more financially sensible choices in 2017 than in 2015. But health insurance choices in our context, and many others, are not just about financial aspects. There are a variety of other attributes of insurance plans that matter to individuals in making these decisions. Indeed, a key concern with algorithmic decision-support tools is that they will lead agents and consumers to over-emphasize the plan attributes included in the algorithm but under-emphasize the welfare-relevant aspects excluded from the algorithm. In the notation of our model in Section 3.2, financial plan dimensions are the variables X_{kj}^i (observed by agents and included in algorithm) while non-financial dimensions such as brand and network are the variables X_{kj}^g (observed by the agents but not included in the algorithm).

Our results in Table 3, however, show that this does not appear to be the case in the setting we study. Non-financial aspects of preferences continue to be valued - and appear to be more appropriately valued - in 2017. First, the coefficient on network breadth remains similar, if somewhat larger, in 2017 to what it is in 2015, suggesting that this important attribute continues to be weighted heavily despite not being included in the recommendation algorithm. If anything, the increase in magnitude suggests the weight on network is more aligned with enrollee preferences for a broader network, all else equal. This implies that the weights agents and consumers place on unused observables (Θ_g) are clearly non-zero on this important domain.

Second, brand preferences for insurers who generally offer similar broad network PPO products are (i) lower in magnitude in 2017 relative to 2015 and (ii) are more similar to each other in 2017 relative to 2015. In particular, the ‘branded’ fee-for-service carriers, Blue Cross Blue Shield and United, are no longer much preferred to substantively similar but potentially less well known regional carriers in 2017, even though they were strongly preferred in 2015. It is, of course, possible that consumers are over-weighting financial characteristics relative to brand under decision support, if one thinks that the brands provide substantive value relative to one another. However, under the hypothesis that similarly structured plans provide similar value, it seems clear that algorithmic

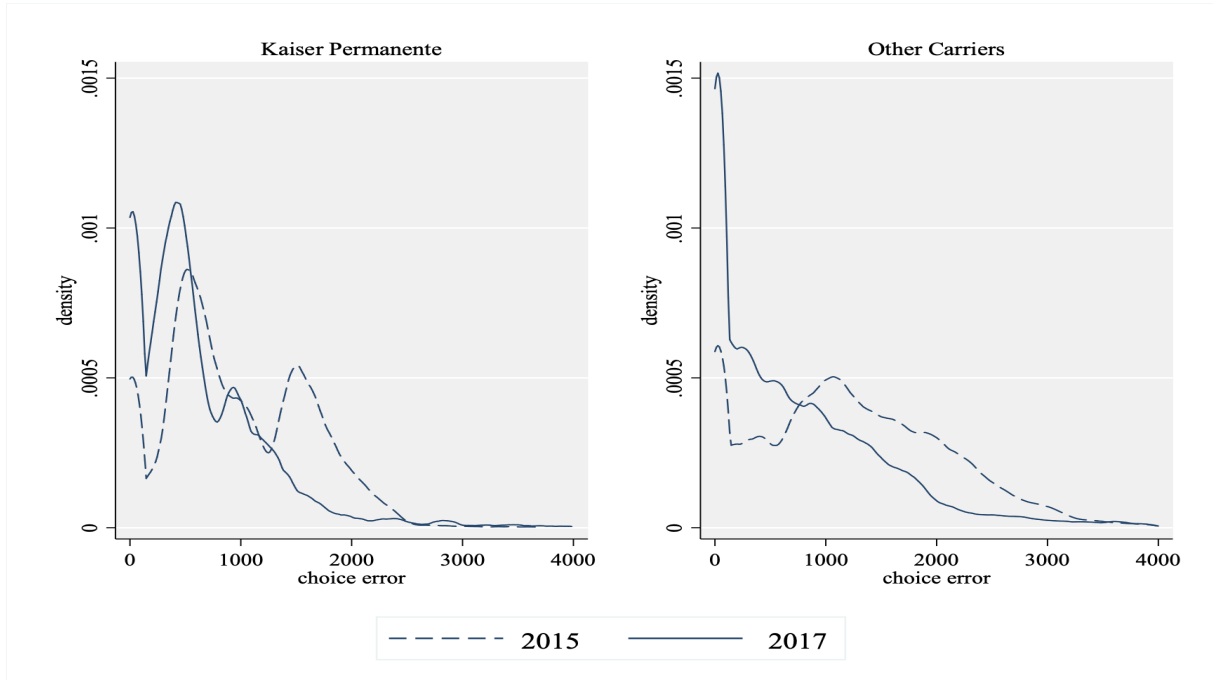


Figure 3: Choice Error by Plan Enrollment (Kaiser Permanente vs. Other) in 2015 and 2017

decision support reduces the emphasis on brands and increases the emphasis on financial value.

Another test of whether the algorithm moves consumers away from valuable brands is how preferences for Kaiser Permanente change with the introduction of AI-based decision support. Unlike the brands mentioned above, Kaiser has a separate network and an integrated care delivery model that differentiate its care substantively from other insurers.¹⁹

In 2015 there are high brand preferences estimated for Kaiser plans, an order of magnitude higher than for other brands. Once decision support is introduced in 2017, Kaiser brand preferences are reduced, but still high in value, reflecting the fact that agents and consumers do not blindly follow the decision support and pick one of the top recommended options if the agent/consumer has strong preferences for the non-financial aspects of Kaiser plans.

Without decision support in 2015, estimated brand preferences could reflect ‘true’ preferences that are above and beyond the expected financial outcomes under each plan. Alternatively, they could reflect agent and/or consumer biases and heuristics that can be overcome with sophisticated decision support. We investigate this in more depth by assessing whether or not the marginal switchers away from Kaiser are at the bottom end of the distribution for Kaiser plan financial value, as one would expect if agents are combining the algorithmic recommendations and information about consumer preferences for Kaiser in a sophisticated manner.

The left panel of Figure 3 presents the distribution of monetary loss for those enrollees choosing

¹⁹KP offers MA-PD plans that are available in a limited set of markets where KP operates. KP is a vertically integrated, closed network managed care plan. Despite the limited choice, KP has been demonstrated to provide high quality health care (see, e.g., McHugh et al. (2016)). These preferences, however, were not available in the algorithm and due to the financial structure of many of the plan offerings KP typically had a relatively high expected OOP cost as well as a low plan score.

Kaiser in 2015 and those choosing Kaiser in 2017. The right panel of the figure does the same for all other plans offered. In 2015 we see that both Kaiser enrollees and those in the PPO plans leave meaningful sums on the table. However, the Kaiser distribution is different in that it has two masses at both \$750 and again at \$1,600.

When decision support is integrated in 2017, the monetary loss for PPO plans uniformly shifts to the left, reducing the foregone savings associated with choosing these relatively homogeneous plans. For Kaiser, on the other hand, the large mass at \$1,600 has disappeared, but there remains a sizeable mass at a valuation of around \$500. That is, many individuals are willing to forgo significant savings to choose Kaiser, but those who were leaving the most on the table have been dissuaded. Indeed, among those who have a monetary loss of more than \$1,000 from choosing Kaiser, 61% do so in 2015 - but only 25% do so in 2017.

Taken together, this evidence suggests that (i) an agent/consumer pair is willing to overrule the algorithm if there is a meaningful consumer preference for Kaiser and (ii) the willingness of an agent/consumer pair to overrule the algorithm's recommendation depends on the magnitude of the loss — the cost of overruling. This is consistent with agents who integrate their information excluded from the algorithm (Kaiser brand preferences) with AI-based recommendations (cost error) efficiently by trading off at the margin.

Overall, the choice model estimates suggest that agents/consumers continue to value non-financial plan attributes that are likely welfare relevant, such as network breadth and the Kaiser delivery model, but do not continue to value non-financial aspects that are likely not welfare-relevant, e.g. the brands of relatively similar broad network PPO carriers.

4.3 Impact on Enrollee Experience

While our choice model estimates are strongly suggestive of the positive impacts that AI-based decision support has on plan choices, they are not dispositive on the normative implications of experts using AI. In this section we focus on the impact of decision support on subsequent plan turnover, a key ex-post measure of enrollee experience that is used for example in Ketcham et al. (2012) and Ketcham et al. (2015). If decision support provides people with valuable information that allows them to make better plan choices, we expect to see lower turnover rates among people who enrolled in recommended plans.²⁰

The effect is apparent in Figure 4, which shows 2018 switching rates for people who enrolled in MA-PD plans in 2017, based on the Score of their 2017 plan. Only 3.2% of people who enrolled in the top scoring plan in 2017 switched plans in 2018, compared to 6.8% and 8.4% of people who enrolled in Yellow (scores between 75 and 89) and Red (scores below 75) plans, respectively. If we look at switching among those who enrolled in a Green (scores of 90 or higher) plan we see switch rates of 4.0% compared to 7.1% among enrollees in lower color tier plans. Taken together, those

²⁰The high level of consumer inertia documented in the literature (see, e.g., Ericson (2014), Handel (2013), Polyakova (2016)) does not fundamentally change the link between turnover and satisfaction. Even if inertia causes fewer consumers to switch in general, plan switching is indicative of dissatisfaction with one's current plan. If anything, inertia implies that the gap in switching rates we document is a likely underestimate of what we might observe if consumers all made active (non-inertial) choices.

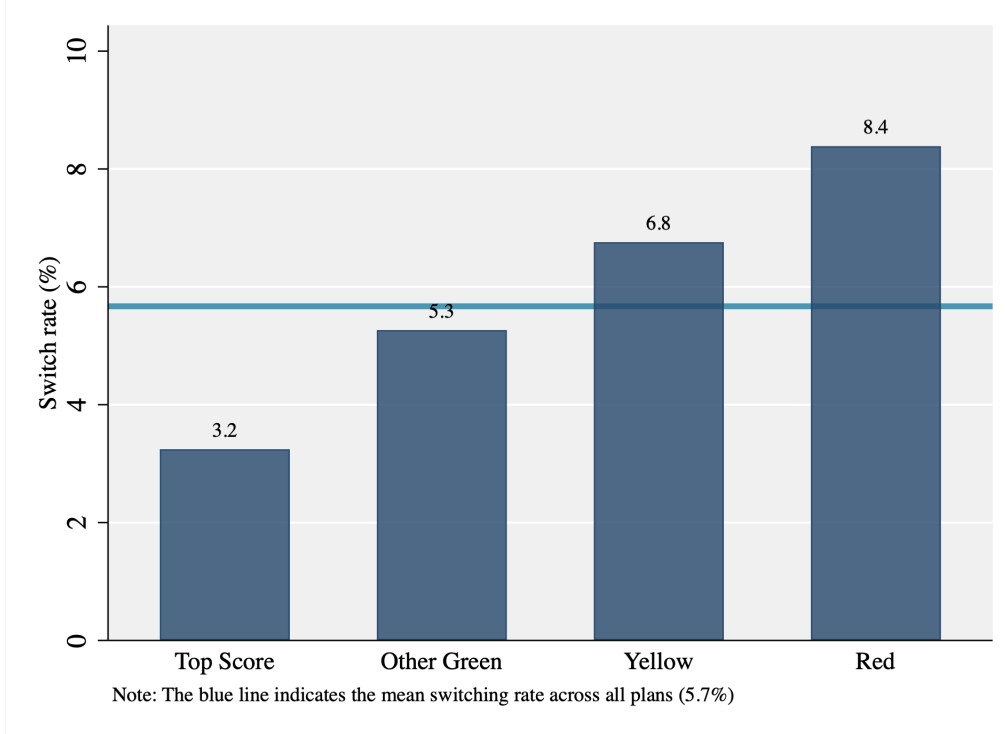


Figure 4: 2018 switching rates by 2017 plan score

following the AI recommendations were meaningfully less likely to switch than those who chose poorly rated plans.

