

The Response to Dynamic Incentives in Insurance Contracts with a Deductible: Evidence from a Differences-in-Regression-Discontinuities Design*

Tobias J. Klein, Martin Salm, and Suraj Upadhyay[†]

July 2020

Abstract

We develop a new approach to quantify how patients respond to dynamic incentives in health insurance contracts with a deductible. Our approach exploits two sources of variation in a differences-in-regression-discontinuities design: deductible contracts reset at the beginning of the year, and cost-sharing limits change over the years. Using rich claims-level data from a large Dutch health insurer we find that individuals are forward-looking. Changing dynamic incentives by increasing the deductible by €100 leads to a reduction in healthcare spending of around 3% on the first days of the year and 6% at the annual level. The response to dynamic incentives is an important part of the overall effect of cost-sharing schemes on healthcare expenditures—much more so than what the previous literature has suggested.

Key words: Patient cost-sharing, health insurance, dynamic incentives.

JEL-classification: I13, H51.

*We would like to thank Sara Abrahamson and Joachim Winter as well as conference participants at the 2018 EuHEA conference in Maastricht, the 2019 Essen Health Conference, the 2019 iHEA conference in Basel, and the 2019 Meeting of the Health Economics Committee of the German Economic Association in Munich for helpful comments and suggestions. This paper is a follow up to work on the PACOMED project on patient cost sharing that was funded by the Netherlands National Institute for Public Health and the Environment (RIVM). There are no conflicts of interest.

[†]Klein: Tilburg University, Department of Econometrics and Operations Research; e-mail: T.J.Klein@uvt.nl
Salm: Tilburg University, Department of Econometrics and Operations Research; Upadhyay: Tilburg University, Department of Econometrics and Operations Research.

1 Introduction

Annual deductibles are a common feature in health insurance contracts in many countries, as well as in other types of insurance contracts. In the United States, 82% of employer-sponsored health insurance plans feature an annual deductible (Kaiser Family Foundation, 2019 Health Benefits Survey). In the Netherlands, annual deductibles are mandatory by law in health insurance contracts for adults. With deductibles, patients have to pay for a specific amount of healthcare out-of-pocket before insurance coverage begins, which gives rise to both static and dynamic incentives. Static incentives exist when patients have to make deductible payments for current healthcare use. Dynamic incentives exist when deductible payments for current healthcare utilization reduce deductible payments for healthcare utilization later in the year. Decisions about healthcare utilization will be affected by dynamic incentives if patients are aware of them, are forward-looking, and value future states sufficiently. Consider for instance the case of a patient with an expensive chronic disease, who will cross the limit for the annual deductible (almost) certainly. Such a patient will pay the full deductible amount anyway and therefore her current healthcare use should not (or should barely) respond to the deductible.

From a policy perspective, knowledge of whether and how individuals respond to dynamic incentives imposed by deductibles (or similar cost-sharing schemes) is crucial for the design of health insurance contracts: setting the amount of a deductible is an essential element in the design of many health insurance contracts, and understanding how people respond to insurance contracts with different deductible limits requires knowledge of how people respond to dynamic incentives. For example, forward-looking patients will respond less to a deductible than myopic (not forward-looking) patients, because forward-looking patients anticipate that higher out-of-pocket spending now might decrease future out-of-pocket spending while myopic patients do not. Related to this, it is important to understand whether different groups in the population respond differently to dynamic incentives because this will lead to differences in utilization even if healthcare needs are the same.

It remains a challenge to quantify the reaction to dynamic incentives in this context because doing so requires a setting where dynamic incentives vary while all other factors—including static incentives—are kept constant. Previous studies have taken two alternative approaches: a reduced-form approach, where quasi-experimental sources of variation are used to test whether individuals respond to dynamic incentives, or structural modeling, where the response to dynamic incentives is quantified using a fully specified structural model. Studies that use reduced-form methods (e.g., [Aron-Dine et al., 2015](#), and [Guo and Zhang, 2019](#)) deliberately abstain from interpreting the magnitude of the effect they estimate, as the estimation equation is not derived from a structural model. Studies that use structural models (e.g., [Einav et al., 2015](#), and [Dalton et al., 2020](#)) rely on functional form and distributional assumptions. An additional challenge researchers face when pursuing a structural approach is the difficulty in accounting for unobserved heterogeneity in healthcare needs and their persistence over time.

The aim of this paper is to quantify the response to dynamic incentives in the context of health insurance with a deductible. For this, we propose an approach that (i) combines advantages of the reduced form and the structural approach, and (ii) can be pursued in other institutional contexts, provided that healthcare usage data are available at a high frequency and deductible limits exogenously vary across multiple years. Related to (i), we combine the advantages of the reduced form and the structural approach by estimating a micro-founded reduced-form equation exploiting quasi-experimental variation. As a consequence of deriving the equation from a full structural model, we can interpret the magnitude of the effects we estimate in terms of utility parameters and (to some extent) perform counterfactual experiments. At the same time, by abstaining from estimating a full structural model and focusing on estimating one key parameter of one structural equation, we can make very explicit which variation we exploit. Related to (ii), our approach is applicable in many contexts, because it uses two standard features of health insurance contracts with a deductible. The first feature is that deductibles reset at the beginning of the year. The second feature is that dynamic incentives change across years, because the deductible limit changes. Changes in deductible limits across years are often set either by law (e.g. in the Netherlands) or by changes in employer policy for employment based health insurance in the United States (e.g. [Brot-Goldberg et al., 2017](#)). We show that the combination of these two features gives rise to a differences-in-regression-discontinuities design.

In our analysis, we proceed in two steps. In a first step, we follow individuals who have crossed the deductible in year y as they experience a reset in the deductible at the beginning of year $y + 1$. This means that they do not face cost-sharing incentives at the end of year y , but face them at the beginning of year $y + 1$. The change in care consumption is local and driven by the change of static *and* dynamic incentives, from the end of year y to the beginning of year $y + 1$. We repeat this approach for other year-pairs. Importantly, by estimating changes locally, around the turn of the year, we control for seasonality in healthcare needs and individual heterogeneity. Moreover, by design, the change in static incentives around the turn of the year is the same for all year-pairs. This means that *differences* in regression discontinuities are directly related to *differences* in dynamic incentives across years. Therefore, in a second step, we relate the sizes of the estimated discontinuities to a commonly used measure of dynamic incentives, the expected end-of-year price at the beginning of the second year of the respective year-pair. The average expected end-of-year price can be estimated as the fraction of individuals who have not hit the deductible limit by the end of the year. It varies across years due to changes in the deductible amount. A lower expected end-of-year price makes it more attractive to consume more care earlier in the year. Forward-looking behavior would imply a negative relationship between the changes of care consumption we estimate around the turn of the year and the expected end-of-year price at the beginning of the second year of the respective year-pair. To test this prediction, we apply both parametric and non-parametric statistical tests. Moreover, we estimate the size of the dependence between the discontinuities and our measure of dynamic incentives, and we

give these estimates a structural interpretation.

We implement this approach using administrative data from a Dutch health insurer for the years 2008 to 2015. The population of insured individuals for which we have data is broadly representative of the Dutch population. It is not limited to certain groups such as the elderly or employees of a specific firm. The Netherlands provides a favorable setting for our purpose. Health insurance coverage and the set of providers patients have access to is comparable across insurance companies. Furthermore, cost-sharing incentives are the same across insurance providers. Deductibles are mandatory for all health insurance contracts for adults, and the minimum deductible amount is set by the Dutch government each year. Deductible amounts are low, which means that dynamic incentives are important. At the same time, deductibles have increased substantially over our study period, from €150 in 2008 to €375 in 2015. This provides variation in dynamic incentives that we exploit with our differences-in-regression-discontinuities approach.

We find that individuals respond to the expected end-of-year price, in a way that is consistent with forward-looking behavior. An increase in the expected end-of-year price reduces daily expenditures at the beginning of the year. Increasing the deductible by €100 leads to a reduction of around 3.0% in daily healthcare expenditures at the beginning of the year when patients have not yet exceeded their deductible. Our results are robust to alternative specifications, and they cannot be explained by changes to the healthcare system other than changes in dynamic incentives.

We also explore whether the reaction to dynamic incentives differs across subgroups in our sample, defined by age, gender, and income. We control for differences in healthcare needs and find that almost every subgroup exhibits forward-looking behavior. The only exception is the group of individuals that are below age 45. For the remaining groups, we find the reactions are very similar. An increase in the size of the deductible by €100 leads to a percentage change in daily expenditures that ranges from -2.5% to -3.5%.

Finally, we show that our framework allows us to conduct interesting counterfactual experiments that incorporate prior knowledge about reactions to static incentives from the literature. We find that increasing the deductible by €100 leads to a reduction in per capita yearly expenditures by approximately €211, which is around 10.55% of total yearly expenditures that count towards the deductible. The reaction to dynamic incentives accounts for a similar share in this reduction as the reaction to static incentives. Generalizing our findings to the Dutch population would imply that the reaction to dynamic incentives alone leads to a reduction in expenditures by approximately €2 billion when the deductible size is increased by €100.

Our study relates to the literature on patients' responses to cost-sharing incentives (see surveys by [Cutler and Zeckhauser, 2000](#), [McGuire, 2011](#), [Einav and Finkelstein, 2018](#)). Earlier theoretical contributions have examined patients' responses to dynamic incentives under a health insurance contract with a deductible. They conclude that, under some assumptions, forward-looking individuals should only respond to the expected end-of-year price and not to static

incentives (Keeler et al., 1977, Ellis, 1986). This forms the basis for using the end-of-year price as a measure of dynamic incentives. We discuss this literature in detail in Appendix A, where we also provide the micro foundation for our analysis.

Several recent empirical studies examine whether and how patients respond to dynamic incentives in the context of patient cost-sharing. For Medicare Part D, some studies test the hypothesis of full myopia through the estimation of a discount factor; a discount factor of 0 would indicate full myopia. Einav et al. (2015) utilize the non-linearity in prices caused by the donut hole structure of Medicare Part D plans to estimate a weekly discount factor of around 0.96, or 0.12 at the annual level, rejecting the hypothesis of full myopia. Dalton et al. (2020), on the other hand, estimate a discount factor of 0—an indication of complete myopia. Similarly, Abaluck et al. (2018) find evidence for substantial myopia.

In the setting of employee insurance in the United States, Brot-Goldberg et al. (2017) find that for high deductible health insurance plans dynamic incentives play only a minor role in determining healthcare utilization. Similarly, Guo and Zhang (2019) study the spending patterns of individuals who have a large expenditure planned in the future (childbirth). Their results are consistent with individuals exhibiting myopic behavior. Aron-Dine et al. (2015) focus on employees that enroll into health plans in different months within the same year. Under the assumption that individuals who enrolled into these health plans in different months are comparable, they relate differences in healthcare utilization to the differences in the expected end-of-year price that are driven by differences in the time at which individuals enrolled during the year. Using a differences-in-differences approach, they reject the hypothesis of full myopia and conclude that dynamic incentives do matter for healthcare utilization.

We contribute in various ways to the literature. On the substantial side, we show that a broad population of individuals strongly react to dynamic incentives and that changing these incentives through increasing deductible limits has quantitatively important effects. Moreover, we show that the reaction to dynamic incentives is similar across groups in the population, with the exception of young individuals.

In addition, our study makes two methodological contributions. First, we show that a combination of a key feature of standard deductible contracts, namely that they reset at the turn of the year, and exogenous variation in the deductible amount across years give rise to a differences-in-regression-discontinuities design. We demonstrate how this can be used to estimate the effects of dynamic incentives on patient behavior.

Second, we explicitly derive our estimation equation as a reduced form from an economic model. This allows us to give a structural interpretation to the response to dynamic incentives at the beginning of the year that we estimate, and to perform counterfactual experiments. Moreover, we show that the expected end-of-year price is either a valid measure or a good proxy of dynamic incentives for a broad class of models that includes the model by Keeler et al. (1977) and more recent models with quasi-hyperbolic discounting, such as the model by Abaluck et al. (2018).

Our study continues as follows: we describe our empirical approach in Section 2. Then, in Section 3 we provide details on the institutional background and the data. We discuss our empirical implementation in Section 4. Results of our main analysis are presented in Section 5. In Section 6 we present robustness checks. In Section 7 we perform a counterfactual experiment and quantify the effect of an increase of the deductible on annual expenditures, and the contribution of dynamic incentives to this effect. Section 8 concludes.

2 Empirical approach

2.1 Financial incentives and care consumption

The primary aim of this paper is to measure patients’ responses to dynamic incentives. We do so in the context of a deductible contract for health insurance, where a patient pays for the first euros of care consumption herself and faces no out-of-pocket payments after exceeding the deductible limit. A key institutional feature of almost all deductible contracts is that they reset at the beginning of the calendar year. In this section, we explain how our differences-in-regression-discontinuities approach exploits this to identify the reaction to dynamic incentives and to disentangle it from the reaction to static incentives.

It is useful to think of deductible contracts in terms of prices. If individuals have to pay for the last unit of care in a given period, then we say that the current price in that period is 1, and 0 otherwise. The current price in period t is specific to person i , as it depends on individual care consumption since the beginning of the year.

As noted already by Keeler et al. (1977), when deciding how much care to consume in t , individuals should take into account that out-of-pocket spending today can be seen as an investment that is associated with what they call a bonus: in expectation, out-of-pocket spending today will lower the price of care tomorrow. Based on a contribution by Ellis (1986), who characterizes optimal behavior for that model, a commonly used measure of dynamic incentives in this context is the expected end-of-year price.¹ For the standard deductible contracts we consider in this paper, at any point in time, this is equal to the probability that the patient will have to pay for the last unit of care in the year. This is also the measure we use. So, our aim is to measure how care consumption depends on the probability to pay out-of-pocket for the last unit of care in the year, controlling for medical needs and the current price.

Denote current prices by P_{it}^c and the expected end-of-year price by P_{it}^e . The superscript “c” stands for current and the superscript “e” stands for expected. Our reduced-form estimation equation makes the dependence of healthcare consumption on medical needs and both prices explicit by writing it as the sum of three parts: baseline consumption κ_{it} that is specific to

¹See for instance Keeler and Rolph (1988), Aron-Dine et al. (2015), Brot-Goldberg et al. (2017), and Abaluck et al. (2018). See Appendix A and in particular Appendix A.5.1 for a detailed discussion.

individual i in period t , and two price effects:

$$c_{it} = \kappa_{it} - \gamma^c \cdot P_{it}^c - \gamma^e \cdot P_{it}^e. \quad (1)$$

So, κ_{it} is the consumption of care when care is free. Our aim is to estimate the dependence of care consumption on dynamic incentives, γ^e .

2.2 Micro foundation

The reduced form equation (1) can be derived from a dynamic structural model of healthcare consumption. The advantage to providing such a structural foundation is that it gives the parameters γ^c and γ^e in the reduced form equation (1) a structural interpretation. Moreover, it justifies basing counterfactual experiments on estimates of the parameters of the reduced-form equation, as we can think of such a structural relationship as being stable under policy variation.

We now briefly describe the model setup. Appendix A contains all derivations for a more general version of the model and a discussion of various technical points. It explains in detail why care consumption can be written as a function of the two prices. The appendix also discusses how the model relates to other models in the literature and in which sense the expected end-of-year price is a good measure of dynamic incentives.

In the model, patient i knows how much care she has consumed up to t . She learns in period t about her medical needs λ_{it} and forms expectations on the likelihood to hit the deductible limit by the end of the calendar year. Her flow utility is quasi-linear in money and quadratic in the difference between care consumption and medical needs. Patients are quasi-hyperbolic discounters (O'Donoghue and Rabin, 1999) and dynamically optimize.

The utility function is specified such that if a patient has to pay for care consumption in the last period of the year (so that choice is static), then care consumption will be equal to the medical need λ_{it} . Conversely, if care is free, then patients consume $\lambda_{it} + \omega$. This means that ω is the additional care patients consume when it is free, a price effect. The parameters γ^c and γ^e in (1) are both functions of ω and β . β is a measure of present bias. It is between 0 and 1. If patients are fully aware of dynamic incentives, are forward-looking, and are not present-biased, then $\beta = 1$ and they will only react to P_{it}^e . The lower β the more patients react to P_{it}^c and the less they react to P_{it}^e . The most important two results that are derived in Appendix A are that (i) changes in dynamic incentives can be well measured by changes in the expected end-of-year price and that (ii) the effect of changes in dynamic incentives, γ^e , is equal to $\omega \cdot \beta$.² Next we show how we exploit the differences-in-regression discontinuities design to estimate this parameter of interest.

²We also show that $\gamma^c = \omega \cdot (1 - \beta)$.

2.3 Differences-in-regression-discontinuities design

We treat both prices, P_{it}^c and P_{it}^e , as known. A challenge for estimating the effect of prices on care consumption is that prices are endogenous: higher care consumption earlier in the year is associated with lower current and expected end-of-year prices, and at the same time likely positively correlated with medical needs later in the year. This means that prices will be negatively correlated with medical needs so that a regression of medical care on prices yields negatively biased coefficient estimates. The approach we pursue in this paper is to identify γ^e by exploiting a differences-in-regression-discontinuities design that allows us to use variation in P_{it}^e across years while at the same time controlling for static incentives and time effects.

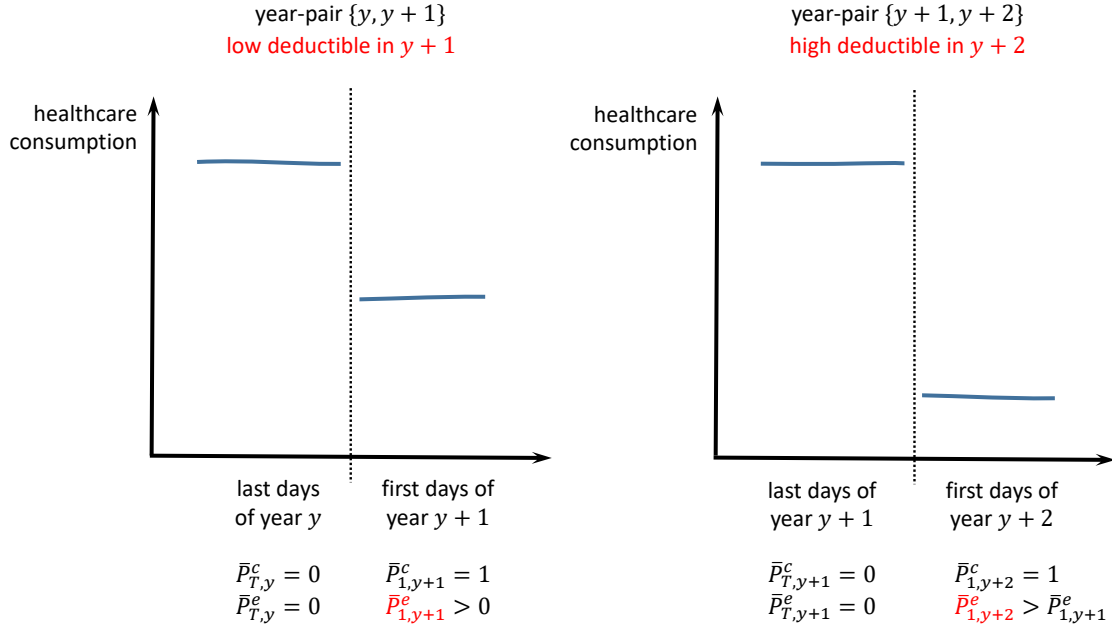
In our analysis we work with samples for year-pairs. The samples consist of observations for the last periods of the first year of the year-pair and the first periods of the second year. We restrict the samples to individuals who have reached the cost-sharing limit after January and before September of the first year.³ Thus, all individuals in our samples face a current price of 0 by the end of the first year of the year-pair, and therefore their expected end-of-year price is 0. Deductibles reset each year and therefore all individuals face a current price of 1 right at the beginning of the second year in the year-pair. This is the same for all years. Importantly, the expected end-of-year-price at the beginning of the second year in the year-pair is different across years due to changes in deductible amounts. From now on, we denote averages over individuals in our samples in period t of year y by a bar indexed by t and y . We can use averages because care consumption in (1) is a linear function of P_{it}^e . This implies that average care consumption is a function of the average expected end-of-year price.

Figure 1 illustrates our approach with an example. The left side shows average spending around the turn of the year for individuals who have crossed the deductible in year y . The solid blue line denotes average spending before the turn of the year, when $\bar{P}_{T,y}^c = 0$ and $\bar{P}_{T,y}^e = 0$. When individuals enter the new year, they face a current price of $\bar{P}_{1,y+1}^c = 1$ and a positive average expected end-of-year price $\bar{P}_{1,y+1}^e$, and they respond to the increase in prices by reducing their average spending. The right side shows care consumption around the turn of the year from $y + 1$ to $y + 2$ for a comparable sample of individuals who have exceeded the cost sharing limit in $y + 1$. For this group we have $\bar{P}_{T,y+1}^c = 0$ and $\bar{P}_{T,y+1}^e = 0$. In the figure, consumption just before the end of $y + 1$ (in the right year-pair) is the same as just before the end of y (in the left year-pair). At the beginning of the year, the current price is again $\bar{P}_{1,y+2}^c = 1$ in $y + 2$. However, the expected end-of-year price is higher than in $y + 1$, $\bar{P}_{1,y+2}^e > \bar{P}_{1,y+1}^e$, because the deductible is higher (and hence it is less likely that comparable individuals exceed the cost sharing limit). Therefore, when they enter the new year, they face a higher $\bar{P}_{1,y+2}^e$ and reduce their spending by a larger amount than in $y + 1$.

In the following, we formalize the intuition illustrated in Figure 1 and describe how we can estimate γ^e from differences in regression discontinuities. Formally, it follows from (1) that the

³We select the samples in a way so that the level of medical needs is comparable across year-pairs. Details are provided in Section 3.2 below.

Figure 1: Differences-in-regression-discontinuities design



Notes: The figure illustrates the intuition of our empirical approach for hypothetical year-pairs $\{y, y+1\}$ and $\{y+1, y+2\}$. The vertical lines depict the respective turn of the year. For both year-pairs the solid line depicts average care consumption at the end of the year (to the left of the vertical line) and average care consumption at the beginning of the new year (to the right of the vertical line). Individuals in the figure on the right reduce their average spending by a larger amount than individuals in the figure on the left because they face a higher expected end-of-year price at the beginning of the new year. The notation in the bottom is explained in the main text.

discontinuity in care consumption around the turn of the year from year y to $y+1$ is

$$\bar{c}_{1,y+1} - \bar{c}_{T,y} = \bar{\kappa}_{1,y+1} - \bar{\kappa}_{T,y} - \gamma^c \cdot \bar{P}_{1,y+1}^c - \gamma^e \cdot \bar{P}_{1,y+1}^e,$$

since our samples include only individuals for whom $\bar{P}_{T,y}^c = 0$ and $\bar{P}_{T,y}^e = 0$. For the discontinuity from year $y+1$ to $y+2$ we obtain a similar expression. Combining those, the difference in the discontinuities is given by

$$\begin{aligned} (\bar{c}_{1,y+2} - \bar{c}_{T,y+1}) - (\bar{c}_{1,y+1} - \bar{c}_{T,y}) &= (\bar{\kappa}_{1,y+2} - \bar{\kappa}_{T,y+1}) - (\bar{\kappa}_{1,y+1} - \bar{\kappa}_{T,y}) \\ &\quad - \gamma^c \cdot (\bar{P}_{1,y+2}^c - \bar{P}_{1,y+1}^c) - \gamma^e \cdot (\bar{P}_{1,y+2}^e - \bar{P}_{1,y+1}^e). \end{aligned}$$

Since $\bar{P}_{1,y+2}^c = \bar{P}_{1,y+1}^c = 1$ at the beginning of the year,⁴ we obtain

$$(\bar{c}_{1,y+2} - \bar{c}_{T,y+1}) - (\bar{c}_{1,y+1} - \bar{c}_{T,y}) = (\bar{\kappa}_{1,y+2} - \bar{\kappa}_{T,y+1}) - (\bar{\kappa}_{1,y+1} - \bar{\kappa}_{T,y}) - \gamma^e \cdot (\bar{P}_{1,y+2}^e - \bar{P}_{1,y+1}^e).$$

⁴It is possible that some individuals already exceed the deductible limit in the first period of the year. We discuss this in Section 6.1.

This equation differences out different levels in baseline care consumption across year-pairs. Such differences could arise, for example, because the flu was particularly severe in some winters. If the flu season was particularly severe around the turn of the year from y to $y + 1$, then this would affect both $\bar{\kappa}_{T,y}$ and $\bar{\kappa}_{1,y+1}$. By taking differences between regression discontinuities we control for such differences in seasonal effects across year-pairs.

Our main identifying assumption is that

$$(\bar{\kappa}_{1,y+2} - \bar{\kappa}_{T,y+1}) = (\bar{\kappa}_{1,y+1} - \bar{\kappa}_{T,y}) \quad (2)$$

This assumes that the difference in mean baseline consumption across the turn of the year is the same for both year-pairs. In Section 6, we discuss the plausibility of this assumption, and we conduct a number of robustness checks. Under this assumption, our parameter of interest is given by the ratio between the difference in regression discontinuities and the difference in expected end-of-year prices from the perspective of the beginning of the year,

$$\gamma^e = \frac{(\bar{c}_{1,y+2} - \bar{c}_{T,y+1}) - (\bar{c}_{1,y+1} - \bar{c}_{T,y})}{(\bar{P}_{1,y+2}^e - \bar{P}_{1,y+1}^e)}.$$

The derivation so far was for two year-pairs. For multiple year-pairs, as we have in our study, we can obtain the discontinuity for each year-pair and then plot it against the average expected end-of-year price in the respective second year of each year-pair. We can use this as the basis for testing whether the relationship is monotonic and to estimate γ^e using a linear regression of the discontinuity on the expected end-of-year price. See Section 4.2.

3 Institutional background and data

3.1 Institutional background

In the Netherlands, health insurance is mandatory. Patients have to buy insurance that covers a “basic package”. Since 2008, by law it has to feature a deductible for all residents who are at least 18 years old. For each calendar year, the minimum, mandatory deductible amount is set by the Dutch Government to a baseline amount. Individuals are allowed to opt for a higher deductible.⁵ There have been substantial increases in the amount of the mandatory deductible over time: for example, the deductible was €150 in 2008; in 2015, it was €375.⁶ The deductible resets annually, regardless of how much healthcare was consumed in the previous year.

⁵This voluntary deductible can be up to €500 above the mandatory deductible. Very few people in our data choose a voluntary deductible and we omit them from our sample.

⁶The mandatory deductible was €155 in 2009, €165 in 2010, €170 in 2011, €220 in 2012, €350 in 2013 and €360 in 2014. The increases in deductibles in the years since 2010 reflect the political preferences of governments under prime minister Rutte, head of the centre-right VVD party. The size of deductibles is a contentious political topic in the Netherlands, with more right-wing parties in favor of higher deductibles and left-wing parties in favor of lower (or no) deductibles.

Some services are exempt from the deductible, such as consultations by General Practitioners (GPs), maternity care, and medical equipment on rent (e.g., wheelchairs).

The contents of the basic package are determined by law, and they are adjusted annually.⁷ Treatments are largely billed in terms of diagnosis treatment combinations (DTCs). A DTC compensates for all care administered within an episode of treatment, including follow-up visits. Compensation for DTCs are determined through bargaining between insurers and providers. Providers send a bill to health insurers, who then determine how much patients have to pay out-of-pocket depending on their remaining deductible.

3.2 Data

We use claims data from a large Dutch health insurer for the years 2008 to 2015 (see [Hayen et al., 2015](#) for details).⁸ We restrict our analysis to types of care that are part of the basic package and count towards the deductible.⁹ Our data include the amount paid for a claim, the date of claim initiation, the type of claim, and demographic information, such as age and gender of the enrollees.

We construct separate samples for each year-pair. Each sample consists of individuals who cross the deductible at any point between February and August (inclusive) in year y and we follow these individuals around the turn of the year, from y to $y + 1$.¹⁰ Conditioning on individuals who cross the deductible in a given year leads to sicker samples in years with a larger deductible, as individuals have to spend more to be included in the sample of crossers for these years: for example, in 2008, an individual would have to spend only €150 to cross the deductible, while in 2014 she would have to spend €360. To address such concerns, we utilize a percentile-matching strategy based on [Brot-Goldberg et al. \(2017\)](#). Specifically, across all years we include only the top 38% of cumulative spenders by the end of August in our sample. 38% is the share of individuals who exceed the deductible limit between February and the end of August in the year 2014, the year with the lowest share of such individuals in our data. In 2008, for example, even though 47% of the sample crossed the deductible between February and the

⁷A list of the changes can be found (in Dutch) on <https://www.cbs.nl/nl-nl/onze-diensten/methoden/onderzoeksomschrijvingen/aanvullende%20onderzoeksbeschrijvingen/pakketwijzigingen-zorgverzekeringswet>.

⁸Our data were obtained under a pilot project in which a new payment model for GPs was evaluated. The pilot project started in July 2014, and the data cover multiple years before the start of the pilot project. While the pilot project was not related to patient cost-sharing we cannot exclude the possibility that the results for the last year-pair {2014,2015} could be influenced by the treatment in the pilot project. In order to take this into account, we conduct a sensitivity analysis in which we restrict our analysis to year-pairs before the start of the pilot project. See Appendix D.5.

⁹This implies we do not look at consultations at a GP or maternity care, for example.

¹⁰We exclude individuals that cross the deductible in January of year y since they have a very different pattern of healthcare expenditures when compared to the rest of the sample. We exclude individuals who cross the deductible after August of year y because healthcare expenditures can be autocorrelated over time, and we aim to limit the influence of healthcare expenditures associated with crossing the deductible in the previous year on expenditures across the turn of the year.

end of August, we only include the top 38% of total spenders in our sample.¹¹ The total number of individuals in our final sample, after applying percentile-matching, and the total number of individuals without matching in each year are reported in the last two rows of Table 1.

Our goal is to compare healthcare utilization at the end of year y with healthcare utilization at the beginning of year $y + 1$. For this, we specify the date of treatment as the date the claim was initiated and aggregate our claims data to the daily level.

Non-emergency care is limited in availability on weekends and during the Christmas break. Therefore, we only use week days for our analysis. Moreover, for each year-pair we omit a set of days that include the Christmas break. For this, we specify a last regular day before the start of the Christmas break (day $t = T$ in the first year of the year-pair) and a first regular day (day $t = 1$ in the second year of the year-pair) after the end of the Christmas break. The resulting empirical setup is commonly referred to as a donut regression discontinuity (RD) design (Barreca et al., 2011).¹² Gerfin et al. (2015) have used it earlier to estimate the size of the discontinuity in healthcare expenditures around the turn of the year for one year-pair.¹³

The distribution of healthcare expenditures is characterized by a heavy right tail. Although our sample size is not small, outliers can still have a large influence on estimated coefficients. To ensure that our results are not influenced by high healthcare expenditures far above the deductible amount, we pseudo-censor our expenditure variable, i.e., we code any daily expenditure above €500 as €500. This cutoff is higher than the highest deductible in our study period. In Appendix D we show that our results are robust to changes in this cutoff amount and to an alternative transformation of the expenditure variable from levels to logs, adding a constant of 1 to account for the zeros present in the data.