One concern in interpreting these results is the endogeneity of the decision to take the AI-based recommendation. To address this issue we develop an IV strategy that takes advantage of random assignment of enrollees to agents. To do this, we estimate heterogeneity in predicted agent plan score and then investigate whether plan turnover directly relates to agent-specific fixed effects. This approach is comparable to a ‘judges design’ (see, e.g., Kleinberg et al. (2017)), and allows us to isolate the impact of a plan that would be scored higher by the AI-based tool solely due to being randomly assigned to a particular agent.

For our first stage of this analysis we use the following fixed effects specification:

$$PlanScore_{kj} = Age_k + Female_k + Cost_k + Agent_b + e_{kj} \quad (11)$$

Recall from our discussion in Section 2 that plan score is a measure of plan financial value ranging from 0 to 100, with 100 being the top end of the scale and 0 being the low end. Age_k and $Female_k$ indicate the age at the time of enrollment and whether individual k is female. We also include $Cost_k$, which assigns individuals to one of five quintiles based on predicted costs, in order to account for the possibility that higher cost individuals have more complicated cases. We recover the fixed effects for each agent indexed by b .²¹

²¹We perform this analysis only for agents with 20+ customers in 2017, to reduce concerns about statistical noise with the estimated fixed effects. See Section 5 for greater detail on our analysis of agent heterogeneity and see

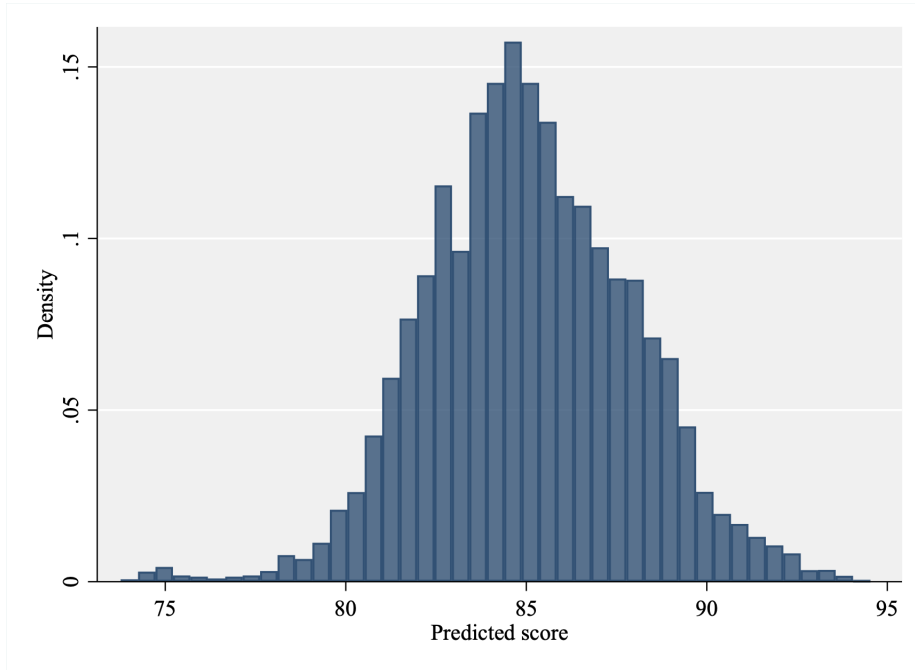


Figure 5: Distribution of predicted plan scores for 2017, across agents who work with over 20+ customers.

Figure 5 presents the results of this first stage, showing some meaningful variation in predicted 2017 Picwell score, which has a mean in this sample of approximately 84 and a standard deviation of approximately 3 (which is reasonably large since most scores are concentrated near the top end of the range 0 to 100). Thus, being assigned to agents one standard deviation apart from each other means that, on average, your expected plan score falls by 3 percentage points.

We then instrument for the enrolled plan score using the agent fixed effect estimate, essentially asking: if you are randomly assigned to an agent who is more tool compliant, are you more satisfied with your ultimate choice of plan? Table 5 presents IV estimates for the impact of a higher quality plan according to the AI tool on subsequent turnover. For an agent who is one standard deviation better (3 plan score points) than another in terms of predicted plan score, consumers are 7 percentage points less likely to switch plans the following year. This large effect is equivalent to nearly the full propensity to switch plans in the overall sample. This shows that customers who are randomly assigned to agents who recommend higher ranked plans are more likely to stick with those plans for multiple years, a strong indication that they are better off in those plans and that the tool is leading to better ex-post outcomes.²²

One concern with this analysis is that the agents who are most “tool compliant” in 2017 may be better for other, unobserved, reasons - so that it is not adherence to the tool which is lowering switching, but other aspects of agent behavior. Ideally we would address this by using our IV to

Appendix C for greater detail on our statistical mean reversion robustness analysis.

²²For additional robustness, we also estimate the same regression limiting the sample to the subset of enrollees who are 65. These “age-ins” have no inertia with respect to plans and are experiencing their chosen plan for the first time. The impact of plan score on turnover is nearly identical in this population.

Table 5: Switch IV

	2017 to 2018 Switching (1)
Agent Level Score	-0.0221*** (-0.0066)
Age Group	
<=65	-
66-70	-0.157*** (-0.0439)
71-75	-0.210*** (-0.0476)
76+	-0.206*** (-0.0443)
Brand	
Regional carrier	-
Aetna	0.837*** (-0.0569)
Blue	0.430*** (-0.0622)
Humana	0.183** (-0.07)
Kaiser Permanente	-0.272** (-0.102)
United	0.225*** (-0.0672)
Constant	0.175 (-0.569)
Observations	20,147

Standard errors in parentheses,
* p<0.05 ** p<0.01 *** p<0.001

predict switching rates for 2015 choices, as revealed in 2016, as a function of the instrument values from 2017. Unfortunately, the weak data available for 2016 means we can't rely on this specification test. Instead, we turn back to our ex-ante measures of foregone savings in Appendix Table 12.

To do so, we restrict ourselves to the set of agents who are in our data for both 2015 and 2017. We then estimate these same fixed effects model for this restricted set of agents in 2017, and use those fixed effects to instrument for ex-ante savings in both 2015 and 2017. Unsurprisingly, we find that those enrollees randomly assigned to agents more likely to follow the tool recommendations in 2017 have lower foregone savings in 2017. But we also show that enrollees using the same 2017-compliant agents when making their choices in 2015 have no differential foregone savings in 2015. If these agents were systematically “better”, we would expect it to show up in the ex-ante 2015 measure. Thus, it appears that tool compliance in 2017, and not underlying agent skill, is driving the increased enrollee satisfaction that we see in Table 5.²³

5 Agent Heterogeneity

Our results thus far show that (i) skilled agents alone are not sufficient to address choice inconsistencies with respect to financial plan dimensions, (ii) algorithmic decision support meaningfully improves decisions on the financial dimensions included in the algorithm, (iii) that money left on the table decreases with greater agent use of AI, and (iv) that agents/consumers still consider factors excluded from the algorithm when making recommendations after decision support. In this section we investigate the heterogeneous treatment effects of decision support, with a focus on effects by baseline agent skill (in terms of money left on the table). The correlation between these treatment effects and baseline skill bring direct evidence to bear on whether decision support is a complement or substitute for human skill.

We start by presenting estimates on heterogeneity in the quality of agents' recommendations prior to the introduction of AI-based decision support. We estimate the following simple model in 2015, prior to the introduction of decision support:

$$ChoiceError_{kj} = Age_k + Female_k + Cost_k + Agent_b + e_{kj} \quad (12)$$

This model is very similar to the fixed effects specification estimated in equation 11 in Section 4.3: the one difference is the dependent variable is now $ChoiceError_{kj}$, defined as the difference in expected total cost ($premium + E(OOP)$) for individual k who enrolls in plan j relative to the plan in their choice set with predicted lowest cost for that person. Under the assumption of random

²³We also perform an analysis, presented in the appendix, where we investigate how the worst agents in 2015 perform in terms of switching rates for 2018. Figure 13 shows that the worst quintiles of agents in 2015 in terms of foregone savings are more likely to have chosen plan with low algorithm plan scores in that year, by a large magnitude. We have shown in this section that from 2017 to 2018, there is much more turnover in the lower score plans. Figures 14 and 15 in the appendix show that (i) the worst performing 2015 agents choose plans of similar scores in 2017 to the best performing 2015 agents and (ii) that conditional on those plan scores, 2018 turnover is similar across the distribution of 2015 quintiles. This analysis also suggests that our switcher IV analysis results do not stem from persistent unobserved heterogeneity in agent performance.

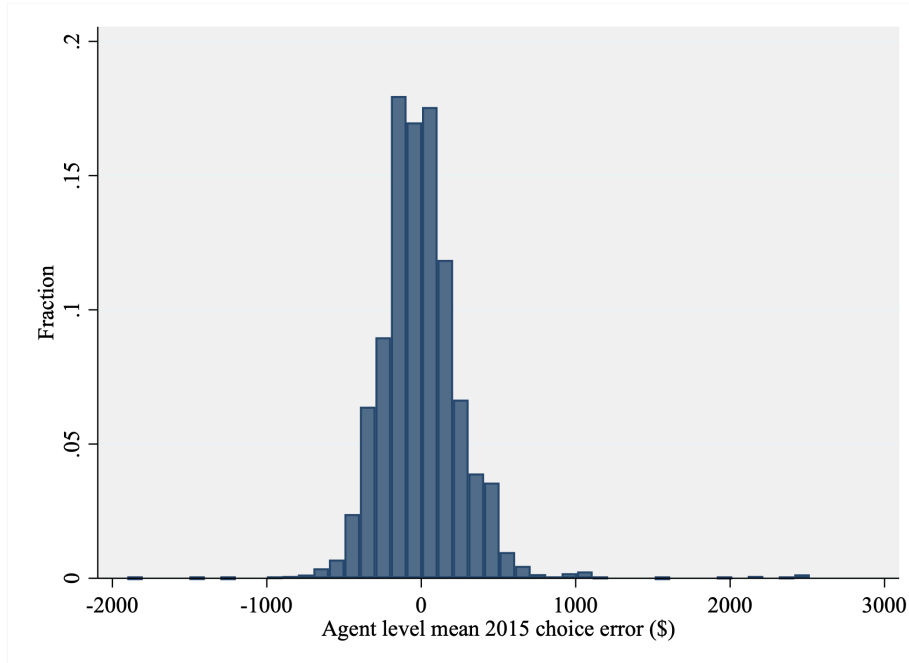


Figure 6: Agent fixed effects in 2015

assignment to enrollees these fixed effects represent the causal “value add” associated with each agent for choice quality, in terms of money left on the table. Figure 6 presents the distribution of fixed effect estimates.

The estimates in Figure 6 demonstrate that heterogeneity in agent skill has a material impact on the quality of health insurance choices. We see a mass of average agents but there are also many high-performing and low-performing agents in terms of their ability to match a particular enrollee to a plan in regards to money left on the table. Moving from the 25th percentile of the distribution to the 75th percentile improves choice quality by \$350 per enrollee per year, on average.

To better understand the nature of this heterogeneity in recommendations, we divide agents into quintiles of money left on the table based on the estimated fixed effects that we show in Figure 6. Table 1 presents agent and customer characteristics by agent skill. As we would expect to see when assignment is random, we do not see significant differences in the mean characteristics of customers across the different agent skill levels. The percent of female customers, age and number of prescriptions are all similar.