Table 1 shows summary statistics, separately for all 7 year-pairs, $\{y, y + 1\}$, in our data. The first three rows report the average age, proportion female, and average income at the 6-digit postal code level in each year y .¹⁴ The next two rows show the deductible in year $y + 1$ and the in-sample expected end-of-year price $\bar{P}_{1,y+1}^e$ at the beginning ($t = 1$) of year $y + 1$.¹⁵ Generally speaking, there is a positive relationship between the two, but this is not always the case: for

¹¹The underlying idea is that there is a monotone relationship between healthcare needs and spending. Then, by selecting the top 38% spenders in all years, one makes sure that the sample populations are comparable across years in terms of needs. The exact assumptions underlying this approach are discussed in Brot-Goldberg et al. (2017).

¹²The dates we use in our preferred specification are detailed in Appendix B. We use a donut hole of 20 days across all year-pairs, except for the year-pair $\{2011, 2012\}$, where the donut hole is 27 days long. Section 6.4 delves into the robustness of our primary results with respect to changes to the days—and the distance between the days—for which utilization changes are compared.

¹³They also document in detail, in their Figure 2, that care consumption is lower on weekends and during the Christmas break.

¹⁴In comparison, average age in the general adult population is 48.75 years. The proportion female is 0.52. Average income is €2188. See the second column of Table 1 in Hayen et al. (2019).

¹⁵We compute the expected end-of-year price as $1 - \Pr(\text{cross}_{y+1})$, where $\Pr(\text{cross}_{y+1})$ is the proportion of individuals in the sample who cross the deductible by the end of year $y + 1$.

Table 1: Summary statistics

	year-pair $\{y, y+1\}$									
	2008, 2009	2009, 2010	2010, 2011	2011, 2012	2012, 2013	2013, 2014	2014, 2015			
demographics in y										
average age	54.57	54.49	54.68	55.35	55.95	57.05	57.93			
percentage female	.575	.575	.572	.565	.568	.566	.567			
average income ^a	2123	2121	2141	2121	2123	2114	2114			
dynamic incentives at the beginning of $y+1$										
deductible in $y+1$	155	165	170	220	350	360	375			
expected end-of-year price $(\bar{P}_{1,y+1}^e)^b$.1452	.1443	.1534	.1691	.2170	.2124	.2200			
average daily expenditure ^c										
September to December year y	11.74	13.14	12.72	13.94	13.03	14.60	14.09			
January of year $y+1$	12.92	10.56	10.50	10.93	11.81	12.32	11.32			
average daily expenditure (PC) ^{c,d}										
September to December year y	4.50	4.83	4.98	5.23	5.48	5.96	6.01			
January of year $y+1$	4.41	4.76	4.55	4.78	5.11	5.38	5.31			
prob. of any daily expenditure ^c										
September to December year y	.057	.062	.064	.069	.070	.076	.078			
January of year $y+1$.055	.063	.061	.064	.066	.074	.074			
no. of individuals without matching ^e	29289	32953	33985	36329	33845	31282	29626			
no. of individuals in estimation sample $(N)^f$	22798	25090	27473	26706	28045	30864	29626			

Notes: This table shows summary statistics for our baseline estimation sample.

^a Average income is at the 6-digit postal code level.

^b $\bar{P}_{1,y+1}^e$ is computed as 1 minus the proportion of individuals in the sample who cross the deductible by the end of year $y+1$.

^c We compute average utilization for regular days only, i.e., we do not include utilization on weekends and days around holidays in computing this average.

^d PC refers to “pseudo-censoring”. The threshold used to pseudo-censor was 500.

^e This is the number of individuals who exceed the deductible between February and August in year y and are observed in year $y+1$.

^f This is the number of individuals in our estimation sample. They exceed the deductible between February and August in year y , are observed in year $y+1$, and are the top 38% of cumulative spenders by the end of August (percentile matching). See Section 3.2 for details.

example, between 2008 and 2009 the deductible increased by €10, while $\bar{P}_{1,y+1}^e$ decreased.

Table 1 also shows different measures of average healthcare utilization, on regular days, for the last four months in year y and the first month of year $y + 1$. We look at mean daily expenditures, mean pseudo-censored daily expenditures, and the probability of having any daily expenditure. We see that utilization, across all measures, increases across years both at the end of year y and at the start of year $y + 1$. We also see that spending at the start of year $y + 1$ is generally lower than at the end of year y , suggesting that patients lower healthcare spending in response to financial incentives, and we see that the difference between spending at the end of year y and spending at the beginning of year $y + 1$ is generally larger for year-pairs with a larger $\bar{P}_{1,y+1}^e$, pointing to forward-looking behavior.

4 Empirical implementation

4.1 Estimation of discontinuity sizes

For each year-pair, we estimate the change in care consumption around the turn of the year using separate local linear regressions before and after the turn of the year. We use a triangular kernel and optimal bandwidths, as detailed in [Calonico et al. \(2014\)](#).¹⁶ The underlying estimation equations are

$$\begin{aligned} c_{it} &= \alpha_0 + \gamma_0(t - (T + 1)) + \varepsilon_{it} \text{ for year } y \text{ and } t \leq T \\ c_{it} &= \alpha_1 + \gamma_1(t - 1) + \varepsilon_{it} \text{ for year } y + 1 \text{ and } t \geq 1, \end{aligned}$$

where day-of-the-year t is the running variable (note that day $t = T + 1$ in year y is the same as day $t = 1$ in year $y + 1$ because the days in-between are omitted) and c_{it} is our measure of healthcare utilization. Our estimate of interest is $\hat{\alpha}_1 - \hat{\alpha}_0$.

Figure 2 illustrates our approach for the year-pair $\{2010, 2011\}$. The last day that we use for our analysis, $t = T$ in 2010, is Thursday 16 December. The first day, $t = 1$ in 2011, is Thursday 6 January. The days in-between lie in the donut hole and are omitted from our analysis because non-emergency care is limited in availability on weekends and during the Christmas break. The situation corresponds to the one described in Figure 1: the current price is 0 at the end of 2010, and it is 1 at the beginning of 2011.¹⁷ In Section 6.4 we show that our results are robust to different specifications of the donut hole.

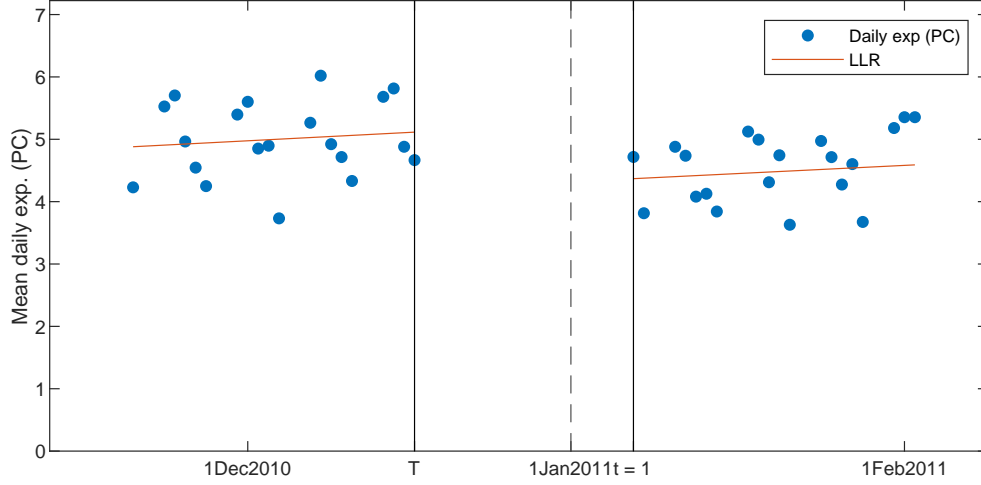
The figure shows that the time trends are similar in both years and, in line with individuals reacting to cost-sharing incentives, healthcare utilization is lower at the beginning of 2011.

There are two interesting aspects of our estimation procedure that we would like to highlight. First, since we follow a balanced panel of individuals, covariates are, by construction, balanced

¹⁶We also performed our empirical analysis using bandwidths of 20 days and found almost no difference.

¹⁷Some individuals exceed the deductible limit on the first days of the year so that current prices are not 1 for all individuals. We show in Section 6.1 that this does not pose a threat for our analysis.

Figure 2: Healthcare consumption around turn of year



Notes: The figure plots mean daily pseudo-censored healthcare spending (denoted by the blue dots) for weekdays. Weekends were omitted from our analysis. T denotes the last day of year y that we use for our analysis and $t = 1$ denotes the first day of year $y + 1$ that we use. See Table A.1 in the appendix for details. The red solid line denotes the local linear regression estimates (LLR). Details of the regression discontinuity estimation are included in Section 4.1.

across both sides of the threshold we compare. Second, since we subtract predicted spending at the end of year y from predicted spending at the beginning of year $y + 1$, we remove any effect of characteristics and influences that are invariant across these two dates, such as, for example, the severity of the flu season.

4.2 Relating discontinuity sizes to dynamic incentives

We repeat the above for all year-pairs, $\{y, y + 1\}$. Our main focus is on estimating the relationship between the changes in expenditure around the turn of the year and $\bar{P}_{1,y+1}^e$. If patients are forward-looking, then the change should be bigger (in absolute terms) in year-pairs with higher $\bar{P}_{1,y+1}^e$.

We use a weighted least-squares regression of the estimated discontinuity size around the turn of the year on $\bar{P}_{1,y+1}^e$ to estimate our parameter of interest, γ^e . Specifically, our estimate of γ^e is given by the slope coefficient of the weighted least-squares regression. The weights take estimation error into account. The underlying assumptions for estimation and inference follow from Hanushek (1973), as outlined in Appendix C.1.

This regression imposes a linear relationship between the estimated discontinuity size and $\bar{P}_{1,y+1}^e$. To test for monotonicity without imposing parametric restrictions, we borrow a non-parametric test for monotonicity from the finance literature, which is described in Patton and Timmermann (2010). The test is based on the ranking of each coefficient w.r.t. $\bar{P}_{1,y+1}^e$: we test

whether the ordering in $\bar{P}_{1,y+1}^e$ is the same as the ordering of the size of the estimated changes. The details of these tests are outlined in Appendix C.2.

To summarize, in our analysis, we proceed in two steps. In a first step, we form year-pairs for comparable samples of individuals and estimate the change in healthcare consumption around the turn of the year. An individual is in the sample for a given year-pair if she exceeds the deductible limit between February and August in the first year of the year-pair and is among the top 38 percent of the spenders (percentile matching). For this sample of individuals, we also compute our measure of dynamic incentives patients face at the beginning of the second year, $\bar{P}_{1,y+1}^e$, as the proportion of individuals in the sample that actually exceed the deductible limit in year $y + 1$. In a second step, we relate the estimated changes in healthcare consumption around the turn of the year to $\bar{P}_{1,y+1}^e$.

Our main assumption is that in the absence of changes in dynamic incentives changes in mean healthcare utilization around the turn of the year would be the same across all year-pairs. This assumption is formally stated in (2) in Section 2.3. Threats to our empirical strategy stem from any possible reason other than the response to dynamic incentives that may explain a negative relationship between the estimated changes in healthcare utilization around the turn of the year and $\bar{P}_{1,y+1}^e$. We discuss such possible violations of our main identifying assumption in Section 6.

5 Results

5.1 Baseline results

We present the results for the empirical framework laid out in Section 4 for two measures of healthcare utilization: mean (pseudo-censored) expenditures and the probability of any claim (extensive margin). We first present estimates for the changes in utilization around the turn of the year, and then we relate these estimates to the respective expected end-of-year prices.

Table 2 reports the estimated changes in healthcare utilization at the turn of the year, for both measures. For all year-pairs, healthcare utilization decreases at the turn of the year, and these decreases are statistically significant at any conventional level. Decreases in healthcare utilization tend to be stronger for later years. For example, for the year-pair {2008,2009} the decrease in pseudo-censored daily expenditures is €0.47, and for the year-pair {2014,2015} the decrease is €1.30. Correspondingly, the probability of any daily expenditure decreases by 0.5 percentage points for the the year-pair {2008,2009} and by 1.3 percentage points for the year-pair {2014,2015}.

Figure 3 plots these estimated changes in healthcare consumption at the turn of the year against $\bar{P}_{1,y+1}^e$. As we have seen in Table 1, $\bar{P}_{1,y+1}^e$ tends to be higher for later years in our sample when the deductible is higher. For example, $\bar{P}_{1,2009}^e$ is 0.15, while $\bar{P}_{1,2015}^e$ is 0.22. Figure 3 reveals a decreasing relationship, providing evidence for forward-looking behavior. Table 3

Table 2: Discontinuity sizes

year-pair $\{y, y + 1\}$	change in daily expenditure (PC) around turn of the year	change in probability of any daily expenditure around turn of the year
2008,2009	-0.737 (0.1368)	-0.010 (0.0010)
2009,2010	-0.473 (0.1252)	-0.005 (0.0011)
2010,2011	-0.827 (0.1270)	-0.008 (0.0009)
2011,2012	-0.955 (0.1242)	-0.011 (0.0010)
2012,2013	-1.345 (0.1682)	-0.016 (0.0014)
2013,2014	-1.296 (0.1345)	-0.013 (0.0013)
2014,2015	-1.577 (0.1712)	-0.013 (0.0011)

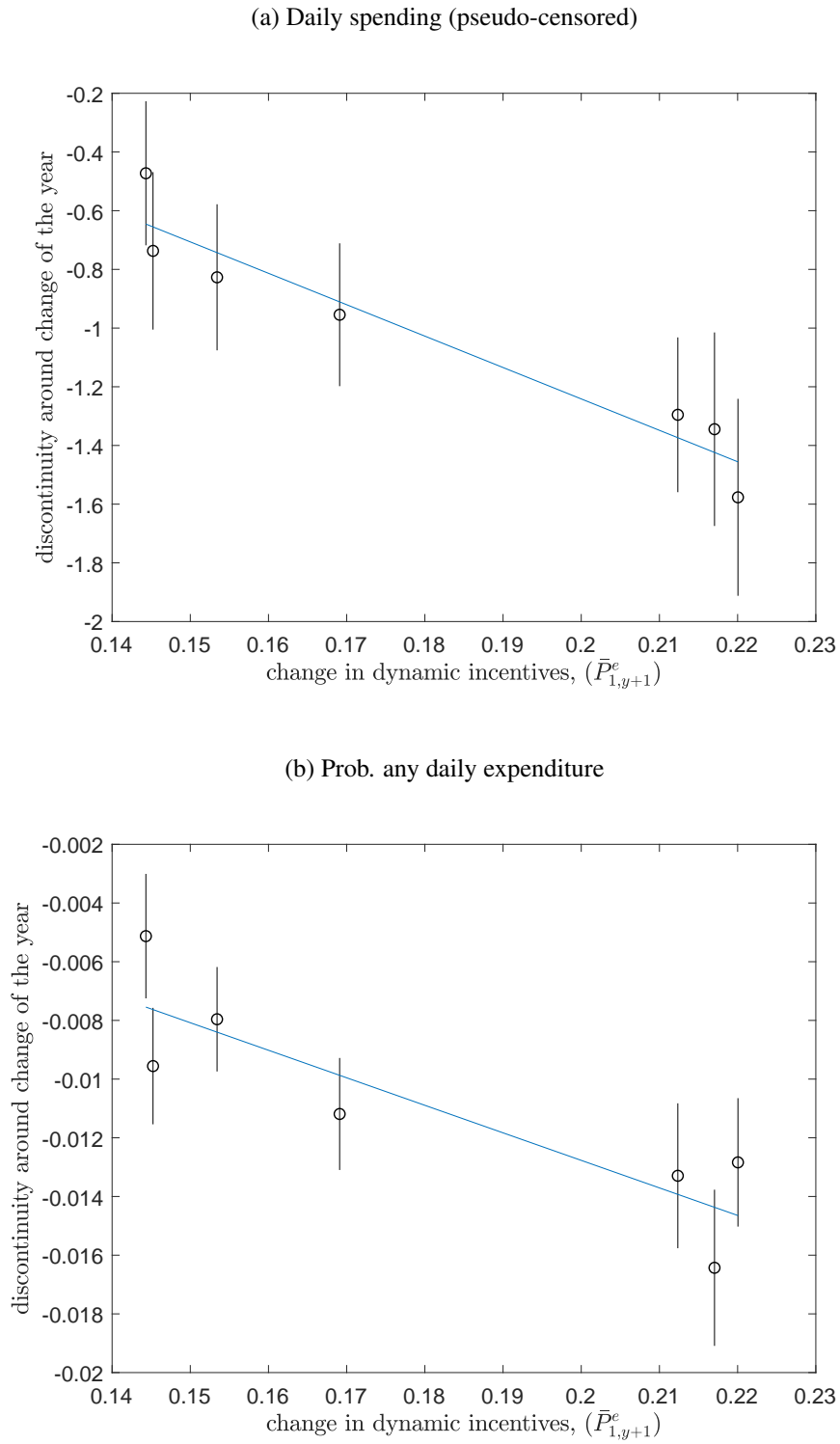
Notes: Expenditures larger than 500 were coded as 500 (hence the abbreviation PC). Changes are estimated using a donut hole regression discontinuity design (see Section 4.1). Robust standard errors that are clustered at the individual and year-pair level are shown in parentheses.

Table 3: Dependence of discontinuity sizes on dynamic incentives

	daily expenditure	any daily expenditure
effect of dynamic incentives (γ^e)	-10.690 (1.4177)	-0.094 (0.0230)
p -value nonparametric MR test	0.000	0.926
p -value nonparametric Up test	0.006	0.000
p -value nonparametric Down test	1.000	0.178

Notes: The effect of dynamic incentives was estimated by regressing the estimated discontinuities reported in Table 2 on $\hat{P}_{1,y+1}^e$, using a weighted linear regression. The weights take estimation errors into account, and they are computed following Hanushek (1973). See Appendix C.1 for details. The MR, Up and Down tests were conducted using 10,000 bootstrap repetitions. See Appendix C for details.

Figure 3: Dependence of discontinuity sizes on dynamic incentives



Notes: These figures plot the relationship between the estimated changes in healthcare consumption around the turn of the year and the computed $\bar{P}_{1,y+1}^e$. See Tables 1 and 2 for exact numbers. The vertical lines denote 95% confidence intervals. The solid line is the OLS regression line.

reports our estimate of the dependence of the discontinuity size on dynamic incentives. Based on a weighted linear regression we find that the slope is equal to -10.7 for daily expenditures, and it is statistically significant at any conventional level. This coefficient is our estimate of γ^e . Our results thus suggest that increasing $\bar{P}_{1,y+1}^e$ by 10 percentage points leads to a reduction in mean (pseudo-censored) expenditures of around €1.07 and a reduction in the probability of any claim of around 0.9 percentage points.

The bottom three rows of Table 3 show non-parametric tests for monotonicity. These tests use the ordering of the estimated changes implied by the ordering of $\bar{P}_{1,y+1}^e$ and take the estimation error resulting from estimating the discontinuities into account. We specified the test statistic such that results aligning with forward-looking behavior produce p -values smaller than a chosen level of significance for the MR test and the Up test only (see Appendix C.2 for details). The Down test should not be significant. For pseudo-censored expenditures, we see that the MR test and the Up test are significant at conventional levels and the Down test is not, indicating a monotonically decreasing relationship between mean expenditures and $\bar{P}_{1,y+1}^e$. For the extensive margin, the MR test is not significant at conventional levels. However, the p -values from the Up and Down test suggest that this result could stem from limited power inherent in the MR test (as discussed in Patton and Timmermann, 2010).

The variation in $\bar{P}_{1,y+1}^e$ comes from changes in deductibles. Therefore, we can relate these changes to changes in the deductible. In the first row of Table 4 we present the effects of a €100 increase in the deductible. This increase leads to a 3.4 percentage point increase in $\bar{P}_{1,y+1}^e$,¹⁸ which implies a reduction of around €0.36 in mean (pseudo-censored) expenditures and a 0.32 percentage point reduction in the probability of making any claim. In percentage terms, these reductions amount to around 3.0% of mean expenditures and 4.4% of the probability of any claim on the first day of the year.¹⁹

5.2 Heterogeneity

Next, we examine whether and how responses to $\bar{P}_{1,y+1}^e$ differ for different groups in the population, defined by gender, neighborhood income, and age. Previous studies such as Manning et al. (1987), Brot-Goldberg et al. (2017) and Farbmacher et al. (2017) have examined heterogeneous responses to cost-sharing incentives for different groups in the population. However, there is still limited evidence on heterogeneous response specifically to dynamic incentives.

Responses to dynamic incentives can differ between groups as a result of two mechanisms that may be at play at the same time. The first mechanism is that the size of the reaction is related to healthcare needs, and that different groups in the population have different healthcare needs. This is plausible, as for instance the possibilities to exert moral hazard may be wider for older patients who are less healthy. The second mechanism is that different groups in the

¹⁸This value was obtained by running a regression of $\bar{P}_{1,y+1}^e$ on the deductible amount across all years for our sample.

¹⁹We use the estimated average spending on the first day of 2015 to calculate these percentage reductions.

Table 4: Difference in reaction to dynamic incentives across groups (controlling for risk score)

	effect of dynamic incentives (γ^e)			range of $\bar{P}_{1,y+1}^e$		utilization at beginning of year ^d		average N	effect of €100 increase in deductible ^b		
	daily expenditure	any daily expenditure	min.	max.	daily expenditure	any expenditure	change in $\bar{P}_{1,y+1}^e$ (in pp)		daily expenditure	any expenditure	% reduction in utilization ^c
baseline sample	-10.690 (1.418)	-0.0938 (0.023)	0.1443	0.2200	12.211	0.0733	27229	3.4	-3.0 %	-4.4 %	
female	-9.232 (2.023)	-0.0781 (0.027)	0.1230	0.2057	11.687	0.0778	15503	3.7	-3.0 %	-3.8 %	
male	-15.147 (4.435)	-0.0959 (0.026)	0.1718	0.2400	13.036	0.0680	11726	3.0	-3.5 %	-4.3 %	
below median income ^d	-12.341 (2.823)	-0.1171 (0.034)	0.1396	0.2121	13.593	0.0839	11456	2.7	-2.5 %	-3.8 %	
above median income	-9.827 (2.983)	-0.0786 (0.022)	0.1485	0.2375	11.401	0.0641	14786	3.7	-3.2 %	-4.5 %	
age below 45	-3.020 (3.822)	-0.0275 (0.019)	0.2408	0.3442	9.320	0.0423	8272	4.2	-1.4 %	-2.8 %	
age 45 and above	-12.103 (1.518)	-0.0911 (0.027)	0.0976	0.1829	13.040	0.0818	18957	3.8	-3.5 %	-4.2 %	

Notes: The effect of dynamic incentives was estimated by regressing the estimated discontinuities on $\bar{P}_{1,y+1}^e$. See Appendix C for details. For this, the discontinuities and $\bar{P}_{1,y+1}^e$ were estimated within each subsample. When doing so, we re-weight observations so that the distribution of the risk score in each re-weighted subsample is equal to the distribution in the baseline sample. To obtain the effect of a 10 percentage point change in $\bar{P}_{1,y+1}^e$ on expenditures, the reported estimates have to be divided by 10. The third and fourth column report the domain of $\bar{P}_{1,y+1}^e$ for each of the subgroups. The last two columns quantify the effect of increasing the deductible by €100 on healthcare utilization for the first day of the new year. We only look at utilization on the first day as all individuals are, by definition, below the deductible.

^a Daily utilization at the beginning of the year was estimated using a local linear regression on the first regular day of 2015 (see 3.2). The observations are weighted such that the distribution of the risk score in each re-weighted sample is equal to the distribution of the risk score in the baseline sample.

^b The effect of an increase in the deductible by €100 is only for the first day of 2015. For the effect of such an increase in the deductible on expenditures throughout the year, refer to Section 7.

^c The estimated average utilization (not pseudo-censored) on the first day of 2015 was used as the base for these percentage change computations.

^d Income is average income at the 6-digit postal code level. Income data were missing for around 1000 individuals across all year-pairs.

population respond differently to dynamic incentives, even if the needs are the same. This could be explained by different information and underlying utility parameters. For example, patients might respond less strongly (or not at all) to dynamic incentives if they are not aware of them, are myopic, or strongly discount the value of future states (a low β in the terms of our model). They might also respond less strongly to dynamic incentives if they exert less moral hazard (a low ω in the terms of our model).

In the following, we focus on the second mechanism. We examine how different groups in the population respond to dynamic incentives while controlling for differences in expected healthcare needs as measured by risk scores.²⁰ For this, we split the sample, alternatively according to gender, average income in the neighborhood below and above median, and ages below 45 years and 45 years and older. For each subsample, we weigh observations according to risk score quintiles such that the distribution of risk scores is the same as for the baseline sample.²¹

Column 1 of Table 4 shows our estimates of γ^e for different subgroups in the population. We find that almost every subgroup exhibits forward-looking behavior. The estimates of γ^e are significantly different from zero for all subgroups except for individuals that are below age 45. Point estimates are similar for males and females, and for individuals living in high and low income neighborhoods.

The results for any daily expenditure, shown in Column 2, show a similar pattern. Table 4 also shows the range of $\bar{P}_{1,y+1}^e$ across years, utilization at the beginning of the year for each subgroup. These numbers are weighted by risk score quintiles.

In addition, we quantify the reduction in utilization (in percentage terms) on the first day of a new year if the deductible were to be increased by €100; changes in $\bar{P}_{1,y+1}^e$ are shown in the third-to-last column of Table 4, and changes in healthcare consumption (in percentage terms) are presented in the last two columns of Table 4. An increase in the size of the deductible by €100 leads to a -1.5% change in daily expenditures for individuals below age 45. For the remaining groups, we find that an increase in the size of the deductible by €100 leads to a percentage change in daily expenditures that ranges from -2.5% to -3.5%.

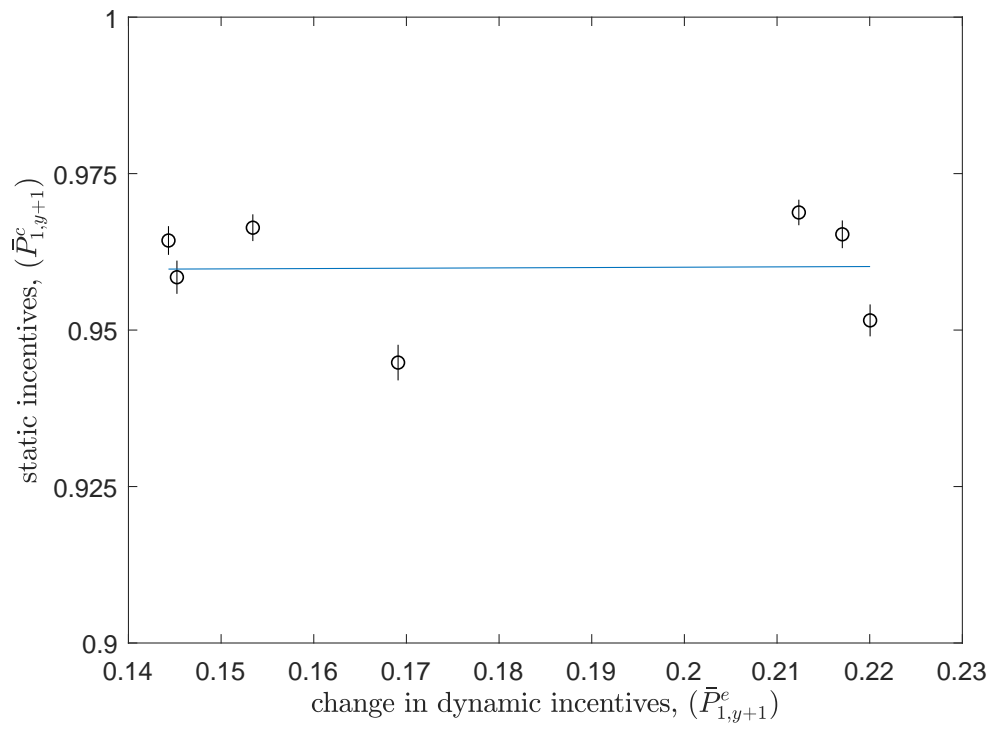
6 Discussion of main identifying assumption and robustness

Our main identifying assumption is that if there are no changes in incentives then changes in healthcare utilization around the turn of the year should (in expectation) be constant across year-pairs. This assumption is formally stated in (2) in Section 2.3, and it allows us to attribute the negative relationship between the changes in healthcare consumption around the turn of the year

²⁰The risk score of an individual is given by her predicted annual expenditures divided by average annual expenditures. The larger the risk score, the more a person is predicted to spend, relative to the average. More details on how this risk score was computed can be found in [Hayen et al. \(2019\)](#).

²¹For each weighted subsample we apply the same estimation approach as for the baseline sample.

Figure 4: Current price changes



Notes: This figure plots the relationship between the current price and $\bar{P}_{1,y+1}^e$. The solid line is the OLS regression line.

Table 5: Robustness

specification	slope		MR test	Up test	Down test
current price	0.005	(0.1131)	1.000	0.000	0.000
inflation corrected	-10.171	(3.0301)	0.136	0.001	0.970
Thursday, 27 day donut	-12.628	(2.7761)	0.484	0.000	0.650
Wednesday, 20 day donut	-8.513	(1.3192)	0.001	0.151	1.000
Wednesday, 27 day donut	-10.742	(2.2300)	0.000	0.049	1.000

Notes: The table reports our estimates for the dependence of the current price (first row) and changes in health care consumption around the change of the year (remaining rows) on $\bar{P}_{1,y+1}^e$. For the last 4 rows the reported slope coefficients are estimates of γ^e . The table also reports our results from non-parametric tests of monotonicity. See Section 6.1 to 6.4 for details. The slope coefficients were estimated by regressing the estimated discontinuities on $\bar{P}_{1,y+1}^e$. The MR, Up and Down tests were conducted using 10,000 bootstrap repetitions. See Appendix C for details.

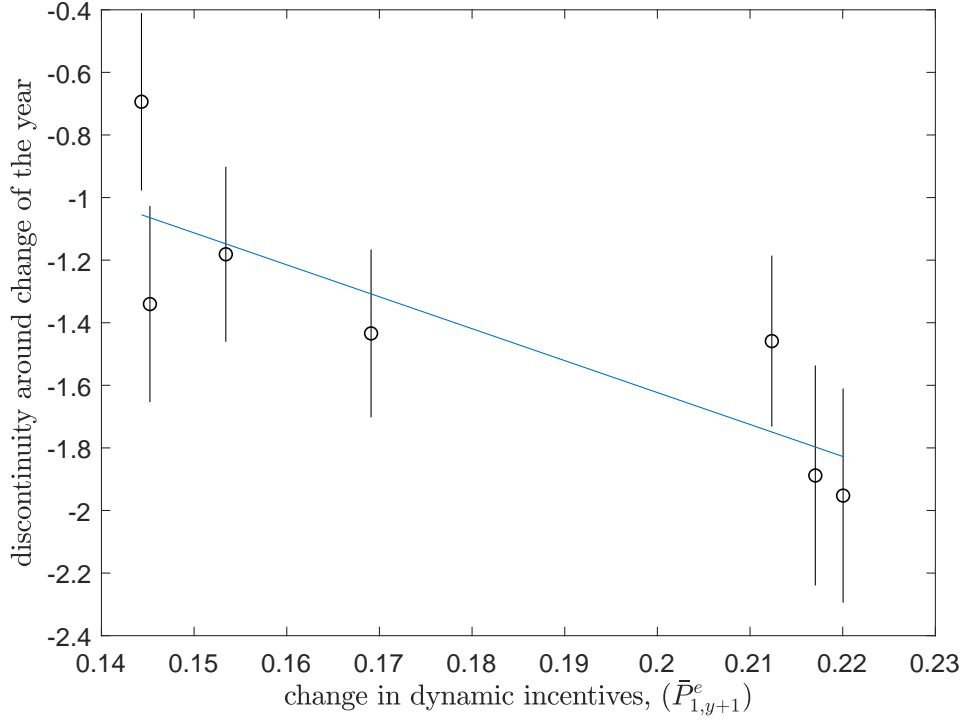
and $\bar{P}_{1,y+1}^e$ to dynamic incentives. Thus, threats to our empirical strategy stem from any other possible reason that could explain such a negative relationship. In the following we discuss three potential threats: 1) Current price changes, 2) changes of provider prices and insurance coverage, and 3) strategic timing of medical care use. We find that these alternative explanations cannot generate our results. We also show that our results are robust to different specifications of the donut hole and the pseudo-censored outcome variable.