Table 6 presents the average money left on the table and call times for each quintile. Despite the balanced characteristics of consumers across quality quintiles, we see a systematic difference in choice quality across the quintile groups in 2015. For agents in the top quality quintile, consumers lose an average of \$893 per year in their chosen plan relative to the best plan for them in terms of expected financial outcome. The analogous amount for the lowest quintile of agents is \$1,734. Thus, the difference in being randomly assigned to a low quality agent as opposed to a high quality agent is almost \$1,000 per year in expected spending. Figure 7 presents the kernel density plot of choice error for each quintile group. For all groups we see substantial variance in money left on the

Table 6: Summary of 2015 agent level choice error, call time and tenure by quality quintile

Agent quality quintile	Choice error		Call time		Yrs. Experience	
	Mean	SD	Mean	SD	Mean	SD
1	\$893	\$786	54.3	33.4	4.2	0.80
2	\$1,086	\$848	53.2	35.3	4.1	0.79
3	\$1,213	\$1,000	48.6	35.6	4.2	0.83
4	\$1,369	\$1,237	53.6	35.2	4.1	0.81
5	\$1,734	\$2,452	56.7	38.1	4.0	0.79

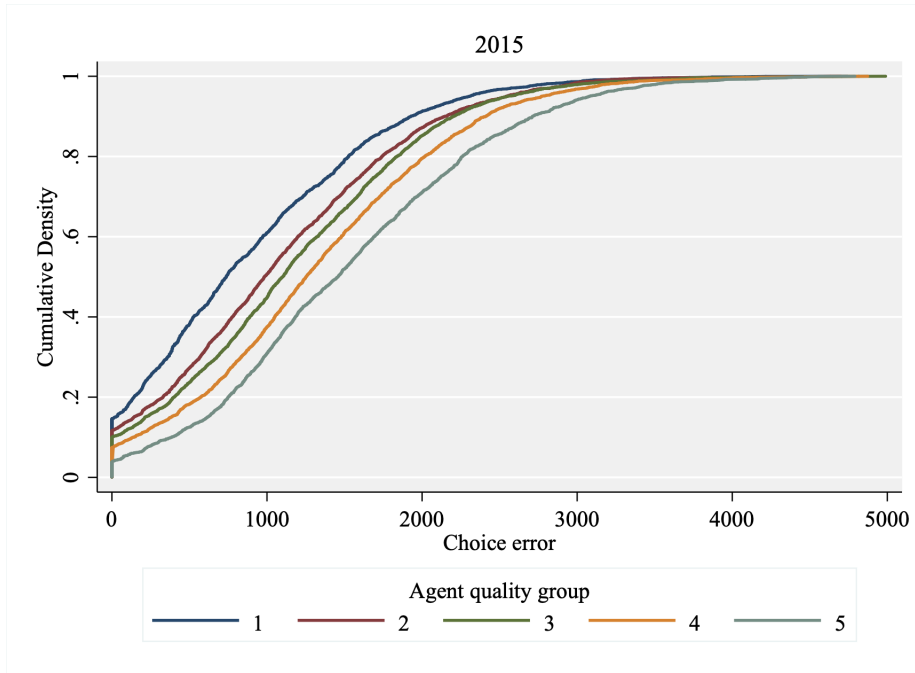


Figure 7: Distribution of 2015 choice error by agent skill quintile

table, but that higher skill agents have distributions that stochastically dominate those of lower quality agents, at each quality level.

Table 6 also shows that (i) agents have the same average tenure across the quality quintiles and that (ii) agents have the same mean call times with customers across the quality quintiles. The former suggests that learning from experience is not a key determinant of quality while the latter shows that agent/consumer call time, one granular measure of agent effort, is not a strong predictor of agent skill.

We use estimates from our structural choice model to assess changes in agent skill over time. To this end, we stratify agents into baseline quintiles using fixed effect estimates from (12) and use the choice model parameters to simulate changes in money left on the table. We simulate 2017 choices using 2015 choice model estimates (including 2015 fixed effects) and compare to 2017 choices using 2017 choice model estimates. Recall that this approach controls for choice set changes over time,

which can meaningfully impact the results.

Table 7: Choice error changes by agent skill

agent	quality metric	mean	p10	p25	p50	p75	p90
1	2015 Sim. - 2017 Act.	-\$55	-\$1,543	-\$388	\$0	\$826	\$1,547
2	2015 Sim. - 2017 Act.	\$118	-\$1,210	-\$293	\$96	\$916	\$1,598
3	2015 Sim. - 2017 Act.	\$133	-\$1,232	-\$262	\$163	\$978	\$1,639
4	2015 Sim. - 2017 Act.	\$305	-\$1,008	-\$29	\$406	\$1,128	\$1,768
5	2015 Sim. - 2017 Act.	\$354	-\$996	\$0	\$501	\$1,255	\$1,898

Table 7 presents the results from this exercise, including mean and quantile effects by baseline agent skill, in terms of money left on the table. The results show that the best agents on this dimension in 2015 perform roughly the same in 2017, while the worst agents in 2015 perform much better in 2017 (saving an average of \$354 per enrollee through improved recommendations). These results reveal (i) meaningful improvement on average after decision support and (ii) large improvements at the low end of the quality distribution but no changes at the top of this distribution.²⁴

We also present a set of kernel density plots for choices in 2015 and 2017 in Figure 8 to show that the net impact of these differential changes in agent skill is that we see both higher and more homogeneous agent skill in 2017 compared to 2015. The left hand panel replicates Figure 7 above for comparison and the right hand panel presents the same plot for plan choices after the widespread adoption of AI in 2017.

The shift in the amount of money left on the table is striking. All distributions shift to the left — moving choice error towards zero. There does, however, remain a long tail of observed recommendation errors potentially reflecting underlying enrollee tastes or private information. The similarity of money left on the table across the different skill levels in 2017 is also striking when compared to 2015. The entire distribution sits on top of one another for all groups in 2017 but is clearly distinct in 2015.

We dive deeper into the heterogeneous impacts of decision support by investigating how some of the micro-foundations estimated in our choice model change as a function of baseline agent skill. Recall that we found systematic mis-weighting of premium relative to $E(OOP)$ in the demand estimates for 2015 (See Table 3). This kind of heuristic decision making was previously attributed, at least implicitly, to consumers choosing and enrolling in plans (see, e.g., Abaluck and Gruber (2016) and Abaluck and Gruber (2011)). Our analysis thus far shows that, even under the guidance of skilled agents, this mis-weighting remains. Here, we ask (i) to what degree does mis-weighting vary by baseline skill and (ii) does offering AI-based decision support correct this mis-weighting to

²⁴One concern is that agents simply improve over time due to the extra experience incurred between 2015 and 2017. This is unlikely to explain the differential changes to quality. First, agents already had an average of 4 years of prior experience in 2015. Furthermore, we find no significant differences in experience by quality quintile (See Table 6). Thus, as of 2015, there appears to be no relationship between experience and the quality of recommendations.

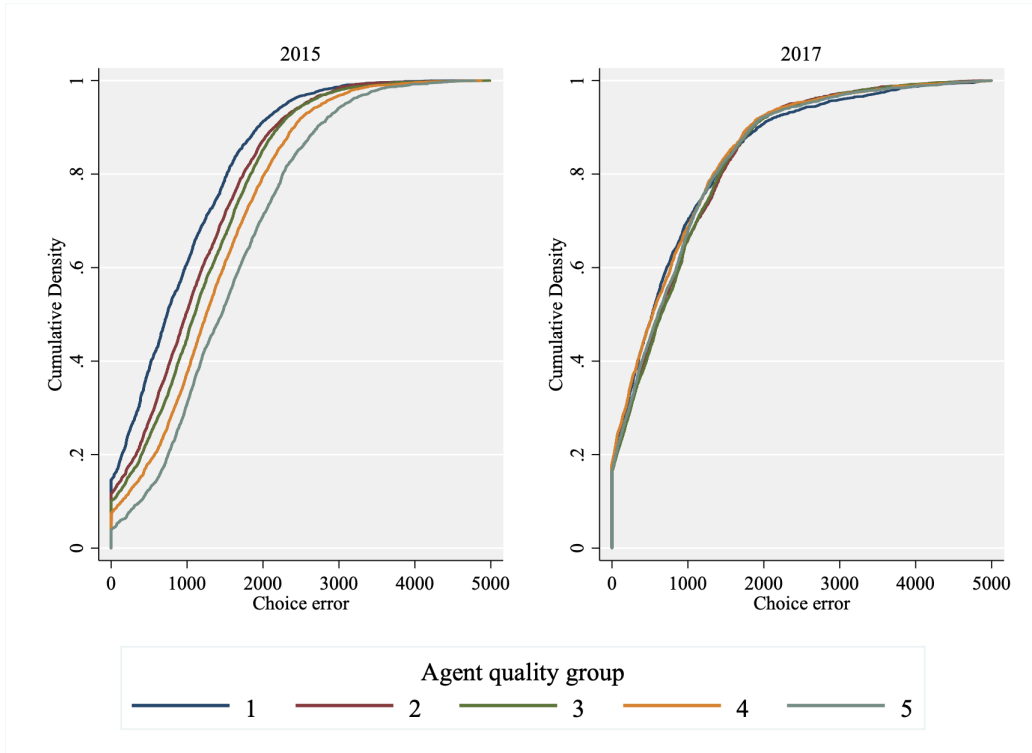


Figure 8: Choice error by agent skill in 2015 and 2017

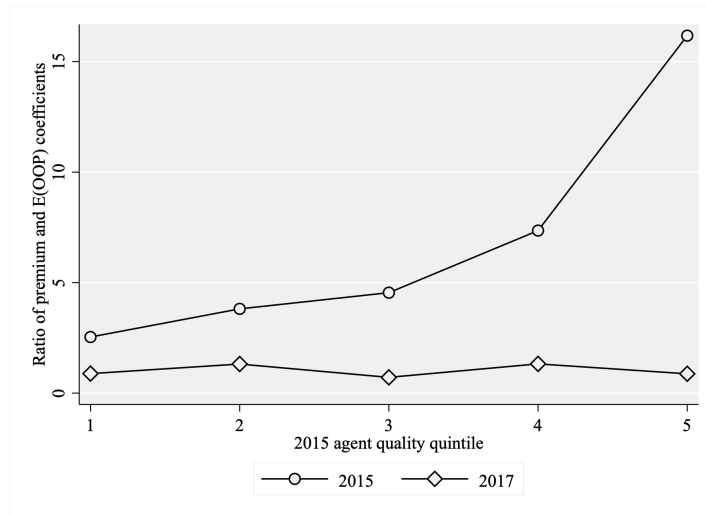


Figure 9: Ratio of coefficient estimates for premium and $E(OOP)$ by year

make lower skilled agents look more like their more highly skilled counterparts?

Figure 9 presents the ratio of weights on premium to $E(OOP)$ (rows 1 and 2 in Table 3) by baseline skill level. In 2015 we see a strong relationship between baseline recommendation quality and relative weights on premium and $E(OOP)$. The worst quintile of agent skill in 2015 has an average ratio of 16 to 1 while the best has a ratio of approximately 2 to 1. This ratio declines monotonically in baseline skill — agents who mis-weight premium relative to out-of-pocket cost

less make better recommendations prior to decision support.

Figure 9 shows that, on the financial dimensions, offering the AI-based tool harmonizes the recommendations of ex-ante different agents. In 2017, all quintiles have ratios near 1 to 1, reflecting a correct weighting of these financial factors. Introducing AI makes the financial recommendations from ex-ante low quality agents similar to the financial recommendations from ex-ante high quality agents. Taken together with our earlier results, we find that the use of the tool improves quality across the board but is also a clear substitute for ex-ante expertise. We note that in some ways this result is mechanical. That is, if agents had varying skill at baseline and they follow the AI recommendations, we would expect to see this result. Nevertheless, it is informative to document empirically that all types of agents — those who are highest skilled/put the most equal weight on premium and E(OOP) and those at the other end of the spectrum — equally adopt the AI-based recommendations and, therefore, improve relative weights in their recommendations to the same level.

One potential concern with our primary findings about the heterogeneous effects of decision support is mean reversion induced by statistical noise. To address this we re-estimate our primary models only for a subset of agents who have substantially more customers in both years and show the results are unchanged. We also note that mean reversion would suggest symmetry in impact between those who appear to perform well (who should get worse) and those who appear to perform poorly (who should improve). We do not find such symmetry in practice. These results are presented in depth in Appendix C.

Our empirical strategy focuses on brand as an “unused observable” measure of agent weights on enrollee preferences. We have already shown how average brand preferences change as a result of decision support. We now turn to how agents heterogeneously emphasize different brands to consumers over time. This is useful both (i) to understand heterogeneous steering and (ii) as a specification check on our primary model, which estimates only a mean brand preference coefficient for each brand.