6.1 Current price changes

One concern could be that the negative relationship between the changes in healthcare consumption around the turn of the year and $\bar{P}_{1,y+1}^e$ shown in Figure 3 can be explained by changes in the average current price instead of changes in $\bar{P}_{1,y+1}^e$. Even though we look at days early in January, it is possible that some individuals exceed the deductible already within the first days of the year. Given that it requires more spending to cross the deductible (and thus have a spot price of 0) in years with a higher deductible, it could be that $\bar{P}_{1,y+1}^c$ is higher in years with a higher deductible. This would result in a positive correlation between $\bar{P}_{1,y+1}^c$ and $\bar{P}_{1,y+1}^e$. In this case, we would falsely ascribe the effects of static incentives to the effects of dynamic incentives.

We can test whether $\bar{P}_{1,y+1}^e$ and $\bar{P}_{1,y+1}^c$ are indeed positively correlated. For this, we compute $\bar{P}_{1,y+1}^c$, the share of individuals who have not exceeded the deductible limit by $t = 1$ in year $y + 1$, for all year-pairs in our data and plot them against $\bar{P}_{1,y+1}^e$. Figure 4 shows that there is little variation in $\bar{P}_{1,y+1}^c$ across years. The variation is not systematically related to $\bar{P}_{1,y+1}^e$ and the slope coefficient of the corresponding regression line, shown in the first row of Table 5, is close to 0 and not statistically significant. Thus, we conclude that our results cannot be attributed to changes in current prices instead of expected end-of-year prices.

Figure 5: Accounting for changes of provider prices and insurance coverage



Notes: This figure plots the relationship between the estimated changes in healthcare consumption around the turn of the year and the computed $\bar{P}_{1,y+1}^e$ after correcting for changes in provider prices and the basic insurance package across years. We first deflated expenditures and then pseudo-censored them at €500. The solid line is the OLS regression line.

6.2 Changes of provider prices and insurance coverage

Our empirical strategy attributes differences in discontinuity sizes across year-pairs to differences in $\bar{P}_{1,y+1}^e$. However, there are other factors that also change at the beginning of the year. For example, costs of claims and the contents of the basic health insurance package are adjusted at the beginning of the year. This could lead to differences in discontinuity sizes across year-pairs that are not related to differences in $\bar{P}_{1,y+1}^e$.

In order to control for changes in expenditures due to changes in the cost of claims and the basic package, we create an annual expenditure deflator using data for periods in each year where the price of healthcare is 0. Recall that individuals were selected into our sample conditional on having crossed the deductible by the end of August in the first year of the year-pair. This implies that these individuals face a 0 price of healthcare from September to December in that first year. We take the ratio of average pseudo-censored expenditures across these months, for a given year, to the average pseudo-censored expenditures for a base year. This produces an expenditure deflator that accounts for differences in expenditures due to changes in the cost of claims and the basic package, i.e., changes that are not influenced by changes in cost-sharing incentives. The values for our expenditure deflator for each year are reported in the appendix in

Table A.2. We report the results from our empirical strategy after deflating expenditures in the second row of Table 5, and depict the relationship between the estimated changes and $\bar{P}_{1,y+1}^e$ in Figure 5. The resulting slope coefficient is almost identical to our baseline specification.²²

6.3 Strategic timing of medical care use

Another threat to our empirical approach could be strategic timing of medical care use.²³ Note, however, that individuals in our estimation sample do not have an incentive to strategically delay care since, for them, care at the beginning of year $y + 1$ is always more expensive than at the end of year y . Instead, they have an incentive for strategic frontloading of medical care, i.e., shifting care from the beginning of year $y + 1$ to the end of year y .

We do not expect frontloading to affect our results, for two reasons. First, unlike for delaying medical care use, the scope for frontloading is likely to be limited since people typically need a medical indication before they can get treatment. Second, the incentive to frontload is the same in all year-pairs $\{y, y + 1\}$. If there is a possibility to frontload then forward-looking patients should always do it. Therefore, we expect the amount of strategic-frontloading to be the same for all year-pairs. If individuals engage in strategic frontloading, then this will not invalidate our empirical approach since it is based on comparing changes in care consumption around the turn of the year across year-pairs. Hence, even though we do not formally model strategic frontloading, it will not affect our estimation results.

6.4 Robustness to changes in the donut hole

In our “donut hole” RD design we omit some days around the turn of the year due to holidays. For example, for the year-pair $\{2008, 2009\}$ we compare expenditures on Thursday December 18 and Thursday January 8. In robustness checks we change the length of the donut hole, and we compare averages across different weekdays. Specifically, we consider increasing the length of the donut hole by one week, and look at Wednesdays instead of Thursdays. Rows 3 to 5 of Table 5 show results for these different specifications of the donut hole. Our conclusions from Section 5 are robust to these different specifications.

6.5 Additional robustness checks

We also implement our empirical strategy for individuals under 18 years of age. Given that they do not face any cost-sharing, they should not exhibit the forward-looking behavior we document for our baseline sample of adults. This is exactly what we find. In fact, in line with minors not facing any cost-sharing, none of the estimated discontinuities are significantly different from zero.

²²Our findings are also robust to an alternative method of correcting for price changes based on a measure of healthcare inflation from Eurostat.

²³For instance, Cabral (2016) finds evidence for strategic delay of medical care use in the context of dental care.

We also repeat our empirical analysis for weekly-level data (as opposed to daily-level data in our baseline specification), a larger cut-off for pseudo-censored expenditures (€5000 instead of €500), and a specification where utilization is measured by the log of daily expenditures plus one. We find that our conclusions do not change.

The results of these robustness checks are presented in Appendix D.

7 The effect of changes in deductibles on annual expenditures

In the previous sections, we have examined the relationship between the expected end-of-year price and daily care utilization around the turn of the year. In this section, we show how our estimate of γ^e can be used to predict the effect of a change in the size of the deductible on healthcare utilization at the annual level. Specifically, we quantify the euro amount saved, per capita in a year, from an increase in the size of the deductible by €100—from €375 to €475.

Recall that according to (1), care consumption is given by:

$$c_{it} = \kappa_{it} - \gamma^c \cdot P_{it}^c - \gamma^e \cdot P_{it}^e.$$

This holds at the individual level at any time t within a year. We can take this as a starting point, aggregate over individuals, and make the dependence on the year explicit. y is our baseline year, 2015, with a €375 deductible. y' is a hypothetical year in which everything is the same, except that the size of the deductible is €475. We can write the resulting difference in average expenditures as

$$\bar{c}_{t,y'} - \bar{c}_{t,y} = - \left[\gamma^c \cdot (\bar{P}_{t,y'}^c - \bar{P}_{t,y}^c) + \gamma^e \cdot (\bar{P}_{t,y'}^c \cdot \bar{P}_{t,y'|P_{it,y'}^c > 0}^e - \bar{P}_{t,y}^c \cdot \bar{P}_{t,y|P_{it,y}^c > 0}^e) \right], \quad (3)$$

where $\bar{P}_{t,y|P_{it,y}^c > 0}^e$ is the average expected end-of-year price in year y for individuals who have a positive current price in year y at time t . $\bar{P}_{t,y'|P_{it,y'}^c > 0}^e$ is defined analogously. Appendix A.4 derives equation (3) and establishes the link to the micro foundation of (1).

(3) shows that the effect of a change in the deductible consists of two parts. $(\bar{P}_{t,y'}^c - \bar{P}_{t,y}^c)$ is the change in the fraction of the population of individuals for whom the current price is 1 in period t of the year, when the deductible is increased by €100. This is multiplied by the effect of the current price on care consumption, γ^c . The product of the two is the change that is due to the change in static incentives. The second part is the effect that is due to the change in dynamic incentives. It arises for all individuals who have a positive current price. The fraction of individuals for whom this is the case is given by $\bar{P}_{t,y}^c$ and $\bar{P}_{t,y'}^c$ in year y and y' respectively. These fractions are multiplied with the average expected end-of-year price conditional on a positive current price in each year. These quantities reflect the fact that when the deductible

increases by €100, both the expected end-of-year price and the fraction of individuals who face this expected end-of-year price increases.

We can decompose the second part of (3) into two separate effects: (i) the effect due to the change in the fraction of individuals who face dynamic incentives (prevalence effect) and (ii) the effect due to the change in the expected end-of-year price (intensity effect) when the deductible is increased by €100:

$$\begin{aligned} & \gamma^e \cdot (\bar{P}_{t,y'}^c \cdot \bar{P}_{t,y'|P_{it,y'}^c > 0}^e - \bar{P}_{t,y}^c \cdot \bar{P}_{t,y|P_{it,y}^c > 0}^e) \\ &= \gamma^e \cdot \left(\bar{P}_{t,y'|P_{it,y'}^c > 0}^e \cdot (\bar{P}_{t,y'}^c - \bar{P}_{t,y}^c) + \bar{P}_{t,y}^c \cdot (\bar{P}_{t,y'|P_{it,y'}^c > 0}^e - \bar{P}_{t,y|P_{it,y}^c > 0}^e) \right) \end{aligned} \quad (4)$$

The first term in (4) relates to the prevalence effect, i.e., the increase in the share of individuals who face dynamic incentives when the deductible increases. This increased share of individuals is given by $(\bar{P}_{t,y'}^c - \bar{P}_{t,y}^c)$ and is multiplied by the average expected end-of-year price in year y' . The second term relates to the intensity effect, i.e., the increase in the expected end-of-year price for individuals who face dynamic incentives in year y . The share of individuals who face dynamic incentives in year y is given by $\bar{P}_{t,y}^c$ and this is multiplied by the change in the average expected end-of-year price when the deductible is increased by €100.

For our prediction of annual expenditures, we use our estimate of γ^e . It measures the response to dynamic incentives. We have not estimated γ^c . Based on the literature, our starting point for the latter is that individuals reduce their expenditures by 40% if the current price increases from 0 to 1, conditional on the expected end-of-year price. 40% is very close to the estimates [Brot-Goldberg et al. \(2017\)](#) report for a sample of employees of a large firm in the U.S. and the estimates of [Hayen et al. \(2019\)](#) for the same data as we use in this paper. We also report results for a current price effect of 90%, 20%, and 0%, respectively.

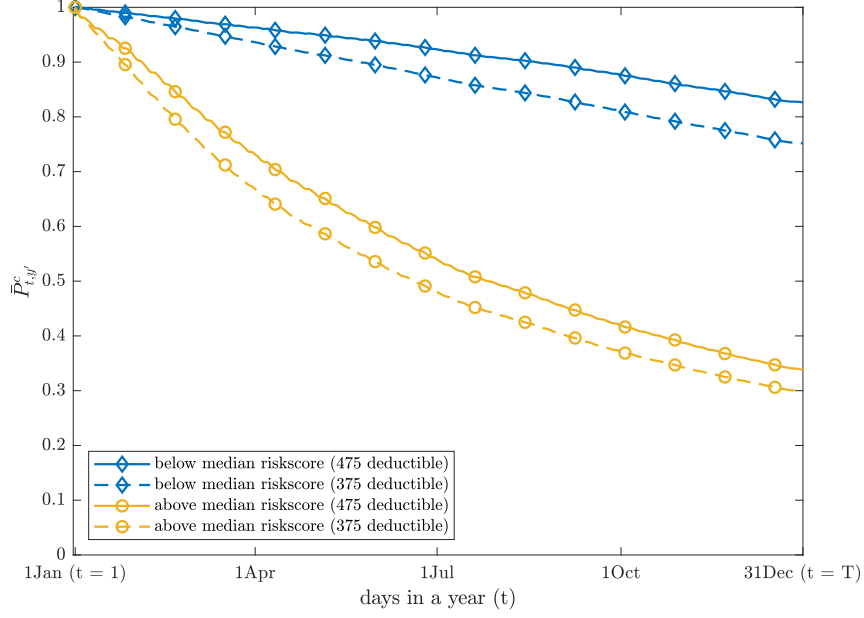
We also need to predict prices. To that end, for each day t , we regress the share of individuals who have crossed the deductible and the expected end-of-year price (conditional on not yet having crossed) on the deductible amount in the corresponding year. Based on this, we predict all 4 prices in (3). The estimation equation is

$$outcome_{t,y} = \theta_{0,t} + \theta_{1,t} deductible_y + \eta_{t,y},$$

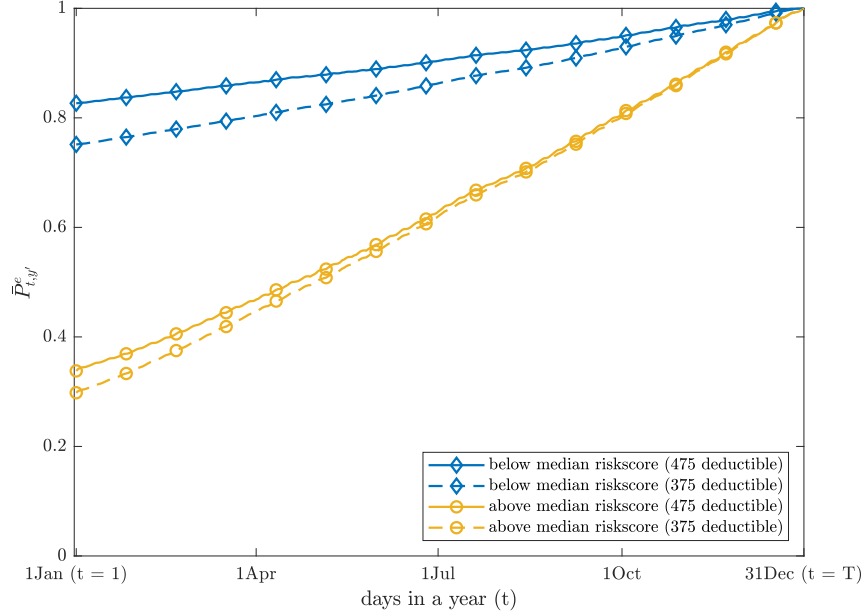
where $outcome_{t,y}$ is either $\bar{P}_{t,y}^c$ or $\bar{P}_{t,y|P_{it,y}^c > 0}^e$. Here, we use the entire sample of individuals in our data from 2008 – 2015. We account for the fact that individuals with higher riskscores react stronger to dynamic incentives by estimating γ^e separately for individuals above and below the median riskscore. We estimate γ^e to be -17.41 and -2.36 for above and below median riskscore groups respectively. We also account for the fact that these different groups may have different values of $\bar{P}_{t,y}^c$ and $\bar{P}_{t,y|P_{it,y}^c > 0}^e$ by performing the underlying regressions for both groups in our sample separately. Figure 6 shows our predictions of $\bar{P}_{t,y}^c$ and $\bar{P}_{t,y|P_{it,y}^c > 0}^e$, for each riskscore group and deductible level. Aggregation over days in the year gives spending effects on yearly

Figure 6: Effects of increasing the deductible by €100 on static and dynamic incentives

(a) Fraction individuals with current price of 1



(b) Expected end-of-year price for those who still have a current price of 1



Notes: Figure (a) shows the share of individuals with a positive current price at time t , $\bar{P}_{t,y}^c$, for two deductible levels, €375 and €475, across riskscore. Figure (b) depicts the predicted average expected end-of-year price for individuals with a positive current price at time t , $\bar{P}_{t,y}^e$, for two deductible levels, €375 and €475, across riskscore.