Figure 10 presents the ratio of the network breadth and different brand preference coefficients to the premium coefficient, as a function of the baseline quality quintile. This figure is similar in spirit to Figure 9 but focuses on the ratio of network and brand to premium preferences, as opposed to the ratio of premium preferences to expected out of pocket cost preferences. The first panel reveals that not only did the average weight on network change little once AI was introduced, this effect was similar across the distribution of agents by baseline quality. Comparing this panel to Figure 9 shows that the agents who improved by the most on the financial aspects of plan choice — the lower quality quintiles — did not differentially change the weight they placed on network breadth.

The remaining 5 panels in Figure 10, each focused on a specific brand demonstrate that, though average brand preferences change after decision support, these changes are not markedly different across baseline quality quintiles, for any of the brands studied. For example, the Kaiser brand preference decreases a little bit relative to the omitted category (regional carries), close to uniformly across the quality quintiles. United, BCBS, and Humana all lose essentially their entire brand advantage relative to the regional carries, across all the quality quintiles. This reveals that, on

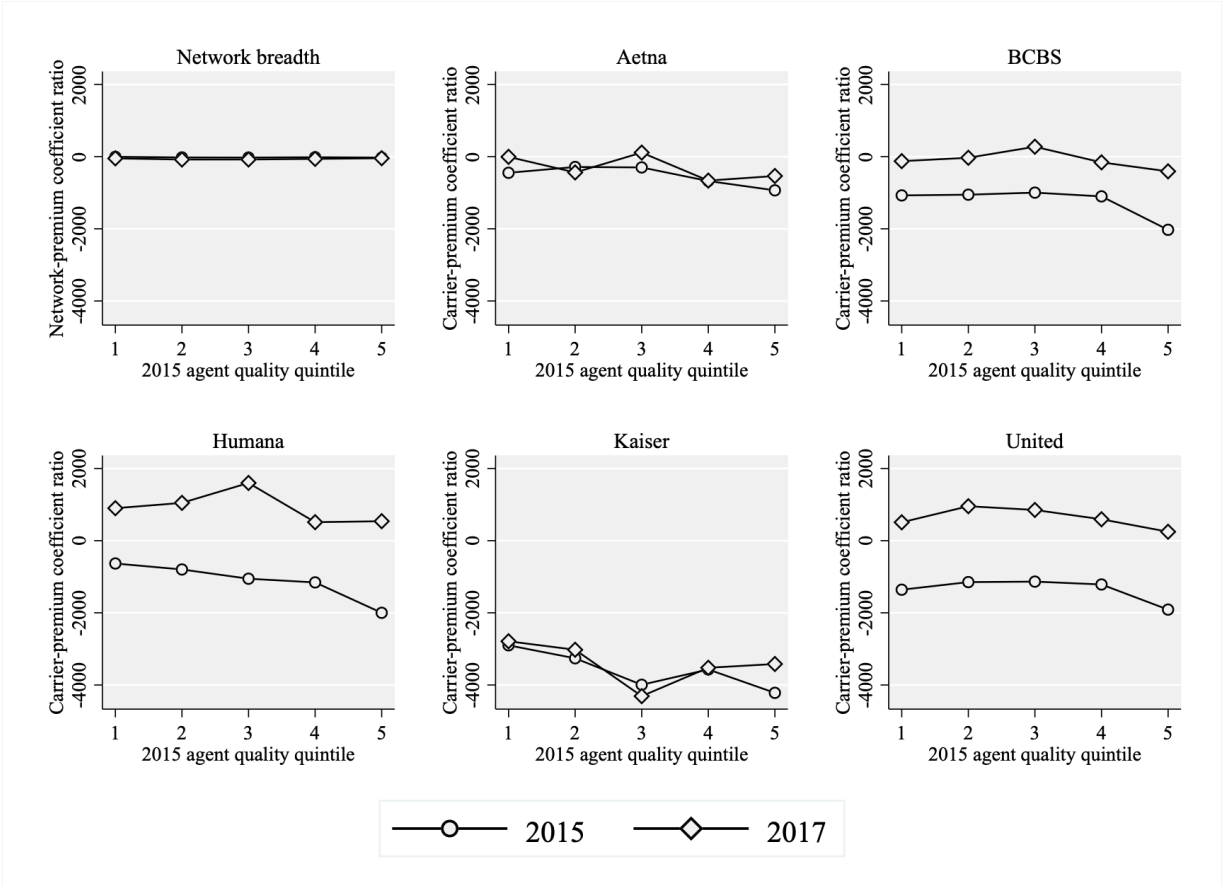


Figure 10: Ratio of coefficient estimates for network, brand preference and premium by year

average, the lower quality agents at baseline were not lower quality due to some clear ex-ante bias in favor of certain brands.

Importantly, the fact that shifts in brand choices/preferences occur, for each brand, across the entire distribution of agents suggests that the mean brand preference effects captured in Table 3 do a good job of capturing the change in brand preferences resulting from decision support roll-out.²⁵

5.1 Productivity

We also investigate agent productivity by using agent call times as a measure of agent effort. While call times could reflect a range of underlying agent-driven or consumer-driven inputs, and the correlation between call time and agent skill is ex-ante ambiguous, we think it is still instructive to assess (i) how call times change with the introduction of decision support and (ii) how call times correlate with the quality of choices pre- and post-decision support. Call times are also of particular interest as they represent the majority of the marginal cost of enrollment. Therefore,

²⁵We also include an exercise in Appendix D that plots the actual brand shares chosen by each agent and compares this with the amount they would have chosen that brand if they randomly chose across all options in their customers' choice sets. This exercise shows that, after decision support is fully integrated in 2017, the shift in brand shares changes across the entire distribution (high vs. low share of a given brand), i.e. the distribution has a level shift not a shape shift. See Figure 17 in Appendix D for more details on this exercise.

from a producer productivity perspective, call time relative to the quality of the plan chosen is a primary consideration.

Table 8: Average call time by agent skill level

Agent Quality	2015	2017
Average	53.27	41.90
1	54.27	38.85
2	53.18	43.64
3	48.59	42.68
4	53.63	41.61
5	56.69	42.74

Table 8 shows mean call times in 2015 and 2017 in the population as a whole and by 2015 quality level quintiles. Across the distribution of agent skill, call times are similar within both years. The average 2017 call time for agents in the top quality quintile are 1 to 3 minutes shorter than the average call times for agents in the lower quality quintiles, but, on the whole, call times are quite similar across the quintiles. This pattern was similar in 2015, prior to AI. Quality is not reflected in call time before or after the introduction of AI.

The main impact of AI was to reduce call times uniformly across the distribution. Average call time fell from 53 to 42 minutes, a reduction of 21% compared to the pre-AI average in 2015. This effect comes alongside an overall improvement in the quality of a call — measured by recommendations quality — and a convergence in agent quality after AI is available.

Combining the results, AI-based decision support allows the lowest skilled agents to make higher quality recommendations than the ex-ante highest skilled agents at significantly lower cost in terms of call time. Put differently, in the cross-section AI appears to be a substitute for skill.

6 The Problem of Improved Choice: Adverse Selection

In this section, we consider the inherent link between choice adequacy and adverse selection. Adverse selection can significantly reduce welfare in insurance markets, and as emphasized by Handel (2013), reducing choice inconsistencies could potentially worsen adverse selection. Whether the selection effect dominates gains from improved choices depends on a variety of setting specific factors and empirical estimates of the impact vary (e.g. Handel (2013), Polyakova (2016)). Handel et al. (2019) demonstrate the equilibrium welfare impact of choice improvements depend on a set of underlying micro-foundations that include consumer (i) costs, (ii) risk aversion and (iii) choice frictions. However, none of the existing papers actually show how reduced choice errors impact adverse selection in practice.

To understand the impact of AI-based decision support on adverse selection we implement a simple “correlation” test (see, e.g., Chiappori and Salanie (2000) or Einav et al. (2010)). We ask whether improved choices in 2017 lead to an increased correlation between risk and chosen plan generosity. We capture plan generosity in three ways. First, we rank plans by premium. A low premium rank implies a plan that provides less financial protection against out of pocket spending,

and vice versa. Second, we consider plans whose premiums are \$0, a common benefit design used to attract customers to plans with less generous coverage. Finally, we rank plans by the actuarial value of the plan overall.

The top panel in Figure 11 shows the correlation between health decile and the premium rank of the enrolled plan. We see an increase in the relationship between plan generosity and choice from 2015 to 2017. The average premium rank chosen is similar across health deciles in 2015. In 2017 it is upward sloping, implying a greater correlation between chosen plan generosity and expected health spending. Particularly striking are the results for the lowest cost decile: these healthiest enrollees were actually more likely to choose a high premium plan in 2015, and are much less likely in 2017.

The second panel in the figure shows the correlation between expected health risk and % of consumers who enroll in a \$0 premium plan. The third panel shows the correlation between expected health risk and the actuarial value rank of an enrolled plan. Both show evidence of increasing correlation between generosity and the health risk of enrollees choosing a plan. Despite the change in sign of the correlation, the relationship between risk and actuarial value rank is still relatively flat in 2017.

These results present evidence of an increase in adverse selection. However, we do not estimate the welfare impact of adverse selection in this marketplace in this paper for a number of reasons. First, the MA program includes robust risk adjustment, an important counter to adverse selection and complement to friction reducing policies such as AI-based decision support (Handel et al. (2019)). Second, the Exchange is only a small share of the MA market overall and in any particular county, which is the market-level at which prices are set. Even for pricing of plans to the exchange operator specifically, insurers offering an insurance product to the operating firm do so across a variety of markets (e.g. Medicare and commercial). Therefore, we do not expect dramatic shifts in plan offerings or pricing due to adverse selection in the MA marketplace alone nor can we easily identify the direct effects in our setting. Third, evaluating the welfare impact of improved choices in a marketplace requires moving beyond static evaluation of adverse selection in which current enrollee marginal cost and willingness-to-pay determine the optimum (e.g. Einav et al. (2010), Hackmann et al. (2015), Handel (2013)). In offering a marketplace to employers who want to cover retirees through the remainder of their lives, the exchange operator needs to provide insurance not only in a current year but against becoming sick: reclassification risk (Handel et al. (2015)). A model of such risk and the associated welfare impacts of information provision and adverse selection therein, while interesting, is beyond the scope of what we study. It does, however, represent an important future avenue for policy relevant research insurance market, particularly as AI-based tools become available.

7 Conclusion

We study insurance choice on a Medicare Advantage exchange platform where (i) consumers receive advice from agents, (ii) consumers are randomized to agents advising them, and (iii) the platform

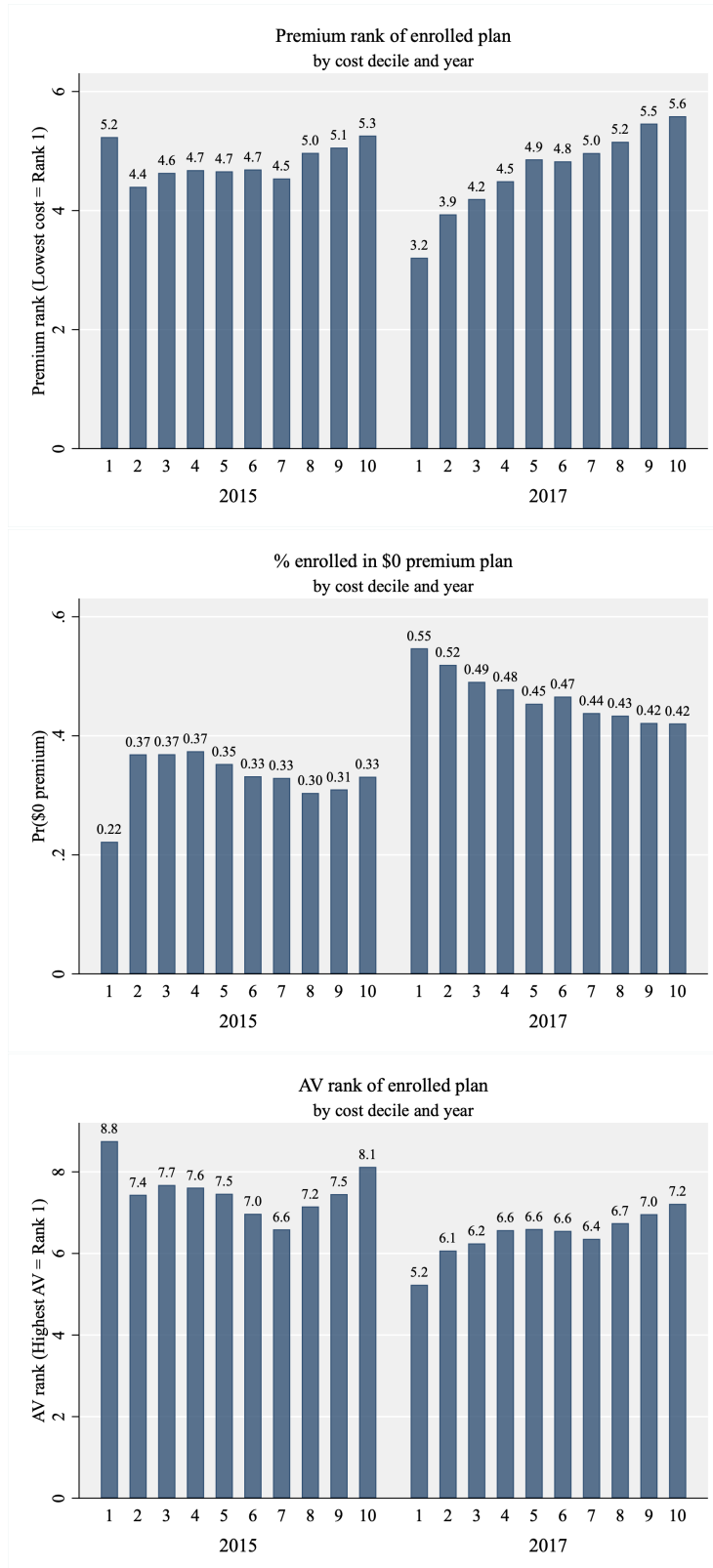


Figure 11: Relationship between choice and cost in 2015 and 2017.

fully integrated an agent-facing decision support tool over time. At baseline, we found that jointly made agent-consumer choices leave a lot of money on the table and exhibit some of the same biases found more broadly in the insurance choice literature, e.g. an emphasis on premiums at the expense of the more complicated to evaluate expected out-of-pocket spending. We find that the introduction of AI-based decision support improves decisions greatly on financial dimensions that the AI is well-suited to address, essentially removing the financial biases found at baseline. Importantly, we also show that (i) agents continue to integrate non-financial dimensions that are excluded from the algorithm into choices in a sophisticated manner and (ii) that consumers have better plan experiences after the introduction of AI-based decision support, as evidenced by reduced plan switching rates for enrollees randomly assigned to more tool-compliant agents and, therefore, choosing plans predicted by the AI tool to be a better match.

We also investigate agent heterogeneity, and find that the top performing agents at baseline, in terms of money left on the table, are helped a little bit by the algorithmic support but that the poorest performing agents are helped substantially, bringing their performance up to, and even slightly above, the level of the top agents at baseline. But we conclude on an important cautionary note, by providing the first direct evidence that addressing choice inconsistencies can lead to larger adverse selection in plan choice.

It is important to highlight some key limitations of our study. First, we are focused on one particular decision support tool, and we are unable to assess whether these results generalize to other forms of decision support. Our results can therefore be taken as demonstrating the potential for decision support to improve choices; future research could usefully assess a range of decision support tools in order to understand the benefits of alternative designs.

Second, in the health insurance literature specifically, quite a few studies have investigated consumer choices and interventions to improve those choices (Chandra et al. (2018)). Yet, there has been surprisingly little evidence of interventions that markedly improve choice quality. In our study, choices do change markedly and seemingly for the better. One key difference between our study and prior studies is that we combine two interventions (i) expert advising and (ii) sophisticated algorithmic decision support. One implication may be that multiple simultaneous interventions are needed to help consumers make better choices in this kind of market. Future work which can more usefully test that implication by considering in the same framework multiple combinations of expert advising and decision support.

Third, future work could usefully use a structural welfare framework to assess empirically this central trade-off highlighted in our findings: better choices (in terms of either reduced financial losses or higher consumer satisfaction) versus more adverse selection. Moreover, one could be concerned that, if decision support became especially prominent, that firms would respond to that by offering products that look best to the algorithm, at the expense of excluded dimensions. The ability of algorithms and advising to jointly integrate heterogeneous preferences and dimensions excluded from decision support are likely crucial determinants of whether such support helps markets function more fluidly or, instead, leads to market capture by firms that game the algorithm.

References

- Abaluck, Jason and Jon Gruber**, “Choice Inconsistencies Among the Elderly: Evidence from Plan Choice in the Medicare Part D Program,” *American Economic Review*, 2011, *101* (4), 1180–1210.
- **and Jonathan Gruber**, “Improving the Quality of Choices in Health Insurance Markets,” December 2016. NBER Working Paper No. 22917.
- **and –**, “Evolving choice inconsistencies in choice of prescription drug insurance,” *American Economic Review*, 2017, *106* (8), 2145–84.
- Acemoglu, Daron and Pascual Restrepo**, “Automation and new tasks: how technology displaces and reinstates labor,” *Journal of Economic Perspectives*, 2019, *33* (2), 3–30.
- Agrawal, Ajay, Joshua Gans, and Avi Goldfarb**, *Prediction machines: the simple economics of artificial intelligence*, Harvard Business Press, 2018.
- Athey, Susan, Kevin Bryan, and Joshua S Gans**, “The Allocation of Decision Authority to Human and Artificial Intelligence,” *Available at SSRN*, 2020.
- Bhargava, Saurabh, George Loewenstein, and Justin Sydnor**, “Choose to lose: Health plan choices from a menu with dominated option,” *The Quarterly Journal of Economics*, 2017, *132* (3), 1319–1372.
- Brown, Jason, Mark Duggan, Ilyana Kuziemko, and William Woolston**, “How does risk selection respond to risk adjustment? New evidence from the Medicare Advantage Program,” *American Economic Review*, 2014, *104* (10), 3335–64.
- Bundorf, Kate, Maria Polyakova, and Ming Tai-Seale**, “How do Humans Interact with Algorithms? Experimental Evidence from Health Insurance,” Technical Report, National Bureau of Economic Research 2019.
- Chandra, Amitabh, Benjamin Handel, and Josh Schwartzstein**, “Behavioral Economics and Health Care Markets,” August 2018. UC Berkeley working paper.
- Chiappori, Pierre and Bernard Salanie**, “Testing for Asymmetric Information in Insurance Markets,” *Journal of Political Economy*, 2000, *108*, 56–78.
- Diamond, Peter A**, “A model of price adjustment,” *Journal of economic theory*, 1971, *3* (2), 156–168.
- Dranove, David and Mark A Satterthwaite**, “Monopolistic competition when price and quality are imperfectly observable,” *The RAND Journal of Economics*, 1992, pp. 518–534.
- Egan, Mark**, “Brokers versus Retail Investors: Conflicting Interests and Dominated Products,” *The Journal of Finance*, 2019, *74* (3), 1217–1260.

- , **Gregor Matvos, and Amit Seru**, “The Market for Financial Adviser Misconduct,” *Journal of Political Economy*, 2019, 127 (1), 233–295.
- Einav, Liran, Amy Finkelstein, and Jon Levin**, “Beyond Testing: Empirical Models of Insurance Markets,” *Annual Review of Economics*, 2010, 2, 311–336.
- Ericson, Keith**, “Market Design When Firms Interact with Inertial Consumers: Evidence from Medicare Part D,” *American Economic Journal: Economic Policy*, 2014, 6 (1), 38–64.
- Fang, Hanming, Michael Keane, and Dan Silverman**, “Sources of Advantageous Selection: Evidence from the Medigap Insurance Market,” *Journal of Political Economy*, 2008, 116 (2), 303–350.
- Finkelstein, Amy and James Poterba**, “Testing for asymmetric information using “unused observables” in insurance markets: Evidence from the UK annuity market,” *Journal of Risk and Insurance*, 2014, 81 (4), 709–734.
- Gambacorta, Leonardo, Luigi Guiso, Paolo Mistrulli, Andrea Pozzi, and Anton Tsoy**, “The Cost of Distorted Financial Advice - Evidence from the Mortgage Market,” EIEF Working Papers Series 1713, Einaudi Institute for Economics and Finance (EIEF) 2017.
- Geruso, Michael and Timothy Layton**, “Upcoding or selection? Evidence from Medicare on squishy risk adjustment,” *NBER Working Paper*, 2015, 21222.
- Hackmann, Martin, Jonathan Kolstad, and Amanda Kowalski**, “Adverse Selection and an Individual Mandate: When Theory Meets Practice,” *American Economic Review*, 2015, 105 (3), 1030–1060.
- Handel, Benjamin**, “Adverse Selection and Inertia in Health Insurance Markets: When Nudging Hurts,” *American Economic Review*, 2013, 103 (7), 2643–2682.
- **and Jonathan Kolstad**, “Getting the Most from Marketplaces: Smart Policies on Health Insurance Choice,” *Hamilton Project Discussion Paper*, 2015.
- **and –**, “Health Insurance For Humans: Information Frictions, Plan Choice, and Consumer Welfare,” *American Economic Review*, 2015, 105 (8), 2449–2500.
- **and Joshua Schwartzstein**, “Frictions or mental gaps: what’s behind the information we (don’t) use and when do we care?,” *Journal of Economic Perspectives*, 2018, 32 (1), 155–78.
- , **Igal Hendel, and Michael Whinston**, “Equilibria in Health Exchanges: Adverse Selection vs. Reclassification Risk,” *Econometrica*, 2015, 83 (4), 1261–1313.
- , **Jonathan Kolstad, and Johannes Spinnewijn**, “Information Frictions and Adverse Selection: Policy Interventions in Health Insurance Markets,” *Review of Economics and Statistics*, 2019, 2 (101), 326–340.

- Heiss, Florian, Daniel McFadden, Joachim Winter, Amelie Wupperman, and Bo Zhou,** “Mind the Gap! Consumer Perceptions and Choices of Medicare Part D Prescription Drug Plans,” *Research Findings in the Economics of Aging*, 2010, pp. 413–481.
- Ho, Kate, Joseph Hogan, and Fiona Scott Morton,** “The Impact of Consumer Inattention on Pricing in the Medicare Part D Program,” *RAND Journal of Economics*, 2017, 48 (4), 877–905.
- Karaca-Mandic, Pinar, Roger Feldman, and Peter Graven,** “The role of agents and brokers in the market for health insurance,” *Journal of Risk and Insurance*, 2018, 85 (1), 7–34.
- Ketcham, Jonathan, Claudio Lucarelli, Eugenio Miravete, and Christopher Roebuck,** “Sinking, Swimming, or Learning to Swim in Medicare Part D?,” *The American Economic Review*, 2012, 102 (6), 2639–2673.
- Ketcham, Jonathan D, Claudio Lucarelli, and Christopher A Powers,** “Paying attention or paying too much in Medicare Part D,” *American economic review*, 2015, 105 (1), 204–33.
- Ketcham, Jonathan, Nicholas Kuminoff, and Christopher A. Powers,** “Estimating the Heterogeneous Welfare Effects of Choice Architecture: An Application to the Medicare Prescription Drug Insurance Market,” October 2016. NBER Working Paper No. 22732.
- Kleinberg, Jon, Himabindu Lakkaraju, Jure Leskovec, Jens Ludwig, and Sendhil Mullainathan,** “Human Decisions and Machine Predictions*,” *The Quarterly Journal of Economics*, 08 2017, 133 (1), 237–293.
- Kling, Jeffrey, Sendhil Mullainathan, Eldar Shafir, Lee Vermeulen, and Marian Wrobel,** “Comparison Friction: Experimental Evidence from Medicare Drug Plans,” *Quarterly Journal of Economics*, 2012, 127 (1), 199–235.
- Lacetera, Nicola, Bradley Larsen, Devin Pope, and Justin Sydnor,** “Bid Takers or Market Makers? The Effect of Auctioneers on Auction Outcomes,” *American Economic Journal: Microeconomics*, 2016, 8 (4), 195–229.
- Mackowiak, Bartosz, Filip Matejka, and Mirko Wiederholt,** “Rational inattention: A disciplined behavioral model,” *Goethe University Frankfurt mimeo*, 2018.
- McHugh, M, L Aiken, M Eckenhoff, and L Burns,** “Achieving Kaiser Permanente Quality,” *Health Care Manage Review*, 2016, 3 (41), 178–188.
- Mullainathan, Sendhil, Markus Noeth, and Antoinette Schoar,** “The market for financial advice: An audit study,” Technical Report, National Bureau of Economic Research 2012.
- Newhouse, Joseph P, Mary Price, Jie Huang, J Michael McWilliams, and John Hsu,** “Steps to reduce favorable risk selection in Medicare Advantage largely succeeded, boding well for health insurance exchanges,” *Health Affairs*, 2012, 31 (12), 2618–2628.