Table 6: Effect of increasing the deductible by €100 on annual expenditures

current price effect	reaction to static incentives	reaction to dynamic incentives		total effect
		prevalence effect	intensity effect	
90%	206.21	86.46	33.13	325.80
40%	91.65	86.46	33.13	211.24
20%	45.83	86.46	33.13	165.41
0%	0	86.46	33.13	119.59

Notes: This table presents annual spending reductions under an assumed current price effect. The last row serves as a baseline: when the current price effect is assumed to be 0, the only reduction in annual expenditures comes from changes in dynamic incentives.

healthcare expenditure.

Table 6 shows the results. Each row is for a different value of the current price effect. The first column shows the current price effect in percentage terms.²⁴ The second column shows the absolute effect on spending that is due to the change static incentives induced by increasing the deductible by €100. This is entirely due to patients facing a price of care of 1 for a longer time. The reduction in average annual expenditures that is due to current price changes is about €92 when we calculate it for our baseline value of the current price effect of 40%.

Next, column 3 and 4 report the effect of dynamic incentives. We find that overall, health-care expenditures decline by about €120 when the deductible size increases from €375 to €475. Given that, on average, an individual has total annual expenditures that count towards the deductible of around €2000 in 2015, this would imply increasing the deductible size by €100 to €475 leads to a 6% reduction in per capita healthcare expenditures due to changes in dynamic incentives.

The last column reports the total effect of dynamic and static incentives. For our baseline value of the current price effect of 40% we find that the reduction in average annual expenditures is €211, or 10.6%. Static incentives account for around half of the total reduction in annual spending.

8 Conclusion

In this paper, we show that a standard feature of deductible contracts, namely that they reset at the turn of the year, in combination with changes in deductible limits across years give rise

²⁴The effects reported in the literature are all relative effects. To conduct the analysis we translate relative effects into absolute effects. These are denoted by γ^c in (1). This is done by multiplying the respective relative effect by the average spending of individuals who have crossed the deductible in 2015. The average is taken over individuals and days after the deductible was hit and is €18.26. The absolute effects are €16.43 for 90%, €7.30 for 40%, and €3.65 for 20%.

to a differences-in-regression-discontinuities design that allows us to estimate the impact of dynamic incentives on healthcare utilization. Our micro-founded approach combines advantages of a model-based, structural approach and a reduced-form approach exploiting a natural experiment.

Using administrative data from the Netherlands, we find that individuals are forward-looking and that the effect of dynamic incentives on healthcare expenditures is quantitatively important: a €100 increase in the deductible reduces our measure of daily expenditures at the beginning of the year, when individuals still have to pay for care themselves, by around 3.0% and reduces the probability of having any claim by around 4.4%.

We also explore whether the reaction to dynamic incentives differs across subgroups in our sample, defined by age, gender, and neighborhood income. Policy makers may be concerned if individuals with the same healthcare needs would react differently to dynamic incentives because they belong to different groups. Controlling for differences in healthcare needs we find that almost every subgroup exhibits forward-looking behavior. The only exception is the group of individuals that are below age 45. For the remaining groups, we find that an increase in the size of the deductible by €100 leads to a percentage change in daily expenditures that ranges from -2.5% to -3.5%.

Based on our model, we can use our estimates to predict the effect of dynamic incentives on annual expenditures. At the annual level, dynamic incentives imply that an increase in the deductible by €100 reduces annual healthcare expenditures (that count towards the deductible) by 6%. For the Netherlands, in the year 2015, this is equivalent to a €2 billion reduction in overall healthcare expenditures. This means that dynamic incentives have a first-order impact on healthcare utilization.

In comparison, we predict that the response to static incentives reduces annual expenditures by 4.6% if the deductible is increased by €100. This prediction is based on an estimate of the effect of static incentives on healthcare expenditures from the literature. The relative size of the two effects suggests that patients' responses to dynamic incentives are an important part of the overall effect of cost-sharing schemes on healthcare expenditures—much more so than what the previous literature has suggested.

References

- Abaluck, J., J. Gruber, and A. Swanson (2018). Prescription drug use under Medicare Part D: A linear model of nonlinear budget sets. *Journal of Public Economics* 164, 106–138.
- Aron-Dine, A., L. Einav, A. Finkelstein, and M. Cullen (2015). Moral hazard in health insurance: Do dynamic incentives matter? *Review of Economics and Statistics* 97(4), 725–741.
- Barreca, A. I., M. Guldi, J. M. Lindo, and G. R. Waddell (2011). Saving babies? revisiting the effect of very low birth weight classification. *Quarterly Journal of Economics* 126(4), 2117–2123.
- Brot-Goldberg, Z. C., A. Chandra, B. R. Handel, and J. T. Kolstad (2017). What does a deductible do? The impact of cost-sharing on health care prices, quantities, and spending dynamics. *Quarterly Journal of Economics* 132(3), 1261–1318.
- Cabral, M. (2016). Claim timing and ex post adverse selection. *Review of Economic Studies* 84(1), 1–44.
- Calonico, S., M. D. Cattaneo, and R. Titiunik (2014). Robust nonparametric confidence intervals for regression-discontinuity designs. *Econometrica* 82(6), 2295–2326.
- Cutler, D. and R. Zeckhauser (2000). The anatomy of health insurance. In J. Newhouse and A. Culyer (Eds.), *Handbook of Health Economics*, Volume 1, Chapter 11, pp. 563 – 643. Elsevier Science.
- Dalton, C. M., G. Gowrisankaran, and R. J. Town (2020). Salience, Myopia, and Complex Dynamic Incentives: Evidence from Medicare Part D. *Review of Economic Studies* 78(2), 822–869.
- Einav, L. and A. Finkelstein (2018). Moral hazard in health insurance: What we know and how we know it. *Journal of the European Economic Association* 16, 957 – 982.
- Einav, L., A. Finkelstein, S. P. Ryan, P. Schrimpf, and M. R. Cullen (2013). Selection on moral hazard in health insurance. *American Economic Review* 103(1), 178–219.
- Einav, L., A. Finkelstein, and P. Schrimpf (2015). The response of drug expenditure to nonlinear contract design: Evidence from Medicare Part D. *Quarterly Journal of Economics* 130(2), 841–899.
- Ellis, R. P. (1986). Rational behavior in the presence of coverage ceilings and deductibles. *RAND Journal of Economics* 17(2), 158–175.
- Farbmacher, H., P. Ihle, I. Schubert, J. Winter, and A. Wuppermann (2017). Heterogeneous effects of a nonlinear price schedule for outpatient care. *Health Economics* 26(10), 1234–1248.
- Gerfin, M., B. Kaiser, and C. Schmid (2015). Healthcare demand in the presence of discrete price changes. *Health economics* 24(9), 1164–1177.
- Guo, A. and J. Zhang (2019). What to expect when you are expecting: Are health care consumers forward-looking? *Journal of Health Economics* 67, 102216.

- Hanushek, E. (1973). Efficient estimators for regressing regression coefficients. *The American Statistician* 28, 66–67.
- Hayen, A., T. Klein, and M. Salm (2019). Does the framing of patient cost-sharing incentives matter? The effects of deductibles vs. no-claim refunds. *C.E.P.R. Discussion Paper* 12908.
- Hayen, A. P., M. J. van den Berg, B. R. Meijboom, J. N. Struijs, and G. P. Westert (2015). Incorporating shared savings programs into primary care: From theory to practice. *BMC Health Services Research* 15(580), 1–15.
- Keeler, E. B., J. P. Newhouse, and C. E. Phelps (1977). Deductibles and the demand for medical care services: The theory of a consumer facing a variable price schedule under uncertainty. *Econometrica* 45(3), 641–655.
- Keeler, E. B. and J. E. Rolph (1988). The demand for episodes of treatment in the health insurance experiment. *Journal of Health Economics* 7(4), 337–367.
- Klein, T. J., M. Salm, and S. Upadhyay (2019). Proposition 1 by Ellis (1986) for a more general model with quasi-hyperbolic discounting and interest. Note, Tilburg University. Tilburg, The Netherlands.
- Manning, W., J. Newhouse, N. Duan, E. Keeler, A. Leibowitz, and M. Marquis (1987). Health insurance and the demand for medical care: Evidence from a randomized experiment. *American Economic Review* 77(3), 251–277.
- McGuire, T. (2011). Demand for health insurance. In T. M. M.V. Pauly and P. Barros (Eds.), *Handbook of Health Economics*, Volume 2, Chapter 5, pp. 317 – 396. Amsterdam: Elsevier Science.
- O’Donoghue, T. and M. Rabin (1999). Doing it now or later. *American Economic Review* 89(1), 103–124.
- Patton, A. and A. Timmermann (2010). Monotonicity in asset returns: New tests with applications to the term structure, the capm, and portfolio sorts. *Journal of Financial Economics* 98, 605 – 625.

Online Appendix

A Micro foundation

In this appendix, we provide a micro foundation for our differences-in-regression-discontinuities analysis and the counterfactual simulations. For this, we propose an empirical model that predicts a reaction to both the current and the expected end-of year price and features quasi-hyperbolic discounting (Appendix A.1). We derive the implied reduced-form equation for care consumption at the individual level (Appendix A.2). Then, we show that our differences-in-regression-discontinuities estimates have a particular structural interpretation (Appendix A.3). Thereafter, we provide a micro-foundation for our counterfactual experiment that we conduct in Section 7 of the paper (Appendix A.4). We also provide details on the relationship of our model to various versions of the [Keeler et al. \(1977\)](#) model that have been used in the literature, with a particular emphasis on the question what constitutes a good measure of dynamic incentives and to what extent our measure is model-dependent (Appendix A.5). We end with a more general discussion (Appendix A.6).

A.1 Model

The year is divided into periods $t = 1, \dots, T$. In each period, patient i faces healthcare needs λ_{it} that are i.i.d. draws from a time-specific distribution F_{λ_t} .

The model and the ensuing analysis are tailored to our setting where cost sharing is implemented using a standard deductible. At the beginning of each time period, i has a remaining deductible R_{it} and learns about her healthcare needs. She then chooses how much care c_{it} to consume. The remaining deductible is $R_{i1} = D$ at the beginning of the year and evolves according to $R_{it} = \max\{0, R_{i,t-1} - c_{i,t-1}\}$. Care consumption leads to out-of-pocket payments

$$C(c_{it}, R_{it}) = \min\{R_{it}, c_{it}\}. \quad (5)$$

It will be useful to define the out-of-pocket price for the last unit of care that is consumed in t as

$$P_{it}^c \equiv \frac{\partial C(c_{it}, R_{it})}{\partial c_{it}}.$$

Here, the superscript “c” stands for “current”.

Flow utility is quasi-linear in money and given by

$$u(c_{it}; \lambda_{it}, R_{it}) = (c_{it} - \lambda_{it}) - \frac{1}{2\omega} (c_{it} - \lambda_{it})^2 - C(c_{it}, R_{it}). \quad (6)$$

This quadratic functional form has been used by [Einav et al. \(2013\)](#) in a different context—patients choosing which health insurance to buy—and can be seen as a quadratic approximation to any utility function that is defined on the difference between healthcare consumption c_{it} and healthcare needs λ_{it} and quasi-linear in money. One advantage of this specification is that the parameter ω is readily interpretable. To see this, it is useful to inspect the first order condition for an

interior solution in the last period,²⁵

$$1 - \frac{1}{\omega} \cdot (c_{iT} - \lambda_{iT}) - \frac{\partial C(c_{iT}, R_{iT})}{\partial c_{iT}} = 0.$$

This implies that the optimal static consumption choice in the last period is

$$c_{iT}^* = \lambda_{iT} + \omega \cdot \left(1 - \frac{\partial C(c_{iT}, R_{iT})}{\partial c_{iT}} \right).$$

The last term in parentheses is the marginal out-of-pocket cost for the last unit of care that is consumed in period T . In the case of a standard deductible that we study in this paper, this cost is either 1 by the end of T the patient will have exceeded the deductible limit, or 0. So, care consumption is given by

$$c_{iT}^* = \lambda_{iT} + \omega$$

for patients above the deductible limit and

$$c_{iT}^* = \lambda_{iT}$$

for patients below the deductible limit. This means that ω is the additional care consumption when individuals do not have to pay out-of-pocket for the last unit of care they consume.

In any earlier period t patients solve a dynamic decision problem when they choose c_{it} . The associated value function is

$$V_t(\lambda_{it}, R_{it}) = \max_{c_{it}} u(c_{it}; \lambda_{it}, R_{it}) + \beta \delta \cdot \mathbb{E} [\tilde{V}_{t+1}(\lambda_{i,t+1}, R_{i,t+1})], \quad (7)$$

$$\text{with } \tilde{V}_t(\lambda_{it}, R_{it}) = \max_{c_{it}} u(c_{it}; \lambda_{it}, R_{it}) + \delta \cdot \mathbb{E} [\tilde{V}_{t+1}(\lambda_{i,t+1}, R_{i,t+1})]. \quad (8)$$

The distinction between $V_t(\lambda_{it}, R_{it})$ and $\tilde{V}_t(\lambda_{it}, R_{it})$ arises because of quasi-hyperbolic discounting with naive individuals: if $\beta < 1$ then they are too optimistic and wrongly expect that from the next period onward they will not suffer from present bias and only be exponential discounters.

The model outlined above is similar to the original theoretical model by [Keeler et al. \(1977\)](#). However, there are two key conceptual differences. First, patients can't save or borrow against future income. As we discuss in Section [A.5.1](#), this is the most straightforward and internally consistent way to build a model in which patients react to the current price. The second conceptual difference is a generalization: patients are quasi-hyperbolic discounters. They discount all

²⁵Keep in mind that the budget set is nonlinear. This means that the first order condition could hold at two values for c_{iT} . Also, notice that $C(c_{it}, R_{it})$ is not differentiable at $c_{it} = R_{it}$. However, (5) implies that it will in general not be optimal to choose the kink point. For this reason, we will abstract from this in the following for the ease of the exposition (but would have pointed it out in the relevant places if it would have led to different conclusions). In the context of Medicare Part D, as is well-understood, one would have to instead carefully take bunching at the kink into account, as has been done for instance by [Abaluck et al. \(2018\)](#).

future utilities by a factor β and, in addition, utility that is τ periods in the future by δ^τ . Related to the discounting, we assume that individuals are naive in the sense of [O'Donoghue and Rabin \(1999\)](#), meaning that they do not foresee that their time preference will change in the future.

Our model is similar to the empirical models in [Einav et al. \(2015\)](#) and [Abaluck et al. \(2018\)](#). Also, these models do not feature savings or borrowing. However, the former assumes $\beta = 1$ and the latter $\delta = 1$. We provide additional formal results for $\beta < 1$ and general values of δ .

We now derive the optimal policy in the form of a reduced-form equation that relates today's healthcare consumption c_{it} to the state variables λ_{it} and R_{it} , and to beliefs about the future.

These beliefs are directly related to the probability of crossing the deductible in future periods. It will be useful to define

$$q_{it}(\tau) \equiv \Pr(R_{i\tau} > 0, c_{i\tau} > R_{i\tau} | R_{it}, c_{it}), \text{ where } \tau > t.$$

In words, this is the probability of exceeding the cost sharing limit in period τ , which refers to the joint event that the remaining deductible at the beginning of period τ is positive and that care consumption in period τ exceeds the remaining deductible. This probability is conditional on information available in t , so it is from the perspective of period t . In particular, we condition on the deductible at the beginning of period t and spending in t . This implies a value for the remaining deductible at the end of period t , which is the key state variable here. Notice that, in line with rational expectations over future healthcare needs and behavior in the future, this probability $q_{it}(\tau)$ can be estimated from data.

With this probability $q_{it}(\tau)$, a consumer who is spending one additional unit of money out-of-pocket in t will spend one unit of money less in τ . As utility is quasi-linear in money, the utility cost of this additional unit of care will then be $1 - \beta\delta^{\tau-t}$. 1 is the cost in the current period when utility is quasi-linear in money and $\beta\delta^{\tau-t}$ are the discounted savings in τ . These discounted savings refer to the chance that higher expenditures in period t might lead to lower expenditures in period τ . If it is uncertain when the individual will exceed the deductible, then the patient has to form expectations about this. The expected utility cost of one additional unit of care, will then be $P_{it} = 1 - \beta \cdot \sum_{\tau=t+1}^T \delta^{\tau-t} q_{it}(\tau)$. P_{it} depends on individual preference parameters β and δ , and it involves beliefs, as it depends on future realization of healthcare needs. It is the correct measure of dynamic incentives in this context.

Furthermore, $P_{it} = 0$ whenever $P_{it}^c = 0$. The expected utility cost of one additional unit of care consumption is zero whenever a patient has already exceeded the deductible limit before the end of year t . Therefore, we can write P_{it} as

$$P_{it} = \begin{cases} 0 & \text{if } P_{it}^c = 0 \\ 1 - \beta \cdot \sum_{\tau=t+1}^T \delta^{\tau-t} q_{it}(\tau) & \text{if } P_{it}^c = 1. \end{cases}$$

A.2 Reduced-form equation for care consumption

We now derive a reduced-form equation for care consumption at the individual level. We summarize our result in the following proposition. The proof resembles the proof by [Ellis \(1986\)](#), as discussed in [Appendix A.5](#) below, and the one by [Abaluck et al. \(2018\)](#) who however impose $\delta = 1$.

Proposition 1. *In the model described in [Appendix A.1](#), optimal care consumption is given by*

$$c_{it} = \lambda_{it} + \omega \cdot (1 - P_{it}), \quad (9)$$

where the relevant price is given by

$$P_{it} = \begin{cases} 0 & \text{if } P_{it}^c = 0 \\ 1 - \beta \cdot \sum_{\tau=t+1}^T \delta^{\tau-t} q_{it}(\tau) & \text{if } P_{it}^c = 1. \end{cases} \quad (10)$$

Proof of Proposition 1. The first order condition for the maximization problem in (7) is

$$\frac{\partial u(c_{it}; \lambda_{it}, R_{it})}{\partial c_{it}} + \beta \delta \cdot \frac{\partial \mathbb{E} [\tilde{V}_{t+1}(\lambda_{i,t+1}, R_{i,t+1})]}{\partial R_{i,t+1}} \cdot \frac{\partial R_{i,t+1}}{\partial c_{it}} = 0. \quad (11)$$

First consider the case when $R_{it} = 0$. In this case, patients solve a static problem, and the second term on the left hand side of (11) is zero. Hence, the first order condition in any period for which $R_{it} = 0$ is given by the derivative of the flow utility in (6) with respect to c_{it} , evaluated at $R_{it} = 0$, being equal to zero,

$$1 - \frac{1}{\omega} \cdot (c_{it} - \lambda_{it}) = 0,$$

which implies $c_{it} = \lambda_{it} + \omega$ for $R_{it} = 0$. By (10) we have that $P_{it} = 0$. So, (9) holds for all t whenever $R_{it} = 0$.

Next consider the case $R_{it} > 0$. Our goal is to show that

$$\frac{\partial \mathbb{E} [\tilde{V}_s(\lambda_{i,s}, R_{i,s})]}{\partial R_{i,s}} = - \sum_{\tau=s}^T \delta^{\tau-s} q_{it}(\tau), \quad (12)$$

for any $t < s \leq T$. If (12) holds for $s = t + 1$ then it can be shown that the first order condition in (11) implies that equation (9) in the proposition is true.

We show that (12) holds by induction. First, we show that it holds for $s = T$. In that period, patients reach the cost sharing limit with probability $q_{it}(T)$ from the perspective of period t . For patients who reach the cost-sharing limit in period T , a remaining deductible that is higher by one unit means that they spend that one unit more out-of-pocket. By the envelope theorem, the derivative of (8) with respect to the remaining deductible in T is

$$\frac{\partial \mathbb{E} [\tilde{V}_T(\lambda_{i,T}, R_{i,T})]}{\partial R_{i,T}} = (-1) \cdot q_{it}(T) + 0 \cdot (1 - q_{it}(T)).$$

So, (12) holds for $s = T$.

It remains to show that (12) holds for s whenever it holds for $s + 1$. By the envelope theorem, we have that the derivative of (8) with respect to R_{is} is

$$\begin{aligned}\frac{\partial \mathbb{E} [\tilde{V}_s(\lambda_{is}, R_{is})]}{\partial R_{is}} &= \frac{\partial u(c_{is}; \lambda_{is}, R_{is})}{\partial R_{is}} + \delta \cdot \frac{\partial \mathbb{E} [\tilde{V}_{s+1}(\lambda_{i,s+1}, R_{i,s+1})]}{\partial R_{i,s+1}} \\ &= -\frac{\partial C(c_{is}, R_{is})}{\partial R_{is}} + \delta \cdot \frac{\partial \mathbb{E} [\tilde{V}_{s+1}(\lambda_{i,s+1}, R_{i,s+1})]}{\partial R_{i,s+1}}.\end{aligned}$$

We have that

$$-\frac{\partial C(c_{is}, R_{is})}{\partial R_{is}} = \begin{cases} -1 & \text{with probability } q_{it}(s) \\ 0 & \text{with probability } 1 - q_{it}(s). \end{cases}$$

Using this and substituting in (12) gives

$$\frac{\partial \mathbb{E} [\tilde{V}_s(\lambda_{is}, R_{is})]}{\partial R_{is}} = (-1) \cdot q_{it}(s) + \delta \cdot \sum_{\tau=s+1}^T \delta^{\tau-(s+1)} \cdot (-1) \cdot q_{it}(\tau) = -\sum_{\tau=s}^T \delta^{\tau-s} q_{it}(\tau).$$

This completes the proof. \square

Before we discuss how equation (9) in the proposition relates to our analysis we write it as

$$\begin{aligned}c_{it} &= \lambda_{it} + \omega \cdot \left[1 - P_{it}^c \cdot \left(1 - \beta \cdot \sum_{\tau=t+1}^T \delta^{\tau-t} q_{it}(\tau) \right) \right] \\ &= \lambda_{it} + \omega \cdot \left[1 - P_{it}^c \cdot \left(1 - \beta \cdot \left(\sum_{\tau=t+1}^T \delta^{\tau-t} q_{it}(\tau) - 1 + 1 \right) \right) \right] \\ &= \lambda_{it} + \omega \cdot \left[1 - P_{it}^c \cdot \left(1 - \beta \cdot 1 - \beta \cdot \left(\sum_{\tau=t+1}^T \delta^{\tau-t} q_{it}(\tau) - 1 \right) \right) \right] \\ &= \lambda_{it} + \omega \cdot \left[1 - \left((1 - \beta) \cdot P_{it}^c + \beta \cdot P_{it}^c \cdot \left(1 - \sum_{\tau=t+1}^T \delta^{\tau-t} q_{it}(\tau) \right) \right) \right]\end{aligned}$$

where the first equality follows from substituting in the expression for P_{it} , the second from adding and subtracting 1 from the expression in the innermost set of parentheses, the third and fourth from rearranging terms.

As a next step, we show that this way to write the optimal policy relates care consumption c_{it} to the current price P_{it}^c and the expected end-of-year price P_{it}^e . To show this, we assume $\delta = 1$. In the main text, we have defined the expected end of year price as the probability that the patient will have to pay for the last unit of care in the year. Formally, we can write

$$P_{it}^e = P_{it}^c \cdot \left(1 - \sum_{\tau=t+1}^T q_{it}(\tau) \right).$$

If $P_{it}^c = 0$, then the expected end-of-year price is 0; if $P_{it}^c = 1$, then the expected end of year price is 1 minus the probability to hit the deductible limit in any future period. So, we have

$$c_{it} = \lambda_{it} + \omega \cdot [1 - ((1 - \beta) \cdot P_{it}^c + \beta \cdot P_{it}^e)]. \quad (13)$$

This way to re-write the optimal policy is useful because it forms a basis for a regression of care consumption on the current and the expected end-of-year price. Such a regression has been carried out for instance by [Brot-Goldberg et al. \(2017\)](#). Furthermore, this equation is equivalent to the first equation in the main text

$$c_{it} = \kappa_{it} - \gamma^c \cdot P_{it}^c - \gamma^e \cdot P_{it}^e.$$

where $\kappa_{it} \equiv \lambda_{it} + \omega$, the parameter on the current price is $\gamma^c \equiv \omega \cdot (1 - \beta)$ and the parameter on the expected end-of-year price is $\gamma^e \equiv \omega \cdot \beta$.²⁶ The relative importance of γ^c and γ^e depends on the parameter for hyperbolic discounting β . If there is no hyperbolic discounting ($\beta = 1$), then patients will exclusively respond to the expected end of year price P_{it}^e . This is the case considered in the original paper by [Keeler et al. \(1977\)](#). If there is hyperbolic discounting ($\beta < 1$) then patients will respond to both P_{it}^e and P_{it}^c .

Finally, observe that the relative size of the coefficients on the current and expected end-of-year price is informative about the extent to which individuals discount all future periods, as summarized by β . In particular, we have that

$$\frac{\gamma^c}{\gamma^e} = \frac{1 - \beta}{\beta}.$$

or

$$\beta = \frac{1}{1 + \gamma^c / \gamma^e}.$$

This means that if γ^c and γ^e are identified, then also β is identified.

A.3 Interpretation of differences-in-regression-discontinuities estimates

In the main part of the paper, we relate changes in care consumption at the turn of the year to the expected end-of-year price. Here we provide a micro foundation. We do so under the assumption that $\delta = 1$. This allows us to follow the literature and measure changes in dynamic incentives by changes in the expected end-of-year price. In [Appendix A.5](#) we further discuss the assumption that $\delta = 1$, and we show that even if we depart from this assumption P_{it}^e is closely related to the correct measure of dynamic incentives.

²⁶This result for $\delta = 1$ is not new. [Abaluck et al. \(2018, p.110\)](#) provide a similar expression.

Denote averages taken over individuals in a given period t of year y by a bar indexed by t and y . Using this notation, we have that the average current price in $t = 1$ of year $y + 1$ is $\bar{P}_{1,y+1}^c$ and the expected end-of-year price at the beginning of the year is $\bar{P}_{1,y+1}^e$.²⁷

Based on (13) we can write average care consumption in $t = 1$ of year $y + 1$ as

$$\bar{c}_{1,y+1} = \bar{\lambda}_{1,y+1} + \omega \cdot [1 - ((1 - \beta) \cdot \bar{P}_{1,y+1}^c + \beta \cdot \bar{P}_{1,y+1}^e)].$$

For each year-pair we use a sample of individuals whose price is zero by the end of year y . In terms of our model average care consumption in the last period of year y is thus

$$\bar{c}_{T,y} = \bar{\lambda}_{T,y} + \omega.$$

As discussed in Section (2), our identifying assumption in (2) is that the change in medical needs around the turn of the year is the same for all year-pairs. This assumption can also be expressed in terms of average medical needs $\bar{\lambda}_{t,y}$ instead of $\bar{\kappa}_{t,y}$. Then we have for all year-pairs $\{y, y + 1\}$ and $\{y + 1, y + 2\}$

$$(\bar{\lambda}_{1,y+2} - \bar{\lambda}_{T,y+1}) = (\bar{\lambda}_{1,y+1} - \bar{\lambda}_{T,y})$$

We discuss the plausibility of this assumption in Section 6. Under this assumption, the change in the discontinuity from year-pair $\{y, y + 1\}$ to, year-pair $\{y + 1, y + 2\}$ is

$$(\bar{c}_{1,y+2} - \bar{c}_{T,y+1}) - (\bar{c}_{1,y+1} - \bar{c}_{T,y}) = -\omega \cdot [(1 - \beta) \cdot (\bar{P}_{1,y+2}^c - \bar{P}_{1,y+1}^c) + \beta \cdot (\bar{P}_{1,y+2}^e - \bar{P}_{1,y+1}^e)]. \quad (14)$$

At the beginning of the year, the current price is one because the deductible resets at the turn of the year. At the end of the first period, it can be zero if an individual experiences a health shock in the first period and consumes more care than the deductible limit. However, in Section 6.1 we show that the difference $\bar{P}_{1,y+2}^c - \bar{P}_{1,y+1}^c$ does not significantly vary across years. Thus, our estimate for γ^e in (1) is an estimate of $\omega \cdot \beta$.

A.4 Micro-foundation for counterfactuals

In Section 7 we show how we can use an estimate of $\omega \cdot \beta$ to make a counterfactual prediction of healthcare expenditures for a different value of the deductible. We also show how one can use prior knowledge of the current price effect to say how much of the change in healthcare expenditures is driven by the current price effect and how much is driven by the reaction to dynamic incentives. In this section, we provide a micro foundation for this.

²⁷Here we do not make a formal distinction between expectations and averages. In our analysis, we interpret estimates that are based on sample averages and are interested in the effect of a change in the deductible on average expenditures.

Maintaining the assumption that $\delta = 1$, it follows from (9) and (10) that

$$\begin{aligned}\bar{c}_{t,y} &= \bar{\lambda}_{t,y} + (1 - \bar{P}_{t,y}^c) \cdot \omega + \bar{P}_{t,y}^c \cdot \omega \cdot \left[1 - \left(1 - \beta \cdot (1 - \bar{P}_{t,y|P_{it,y}^c > 0}^e) \right) \right] \\ &= \bar{\lambda}_{t,y} - \bar{P}_{t,y}^c \cdot \omega \cdot \left(1 - \beta \cdot (1 - \bar{P}_{t,y|P_{it,y}^c > 0}^e) \right),\end{aligned}$$

where $\bar{P}_{t,y|P_{it,y}^c > 0}^e$ is the average expected end-of-year price from the perspective of period t , where the average is taken over individuals with a positive current price in period t of year y ; analogously for y' . Now consider the case in which the only difference between year y and year y' is the size of the deductible. The difference in the expenditure between those two years is

$$\begin{aligned}\bar{c}_{t,y'} - \bar{c}_{t,y} &= - \left[\bar{P}_{t,y'}^c \cdot \omega \cdot \left(1 - \beta \cdot (1 - \bar{P}_{t,y'|P_{it,y'}^c > 0}^e) \right) - \bar{P}_{t,y}^c \cdot \omega \cdot \left(1 - \beta \cdot (1 - \bar{P}_{t,y|P_{it,y}^c > 0}^e) \right) \right] \\ &= - \left[(\bar{P}_{t,y'}^c - \bar{P}_{t,y}^c) \cdot \omega \cdot (1 - \beta) + \omega \cdot \beta \cdot (\bar{P}_{t,y'}^c \cdot \bar{P}_{t,y'|P_{it,y'}^c > 0}^e - \bar{P}_{t,y}^c \cdot \bar{P}_{t,y|P_{it,y}^c > 0}^e) \right]\end{aligned}$$

This can also be written in a way that corresponds to the exposition in Section 7. For this let y' be the year with the higher deductible and denote by $\bar{P}_{t,y|P_{it,y}^c > 0}^e$ the average expected end-of-year price in period t of year y (alternatively y') for individuals who have a positive current price in period t of year y' . Using

$$\bar{P}_{t,y'}^c \cdot \bar{P}_{t,y|P_{it,y'}^c > 0}^e = \bar{P}_{t,y}^c \cdot \bar{P}_{t,y|P_{it,y}^c > 0}^e,$$

we have

$$\bar{c}_{t,y'} - \bar{c}_{t,y} = - \left[(\bar{P}_{t,y'}^c - \bar{P}_{t,y}^c) \cdot \omega \cdot (1 - \beta) + \bar{P}_{t,y'}^c \cdot (\bar{P}_{t,y'|P_{it,y'}^c > 0}^e - \bar{P}_{t,y|P_{it,y}^c > 0}^e) \cdot \omega \cdot \beta \right].$$

We can see that the effect consists of two parts. $(\bar{P}_{t,y'}^c - \bar{P}_{t,y}^c)$ is the change in the fraction of the population of individuals for whom the current price is 1. This is multiplied by the effect of the current price on care consumption, $\omega \cdot (1 - \beta)$, and together gives the change that is due to the change in static incentives. The second part is the effect that is due to the change of dynamic incentives. It arises for all individuals who have a positive current price in year y' , i.e. the fraction $\bar{P}_{t,y'}^c$, and is given by the change in the end-of-year price for those individuals times the reaction to the end-of-year price, $\omega \cdot \beta$.