Polyakova, Maria, “Regulation of Insurance with Adverse Selection and Switching Costs: Evidence from Medicare Part D,” *American Economic Journal: Applied Economics*, 2016, 8 (3), 165–195.

Starc, Amanda and Robert J Town, “Externalities and benefit design in health insurance,” Technical Report, National Bureau of Economic Research 2015.

A Decision Support Performance

The Medicare decision support technology studied here can generate cost predictions based on two different levels of customer data. The base version of the model predicts cost based on age, sex and prescriptions, and another version of the model predicts cost based on these base characteristics plus responses to utilization survey question that ask users to indicate the number of primary care visits, specialist visits and hospital admissions that they experienced in the preceding 12 months. The base version of the decision support model is what was used to generate recommendations during our study period, but we present out-of-sample (OOS) R^2 (generated using 5-fold cross validation) for both the base and enhanced version to demonstrate how performance changes as more information is available to the model and to allow comparisons with other models that include additional information.

Table 9: R^2 of AI-Based prediction models and CMS HCC models for Medicare non-drug spending

AI Prediction Model	
Age+Sex+Drugs	0.069
Age+Sex+Drugs+Utilization survey	0.105
CMS HCC model V21	
New enrollees: age, sex, disability, Medicaid enrollment	0.019
Non-institutional continuing enrollees: age, sex, disability, Medicaid enrollment, ICD-10 codes/HCCs	0.125

In Table 9, we compare the R^2 of both versions of these Medicare prediction models to different versions of the CMS HCC model. In both cases, we consider $OOSR^2$ for models that predict inpatient and outpatient costs, and exclude prescription drug costs. The base decision support model R^2 is about 3.5 times larger than the CMS HCC model for new enrollees, and the decision support model with additional survey questions explains a similar level of variance to the CMS HCC model for continuing enrollees. These values are low in absolute terms, reflecting the high level of variance in medical spending among the Medicare population, but the performance of the Medicare decision support tool relative to the CMS HCC models indicates that a machine learning based model can use limited and easy to collect information to provide useful information to consumers of health insurance (or their agents).

In an alternative setting — recommendations for employer sponsored insurance where prior claims data are available — the same decision support technology predicts total allowed costs, including inpatient, outpatient and prescription drug costs, with an R^2 of 0.32, which exceeds R^2 values of other models that utilize similar inputs in a recent Society of Actuaries review of risk-scoring models (cite Society of Actuaries, 2016).

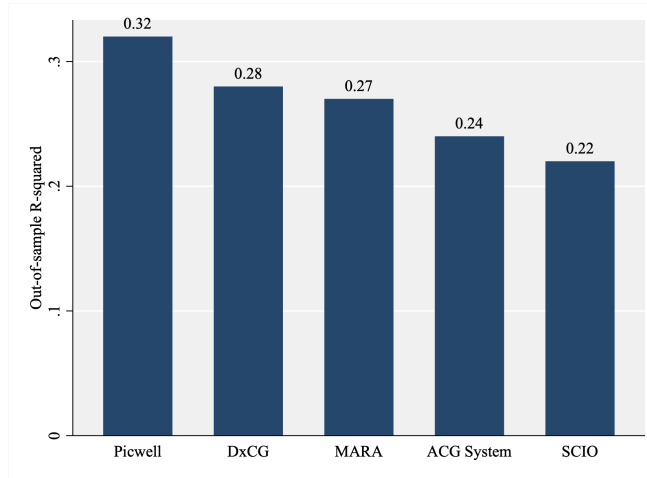


Figure 12: R^2 of health care cost prediction models that use prior claims data

B Additional Specifications

Tables 10 and 11 provide estimates for a version of our primary choice model that censors the top 5% of consumers with the largest estimated choice errors. The results in this specification are similar to those in our primary implementation discussed in the main text.

Table 12 presents the specifications for how foregone savings in 2015 and 2017 respectively are associated with 2017 mean agent plan score. This table is discussed in the main text and is consistent with the 2017 instrument for plan quality being associated with 2017 foregone savings but not 2015 foregone savings, implying that 2017 variation in mean plan score, after the widespread adoption of decision support, is not predictive of 2015 performance but is predictive of 2017 performance. This suggests that the plan switch IV regressions in the main text are not the result of unobservable agent heterogeneity that is correlated with their mean plan score.

Figure 13 shows that the worst quintiles of agents in 2015 in terms of foregone savings are more likely to have chosen plan with low algorithm plan scores in that year. We have shown in the main text that from 2017 to 2018, there is much more turnover in the lower score plans. Figures 14 and 15 in show that (i) the worst performing 2015 agents choose plans of similar scores in 2017 to the best performing 2015 agents and (ii) that conditional on those plan scores, 2018 turnover is similar across the distribution of 2015 quintiles. This analysis also suggests that our switcher IV analysis results do not stem from persistent unobserved heterogeneity in agent performance.

Table 10: Choice Model (censored sample)

	(1)		(2)		(3)	
	2015	2017	2015	2017	2015	2017
Premium (\$100)	-0.0876*** (72.70)	-0.0903*** (67.38)	-0.0874*** (70.55)	-0.0565*** (38.30)	-0.112*** (82.84)	-0.0838*** (37.68)
Predicted OOP (\$100)	-0.0308*** (38.93)	-0.0873*** (53.12)	-0.0312*** (34.79)	-0.140*** (77.24)	-0.0487*** (40.12)	-0.109*** (41.80)
Risk Penalty (\$100)			0.00354 (1.11)	-0.0777*** (50.04)	0.244*** (50.44)	-0.0455*** (18.38)
Actuarial Value					-0.0310*** (12.01)	0.00778* (2.57)
Deductible (\$100)					-0.347*** (47.17)	-0.0433*** (6.86)
Max OOP (\$100)					-0.0429*** (60.98)	-0.0163*** (13.90)
Network Coverage	0.0130*** (17.52)	0.0279*** (37.39)	0.0130*** (17.46)	0.0305*** (40.20)	0.0191*** (25.60)	0.0311*** (39.70)
Plan Type Dummies	X	X	X	X	X	X
Bran Dummies	X	X	X	X	X	X
Pseudo R-squared	0.152	0.132	0.152	0.154	0.192	0.156
Observations	363,367	317,701	363,367	317,701	363,367	317,701

t statistics in parentheses

* p<0.05 ** p<0.01 *** p<0.001

Note: This table estimates demand models on a censored sample in which we exclude individuals who are in the top 5% in terms of estimated choice error.

C Mean Reversion

To address the potential issue of mean reversion we consider two different types of analyses. First, we look at some of our key results conditioning on the set of agents who have a large number of consumers in our data. Figure 16 shows the ratio of our estimated premium coefficient to our estimated out-of-pocket spending coefficient conditioning on (i) agents who have more than 20 MA enrollees in each year and (ii) agents who have more than 50 MA enrollees in each year. As the number of enrollees per agent gets bigger, our results with our primary sample continue to hold, namely that (i) across the distribution of baseline quality quintiles this ratio moves to near 1 to 1 and (ii) the worst quality quintiles have much higher ratios in 2015. With 50+ MA enrollees per agent, it is highly unlikely that mean reversion from within-agent statistical noise impacts the results.

Table 11: Plan Type and Brand Coefficients (censored sample)

	(1)		(2)		(3)	
	2015	2017	2015	2017	2015	2017
Plan Type						
HMO	-	-	-	-	-	-
PPO	1.072*** (50.02)	0.904*** (48.56)	1.068*** (48.42)	1.158*** (57.25)	1.248*** (53.61)	1.226*** (55.62)
Other	-1.711*** (13.46)	-0.824*** (10.13)	-1.714*** (13.49)	-0.775*** (7.83)	-3.200*** (40.95)	-0.734*** (7.25)
Brand						
Regional Carrier	-	-	-	-	-	-
Aetna	0.859*** (26.74)	0.153*** (6.04)	0.859*** (26.69)	0.261*** (10.03)	0.371*** (10.69)	0.211*** (7.85)
BlueCross BlueShield	1.122*** (44.65)	-0.0454 (1.61)	1.125*** (44.33)	0.133*** (4.68)	0.963*** (36.96)	0.0837** (2.90)
Humana	0.714*** (25.63)	-0.621*** (20.12)	0.710*** (25.46)	-0.280*** (8.68)	0.947*** (31.84)	-0.240*** (7.24)
Kaiser Permanente	3.292*** (100.12)	1.935*** (41.85)	3.293*** (100.06)	2.275*** (48.59)	3.283*** (83.05)	2.203*** (44.98)
United	0.618*** (19.68)	-0.299*** (10.62)	0.618*** (19.67)	-0.314*** (11.45)	1.123*** (33.26)	-0.323*** (11.80)
Pseudo R-squared	0.152	0.132	0.152	0.154	0.192	0.156
Observations	363,367	317,701	363,367	317,701	363,367	317,701

t statistics in parentheses
 * p<0.05 ** p<0.01 *** p<0.00.1

Note: This table estimates demand models on a censored sample in which we exclude individuals who are in the top 5% in terms of estimated choice error.

Table 12: 2017 IV: Additional Analysis

	2017 Cost Error	2015 Cost Error
Agent Level Score	-68.73*** (6.128)	6.10 (6.149)
Age Group		
<=65	-	-
66-70	-43.32 (40.52)	34.84 (48.66)
71-75	53.27 (43.64)	155.25* (53.39)
76+	58.87 (38.87)	212.86*** (47.12)
Brand		
Regional carrier	-	-
Aetna	-160.86*** (45.72)	436.94*** (57.95)
Blue	-236.55*** (50.77)	372.23*** (60.46)
Humana	-238.49*** (52.82)	742.30*** (56.42)
Kaiser Permanente	-463.83*** (54.32)	-13.34 (46.42)
United	-207.15*** (50.60)	111.12* (53.08)
Constant	6961.19*** (524.4)	323.79 (534.9)
Observations	10,319	7,977

Standard errors in parentheses

* p<0.05 ** p<0.01 *** p<0.001

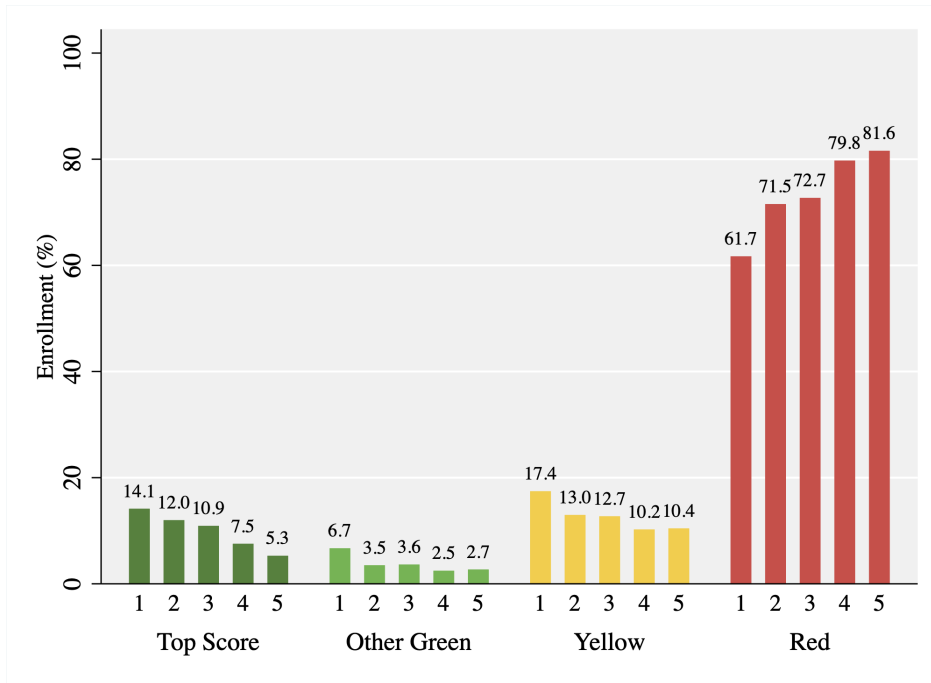


Figure 13: Choice of plan score (categorized by color tier) for 2015, as a function of agent 2015 mean foregone savings quintile (1 is worst, 5 is best).

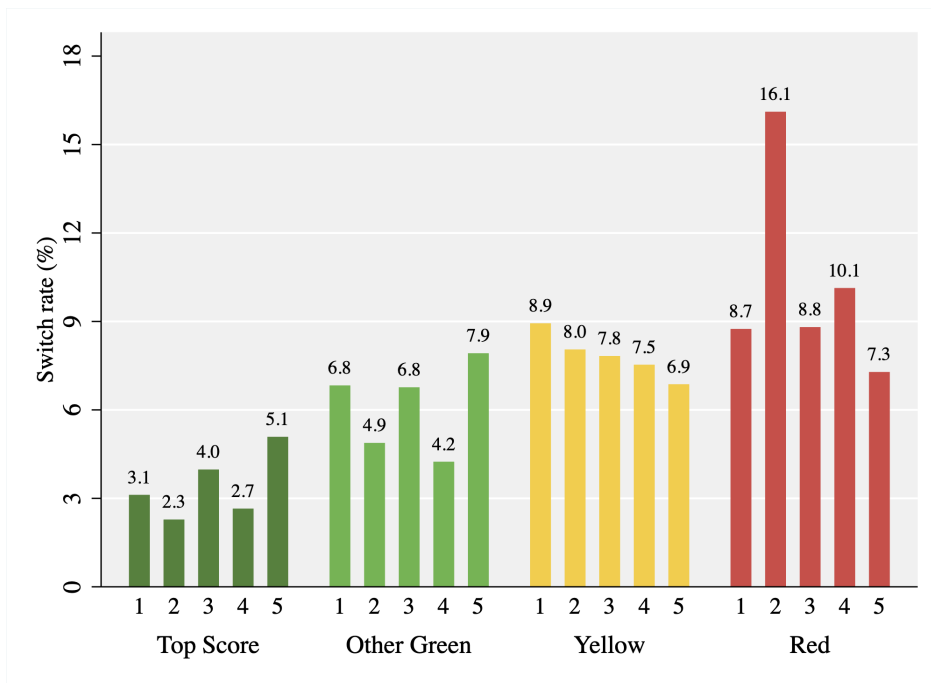


Figure 14: Choice of plan score (categorized by color tier) for 2017, as a function of agent 2015 mean foregone savings quintile (1 is worst, 5 is best).

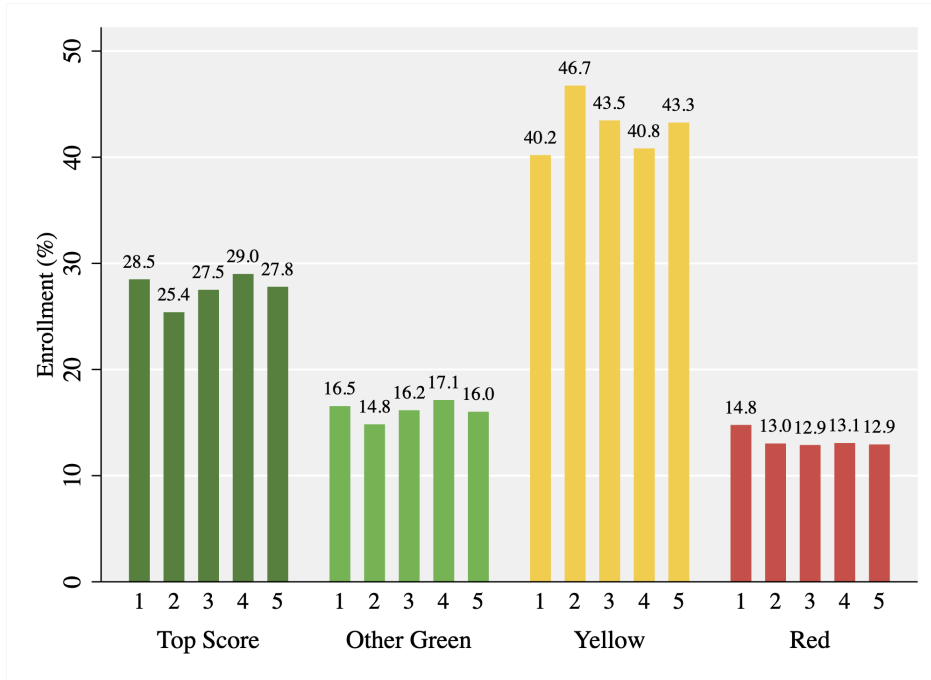


Figure 15: Plan turnover for 2018 as a function of 2015 agent mean foregone savings quintile (1 is worst, 5 is best) and color (score) of plan chosen in 2017.

Second, we note that mean reversion alone would suggest that the better agents in 2015 become worse in 2017, and vice versa. As shown in Figure 8 (as well as throughout our results) this is not the case. The best agents remain similar in 2017 to 2015, both in terms of average money left on the table and, somewhat in terms of their premium to expected out-of-pocket choice model coefficients. The worst agents improve markedly on both fronts. Overall, the distributions of money left on the table by baseline quality quintile are close to homogeneous in 2017 and look very similar to the distribution for the top quintile of agents in 2015. Also, importantly, these asymmetric changes by baseline quality quintile look the same when conditioning on agents with more than 50 consumers in each year. Overall, these results show that it is highly unlikely that mean reversion drives our results concerning the heterogeneous impacts of decision support.

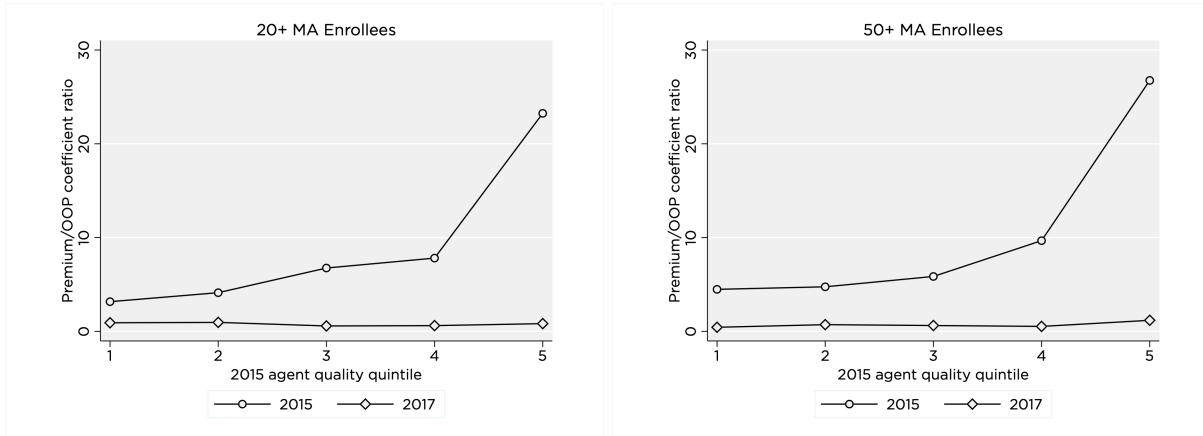


Figure 16: Ratio of coefficient estimates for premium to expected out-of-pocket spending for agents with larger number of MA enrollees.

D Additional Brand Preference Analysis

In the main text we discuss an analysis that compares brand choices over time, pre and post decision support. We compare what share of plans agents actually choose by brand, by year, to what they would have chosen if they randomly chose brands from the choice sets they engage with. To do this exercise, we:

1. Keep top 200 agents by volume in each year to get higher within-agent sample sizes.
2. For each agent in each year, for each brand, only keep consumers that actually had that brand in their choice set when calculating the share for that brand.
3. Compute the share that each agent would have in each brand if choices were made randomly. This controls for choice set size.
4. Subtract the randomly chosen share from the actually chosen share.
5. Plot the histogram of this statistic for all agents in the sample.

Positive values of the statistic for a given brand indicate a positive brand preference, relative to random choice, with negative values indicating a negative brand preference.

Figure 17 plots the results of this exercise for 6 brands, including all regional carriers together as one ‘brand.’ The results show that the changes to the average brand preferences estimated in our structural choice model come from shifts across the entire distribution of agents, rather than from shifts to specific agents who have very strong preferences for certain carriers. Each brand’s histogram reflects a level shift in the distribution to the left or the right, rather than a shift in the shape of the distribution. For example, the Kaiser distribution of shares chosen relative to random choice (i) has a similar shape in 2015 and 2017 (ii) is strongly positive in both years and (iii) is lower in magnitude in 2017 relative to 2015, reflecting the substitution on the margin away from

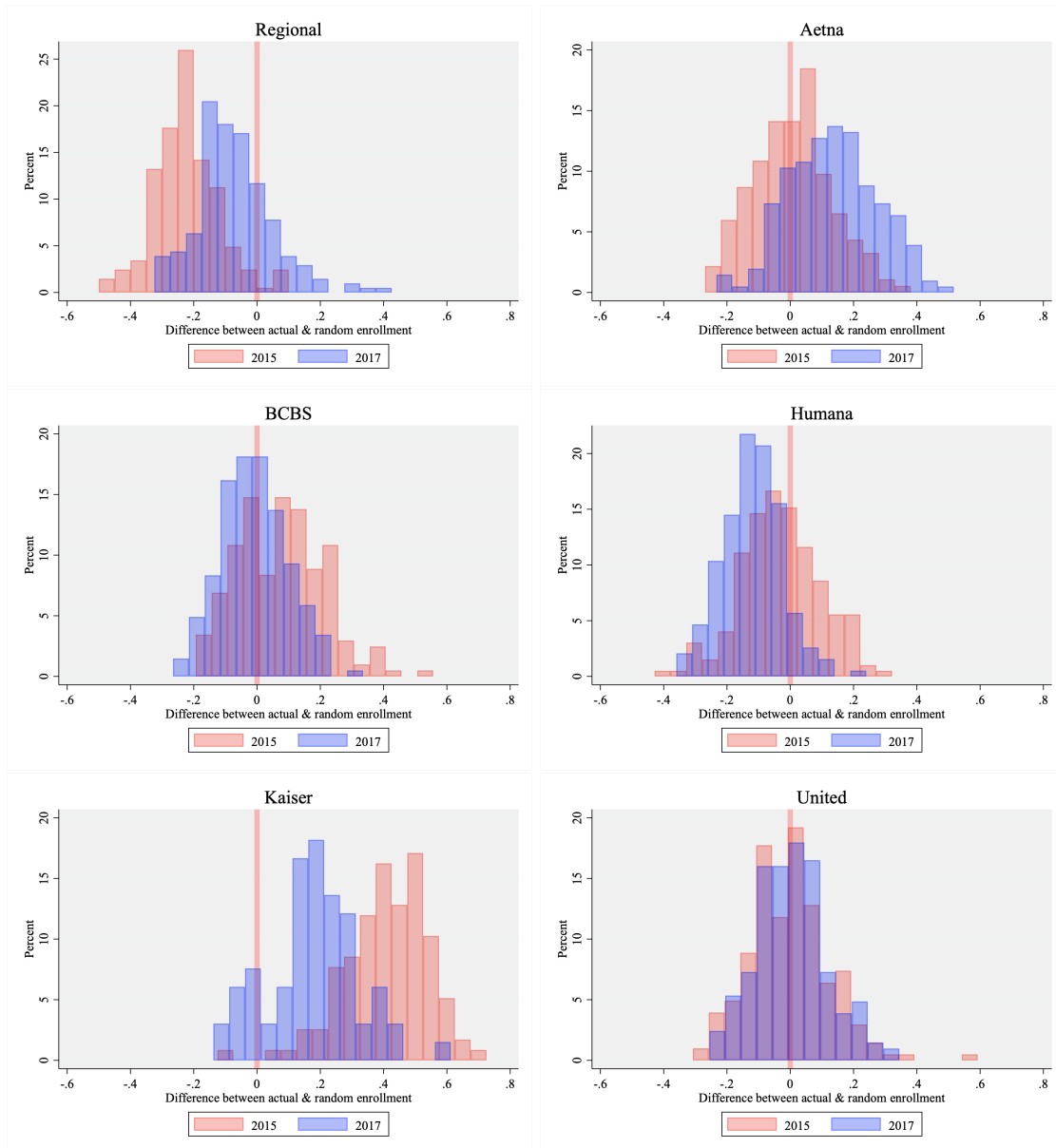


Figure 17: Agent Level Heterogeneity in Brand Choices

Kaiser by the consumers losing the most money from joining Kaiser. The distributions for BCBS and Humana clearly shift to the left, moving from positive on average to negative on average, while regional plans capture a lot of this lost brand equity, moving from quite negative in 2015 to near 0 in 2017. Thus, once decision support is used, the additional financial benefits from regional carriers as opposed to national carriers overcomes some of the brand effects for similar broad network PPO options that had been present in the market. Importantly, the fact that that the shifts occur across the entire distribution of brand preferences for each brand suggests that the mean brand preference effects captured in Table 3 do a good job of reflecting the shift in preferences of the population.

E Measurement Error

Table 13 reports the coefficients from the measurement error analysis. See Section 4.1.1 for a discussion of our approach to assessing measurement error and the implications of our results.

Table 13: Measurement Error Analyses

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
		Rounding		Noise (by Standard Deviation)				
	Picwell OOP	\$500	\$1,000	\$200	\$500	\$1,000	\$2,000	\$3,000
Annual Premium (\$100)	-0.0980 (0.00177)	-0.0975 (0.00177)	-0.0929 (0.00173)	-0.0971 (0.00176)	-0.0949 (0.00173)	-0.0883 (0.00169)	-0.0784 (0.00158)	-0.0724 (0.00156)
Predicted OOP (\$100)	-0.102 (0.00175)	-0.0984 (0.00178)	-0.0994 (0.00177)	-0.0977 (0.00176)	-0.0943 (0.00172)	-0.0808 (0.00161)	-0.0569 (0.00136)	-0.0413 (0.00120)
Deductible (\$100)	0.00465 (0.00605)	0.00808 (0.00601)	0.00169 (0.00602)	0.00653 (0.00603)	0.00130 (0.00600)	-0.0115 (0.00600)	-0.0139 (0.00595)	-0.00747 (0.00593)
Max OOP (\$100)	-0.000719 (0.000750)	-0.00143 (0.000747)	-0.000680 (0.000757)	-0.00130 (0.000750)	-0.000686 (0.000748)	-0.00129 (0.000749)	-0.000743 (0.000736)	-0.00138 (0.000745)
Risk Penalty	0.205 (0.00481)	0.202 (0.00485)	0.206 (0.00489)	0.208 (0.00487)	0.202 (0.00481)	0.196 (0.00480)	0.178 (0.00474)	0.165 (0.00473)
Actuarial Value	-0.0165 (0.00249)	-0.0149 (0.00248)	-0.0161 (0.00250)	-0.0147 (0.00248)	-0.0117 (0.00246)	-0.0144 (0.00247)	-0.0153 (0.00241)	-0.0102 (0.00239)
Network Coverage	0.0187 (0.000736)	0.0187 (0.000737)	0.0185 (0.000740)	0.0181 (0.000735)	0.0183 (0.000733)	0.0178 (0.000723)	0.0173 (0.000707)	0.0157 (0.000696)
Plan Type Dummies								
HMO	-	-	-	-	-	-	-	-
PPO	-3.115 (0.0357)	-3.039 (0.0348)	-3.101 (0.0353)	-3.043 (0.0347)	-2.947 (0.0336)	-2.794 (0.0317)	-2.584 (0.0291)	-2.433 (0.0278)
Other	-0.0899 (0.0627)	-0.119 (0.0623)	-0.195 (0.0632)	-0.162 (0.0632)	-0.0641 (0.0627)	-0.141 (0.0619)	-0.107 (0.0599)	-0.128 (0.0595)
Brand Dummies								
Regional carrier	-	-	-	-	-	-	-	-
Aetna	0.338 (0.0381)	0.293 (0.0385)	0.327 (0.0385)	0.329 (0.0382)	0.299 (0.0378)	0.281 (0.0371)	0.307 (0.0354)	0.282 (0.0351)
Blue	0.953 (0.0239)	0.946 (0.0240)	0.981 (0.0240)	0.954 (0.0240)	0.929 (0.0238)	0.911 (0.0234)	0.843 (0.0230)	0.769 (0.0228)
Humana	0.995 (0.0250)	0.999 (0.0250)	1.016 (0.0251)	1.010 (0.0250)	0.960 (0.0249)	0.915 (0.0247)	0.853 (0.0245)	0.810 (0.0241)
Kaiser Permanente	3.095 (0.0342)	3.093 (0.0340)	3.125 (0.0340)	3.082 (0.0340)	2.989 (0.0338)	2.845 (0.0339)	2.496 (0.0337)	2.407 (0.0334)
United	1.018 (0.0301)	1.025 (0.0300)	1.061 (0.0301)	1.038 (0.0300)	1.001 (0.0298)	0.970 (0.0294)	0.893 (0.0289)	0.860 (0.0286)
Pseudo R-squared	0.277	0.273	0.275	0.273	0.263	0.243	0.206	0.186
Observations	385,804	385,804	385,804	385,804	385,804	385,804	385,804	385,804

Standard errors in parentheses

F Using National Market Share Data to Assess Time Trends in Choice

One way to assess whether there was a more general shift towards better plan choices over the time period we study is to study patterns in Medicare Advantage choices nationally, focusing on a matched comparison group. To do so, we gather data on aggregate MA market shares for each plan, and key plan characteristics, for the same sample of counties that we use in our analysis, for 2015 and 2017. We obtain data on plan enrollment from the Centers for Medicare and Medicaid Services (CMS) for enrollment in MA plans as of June of 2015 and 2017. We also incorporate data on plan characteristics for those years from NBER and the Exchange.

In 2015, there were 14.96 million MA enrollees in plans in the counties served by the Exchange. The Exchange itself serves only 31,090 enrollees in those counties, or 0.21% of the sample. Thus, aggregate shares are essentially independent of the influence of our policy change.

One challenge is that we cannot replicate our full choice model in the aggregate data, since we do not have individual data to estimate expected out of pocket expenses. Instead, we estimate a multinomial choice model of the form:

$$choice_{ij} = \rho P_j + \phi AV_j + \beta n_j + \epsilon_i \quad (13)$$

where $choice_{ij}$ is a binary dependent variable that takes the value 1 if individual i is enrolled in plan j and 0 otherwise, for the set of plans that is available for individual i to choose from. P_j is the premium for plan j and AV_j is the estimated actuarial value of plan j . We include dummy variables to control for the type of plan (i.e HMO, PPO) and brand fixed effects. We additionally run a specification that excludes the percent of providers in-network.

To construct the national sample, we include all MA enrollees in counties served by the Exchange. We also include the enrollees in counties that were not served by the Exchange but that are demographically similar, which we identify using a propensity score matching procedure. Specifically, we match on the average age, gender, race, ethnicity, income, disability, and employment among the Medicare population in each county. Each county that is served by the Exchange (treated counties) is assigned a ‘nearest neighbor’ county that it is most demographically similar to from the group of counties not served by the Exchange (untreated counties). The propensity matching is only used to define which untreated counties should be included in the sample, but not to alternatively weight or oversample them. That is, in the event of several treated counties matching to the same untreated county, we include the untreated county in the sample, but do not assign it a greater weight. The untreated counties that were not selected as ‘nearest neighbors’ to any of the treated counties are excluded from the sample on the basis of being demographically dissimilar. The demographic data are taken from the Integrated Public Use Microdata Series based on 2015 values, and are aggregated at the county level using the 2010 Census Tract to 2010 PUMA Relationship File from the US Census Bureau. We restrict the matching algorithm to only select counties as the closest match if they are not proxied by the exact same combination of Public Use

Microdata Areas as the treated unit.

The demographic composition of our sample is provided in Table 14. The counties served by the Exchange account for approximately 90% of overall MA enrollees, and so the demographics are similar between all counties and just the counties served by the Exchange.

Table 14: Demographics breakdown

	(1) All counties	(2) Counties in Exchange	(3) Matched counties	(4) Final sample
Average age	70.74	70.76	70.69	70.74
Average income	31,167.5	32,375.21	28,945.65	31,365.07
% Male	45.38%	45.09%	45.93%	45.34%
% Female	54.62%	54.91%	54.07%	54.66%
% White	85.75%	84.37%	88.8%	85.67%
% Black	8.95%	9.41%	7.9%	8.96%
% Other Race	5.29%	6.22%	3.31%	5.36%
% Hispanic	5.86%	6.24%	4.88%	5.84%
% Disabled	43.35%	42.28%	45.06%	43.1%
% Employed	14.57%	14.49%	14.59%	14.52%
No. Counties	3,010	1,635	1,068	2,703
No. MA Enrollees (2015)	16,624,852	14,956,195	1,006,457	15,962,652
No. MA Enrollees (2017)	18,843,230	17,001,380	1,142,396	18,143,776

Table 15 reports the results of the demand model estimated for the Exchange enrollees using the demand specification in equation 13. Table 16 presents the results of the demand model for the national sample. Comparing the two, it is clear that the impact of plan characteristics on plan choices evolves differently between enrollees in the Exchange and the national sample between 2015 and 2017. In particular, the impact of actuarial value on choice increases by an order of magnitude in the Exchange population. In contrast, between 2015 and 2017 in the Exchange the coefficients on actuarial value are virtually unchanged and, if anything, declining over the same time period. We can statistically reject a change in preferences for actuarial value in the national Medicare sample that we see in the data on the Exchange.

Table 15: Demand model for exchange enrollees

	(1) 2015	(2) 2017
Annual MA Premium (\$100s)	-0.0713*** (-0.00118)	-0.0773*** (-0.00126)
Picwell Estimated Actuarial Value	-0.00412* (-0.0019)	0.104*** (-0.00277)
Plan Type		
HMO	-	-
PPO	1.059*** (-0.0203)	1.295*** (-0.019)
Other Plan Type	-1.297*** (-0.0818)	-0.856*** (-0.0883)
Brand FE	YES	YES
Enrollees	31,090	27,739

Table 16: Demand model for national sample

	(1) 2015	(2) 2017
Annual MA Premium	-0.03675*** (0.008)	-0.06045*** (0.007)
Picwell Estimated Actuarial Value	0.05249** (0.021)	0.04084*** (0.015)
Plan Type		
HMO	-	-
PPO	-0.37633*** (0.110)	-0.27001*** (0.075)
Other Plan Type	-1.39260*** (0.159)	-1.08919*** (0.126)
Brand FE	YES	YES
Enrollees	15,962,652	18,143,776