A.5 Relevant measure of dynamic incentives across models

A.5.1 The model by Keeler et al. (1977) and the result by Ellis (1986)

We first describe a version of the Keeler et al. (1977)-Ellis (1986) model that is comparable to ours. The main generalization in our model is that patients are quasi-hyperbolic $\beta - \delta$ discoun-

ters.

Flow utility is additive in healthcare consumption and other consumption y_{it} . Instead of (6), we now have

$$u(c_{it}, y_{it}; \lambda_{it}) = (c_{it} - \lambda_{it}) - \frac{1}{2\omega} (c_{it} - \lambda_{it})^2 + u_c(y_{it}). \quad (15)$$

Notice that the costs of healthcare spending, unlike in (6), are not included in (15). These costs, instead, enter the optimization problem through the evolution of wealth across periods:

$$W_{t+1} = W_t - y_t - C(c_{it}, R_{it})$$

Allowing for quasi-hyperbolic discounters and interest, r , on wealth, we have the following Bellman equation:

$$\begin{aligned} V_{it}(\lambda_{it}, W_{it}) &= \max_{c_{it}, y_{it}} u(c_{it}, y_{it}; \lambda_{it}) + \beta \delta \cdot \mathbb{E} [\tilde{V}_{t+1}(\lambda_{i,t+1}, W_{i,t+1})] \\ \text{with } \tilde{V}_t(\lambda_{it}, W_{it}) &= \max_{c_{it}, y_{it}} u(c_{it}, y_{it}; \lambda_{it}) + \delta \cdot \mathbb{E} [\tilde{V}_{t+1}(\lambda_{i,t+1}, W_{i,t+1})] \end{aligned} \quad (16)$$

such that

$$\begin{aligned} W_{it+1} &= (1 + r)(W_{it} - y_t - C(c_{it}, R_{it})) \\ R_{it+1} &= \max(R_{it} - c_{it}, 0) \\ V_{T+1} &= V(\lambda_{iT+1}) + V(W_{iT+1}) \end{aligned}$$

The main difference between (7) and (16) is that the costs of healthcare do not affect flow utility and only affects future values through its effect on W_{it+1} . An important related assumption, given by the terminal condition on V_{T+1} , is that there are no wealth effects in this formulation: the marginal utility of wealth in period $T + 1$ is assumed to be equal to 1. With this, one can show that instead of (11) the first order condition is²⁸

$$\frac{u(c_{it}, y_{it}; \lambda_{it})}{\partial c_{it}} - \beta \cdot \delta^{(T+1)-t} \cdot (1 + r)^{(T+1)-t} \cdot \left(1 - \sum_{\tau=1}^T q_{it}(\tau)\right) = 0. \quad (17)$$

From this we get that optimal care consumption is given by

$$c_{it} = \lambda_{it} + \omega \left(1 - \beta \delta^{(T+1)-t} (1 + r)^{(T+1)-t} \cdot \left(1 - \sum_{\tau=1}^T q_{it}(\tau)\right)\right).$$

The following proposition summarizes the above.

Proposition 2 (Relevant price in Ellis (1986) with discounting and interest). *Assume the model*

²⁸In a separate note, Klein et al. (2019), we derive this result using a more general version of the original setup and the same notation as in Ellis (1986). This note is not meant for publication, but available upon request from the authors.

setup in A.1, with the exception that: (i) patients are endowed with wealth and can freely save and borrow; (ii) utility is not quasi-linear in money, but given by (15); (iii) money earns interest at rate r . Assume that

$0 < \beta \delta^{(T+1)-t} (1+r)^{(T+1)-t} < \infty$. Then, optimal care consumption is

$$c_{it} = \lambda_{it} + \omega \cdot (1 - P_{it}^{KNPE}) \quad (18)$$

and the relevant price is given by

$$P_{it}^{KNPE} \equiv \beta \delta^{(T+1)-t} (1+r)^{(T+1)-t} \cdot \left(1 - \sum_{\tau=1}^T q_{it}(\tau) \right). \quad (19)$$

Here, “KNPE” abbreviates the names of the authors of Keeler et al. (1977) and Ellis (1986). The original result is often cited for showing that the expected end-of-year price is the only relevant price a patient should act on. Proposition 2 shows that this continues to hold under quasi-hyperbolic discounting with interest.

A.5.2 Using the expected end-of-year price as a measure of dynamic incentives

Putting Proposition 1 and 2 side-by-side reveals that in the general case, the relevant price differs across models. It is

$$P_{it} = \begin{cases} 0 & \text{if } P_{it}^c = 0 \\ 1 - \beta \cdot \sum_{\tau=t+1}^T \delta^{\tau-t} q_{it}(\tau) & \text{if } P_{it}^c = 1. \end{cases}$$

in our model and

$$P_{it}^{KNPE} = \beta \delta^{(T+1)-t} (1+r)^{(T+1)-t} \cdot \left(1 - \sum_{\tau=1}^T q_{it}(\tau) \right)$$

in the model by Keeler et al. (1977) and Ellis (1986) with quasi-hyperbolic discounting. This is not surprising, as one can in general not summarize complex dynamic incentives using a scalar measure. Also not surprisingly, there is no difference between the two models once patients exhaust the cost sharing limit. Then, the relevant price is zero.

In either of the two models a patient pays an effective price that is given by out-of-pocket costs minus a “bonus”. The bonus reflects that spending an additional euro today reduces the remaining deductible for all remaining periods by a euro: as a result of higher spending today, the patient might have to pay less in the future. In fact, the patient will only have to pay less in the future if she crosses the deductible in the future; the probability of such an event occurring at time τ , from the perspective of time t , is given by $q_{it}(\tau)$. Thus, a consumer who spends one additional unit of money out-of-pocket in t will spend one unit of money less in τ with probability $q_{it}(\tau)$.

In our model, spending one euro less in period τ , means not incurring a disutility of 1 euro

in period τ with probability $q_{it}(\tau)$. This disutility arises because costs are modeled into flow utility. From the perspective of period t this reduction in disutility in period τ is worth (in expectation) $\beta \delta^{\tau-t} q_{it}(\tau)$. The value of the bonus, in period t , of spending an additional unit is then just the summation of expected disutility reductions over all remaining τ s in the year:

$$b_{it} \equiv \beta \cdot \sum_{\tau=t+1}^T \delta^{\tau-t} q_{it}(\tau).$$

In the model by [Keeler et al. \(1977\)](#) and [Ellis \(1986\)](#), the inclusion of wealth (with no wealth effects) means that spending money today only affects how much wealth the individual carries over to the next period. The only relevant question, thus, is whether or not the patient will have to pay for the last unit of care within a year. So, the bonus is given by the probability that this is the case, times 1. The value of the bonus, from the perspective of t , is thus

$$b_{it}^{KNPE} \equiv \beta \delta^{(T+1)-t} (1+r)^{(T+1)-t} \cdot \sum_{\tau=t+1}^T q_{it}(\tau).$$

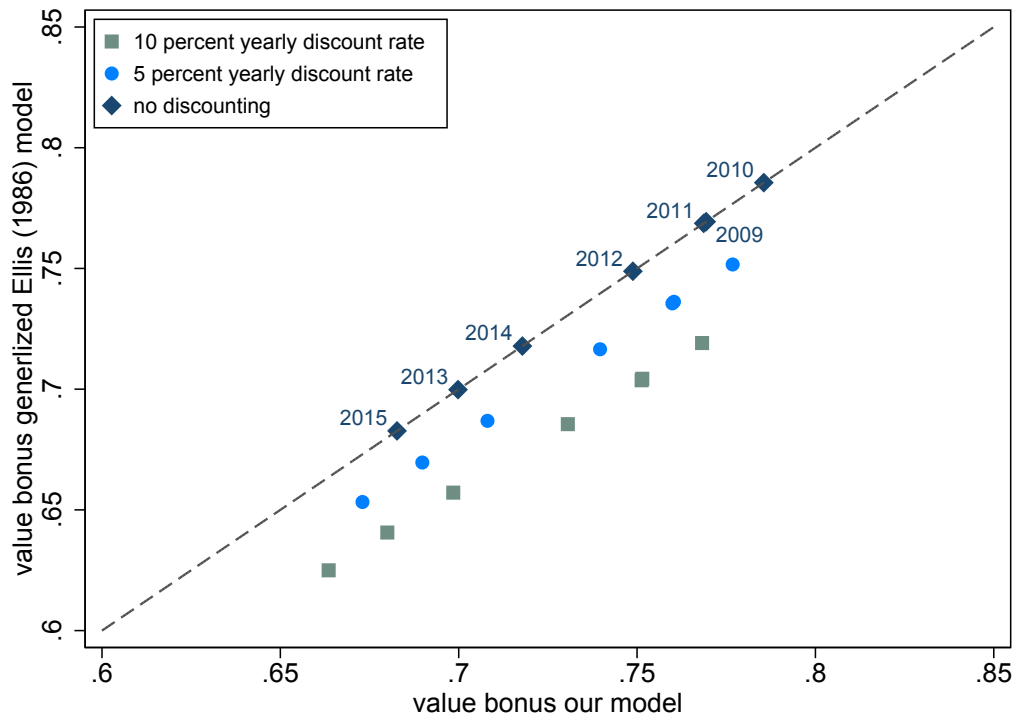
The bonus is the same in both models when $\delta = 1$ and $r = 0$, and the difference between the two is small for the realistic case in which δ is close to 1 and r is small.

When the bonus is the same, then the expected end-of-year price \bar{P}^e is also the same. In our analysis, we relate variation in care consumption to variation in \bar{P}_{it}^e for a given period t . Those changes in \bar{P}_{it}^e are in turn directly related to changes in the true dynamic incentive, b_{it} or b_{it}^{KNPE} . There is no difference between those when $\delta = 1$ and $r = 0$. Also, in both models, according to (9) and (18), the coefficient we would estimate if we would regress healthcare consumption on \bar{P}_{it}^e is $\beta \cdot \omega$, so a reduced-form parameter that measures the reaction to changes in dynamic incentives would have the same meaning in both models.

This remains to hold approximately if $\delta \neq 1$, as long as δ is not too far away from 1 and r is small. To show this, we set $r = 0$ and calculate the average value of the relevant part of the value of the bonus in our model, $\sum_{\tau=t+1}^T \delta^{\tau-t} q_{it}(\tau)$, for January of each year (so for $t = 1$) and using the value of $q_{it}(\tau)$ from our data, and plot it against the relevant part of the value of the bonus in the generalized [Keeler et al. \(1977\)](#)-[Ellis \(1986\)](#) model, $\delta^{-1} \cdot \delta^{(T+1)-t} \cdot \sum_{\tau=t+1}^T q_{it}(\tau)$. We do so for various values of δ . Here, we add the factor δ^{-1} to make the timing comparable: in our model individuals receive the last at most 11 periods in the future, while in the model with savings they formally receive it 12 periods in the future. The factor δ^{-1} thus makes the two models more comparable. We use $\delta = 1$ for reference, $\delta = 0.996$ that corresponds to a 5 percent yearly discount rate, and $\delta = 0.992$ that corresponds to a 10 percent yearly discount rate.

Figure A.1 shows the result. There are 7 years in our data, which means that for each value of δ we obtain 7 data points. We can see that for each of the 3 values of δ , the values of the bonus from both models are co-monotonic. This means that the ordering is preserved.

Figure A.1: Value of bonus for different models and different values of δ



Note: This figure plots the average value of $\sum_{\tau=2}^{12} \delta^{\tau-1} q_{il}(\tau)$, which is directly related to the value of the bonus in our model, against $\delta^{12-1} \sum_{\tau=2}^{12} q_{il}(\tau)$. Done for 3 different discount factors. See text for details.

Moreover, we can see that even for a yearly discount rate of 10 percent, the two values are very close. A regression reveals that a one unit increase in the value of the bonus in our model predicts a 0.894 unit increase in the value of the bonus for the [Ellis \(1986\)](#) model.

A.6 Discussion

In this appendix we have provided a micro foundation for the differences-in-regression-discontinuities analysis we have conducted in the main part of the paper. This shows that our estimate of γ^e has the interpretation of a myopia parameter β times a moral hazard parameter ω . We then show how, based on this, one can perform a counterfactual simulation without solving a structural model.

We have also shown that if we are interested in measuring dynamic incentives, then we will not necessarily have to take a stance on which particular model we prefer. As long as we use a model with quasi-hyperbolic discounting and assume that δ is close to 1 and r is close to zero, our results show that the expected end-of-year “future” price that is also used in the literature based on referencing the particular model by [Keeler et al. \(1977\)](#) is a good measure of dynamic incentives.

Finally, our analysis reveals that a key difference between our model and the original model with wealth by [Keeler et al. \(1977\)](#) is that in our model patients do react to static incentives, with a weight of $(1 - \beta)$, while in the model with wealth this is not the case. This difference is not essential for the purpose of measuring the extent to which individuals react to dynamic incentives, but it is a desirable model property as it is more in line with the empirical finding in [Keeler and Rolph \(1988\)](#) and [Brot-Goldberg et al. \(2017\)](#) that patients react to static incentives even conditional on the expected end-of-year price \bar{P}_{it}^e . The counterfactual analysis we conduct in Section 7 provides a motivation to prefer the model that we describe in the beginning of this Appendix, as it can be used to study the empirically relevant current price effects within the same model framework. Otherwise, we would predict that patients do not react at all to cost-sharing when they are close to being fully myopic, i.e. for small β that are however strictly positive.²⁹

Finally, we show that prior knowledge on the current price effect—which is interesting in many applied contexts, but not the focus of this paper—can be used to separately identify moral hazard effects ω and the myopia parameter β .

B Additional tables and figures

This appendix contains additional tables and figures referred to in the text.

²⁹In the model with savings, the current price does actually matter when patients are fully myopic, but only then. But this means that the limit of c_{it} as a function of β , for $\beta \rightarrow 0$ is not the same as c_{it} when β is exactly zero.

Table A.1: Dates used in analysis

year-pair	Date Before ToY	Date After ToY
2008 – 2009	18 th December 2008	8 th January 2009
2009 – 2010	17 th December 2009	7 th January 2010
2010 – 2011	16 th December 2010	6 th January 2011
2011 – 2012	15 th December 2011	12 th January 2012
2012 – 2013	20 th December 2012	10 th January 2013
2013 – 2014	19 th December 2013	9 th January 2014
2014 – 2015	18 th December 2014	8 th January 2015

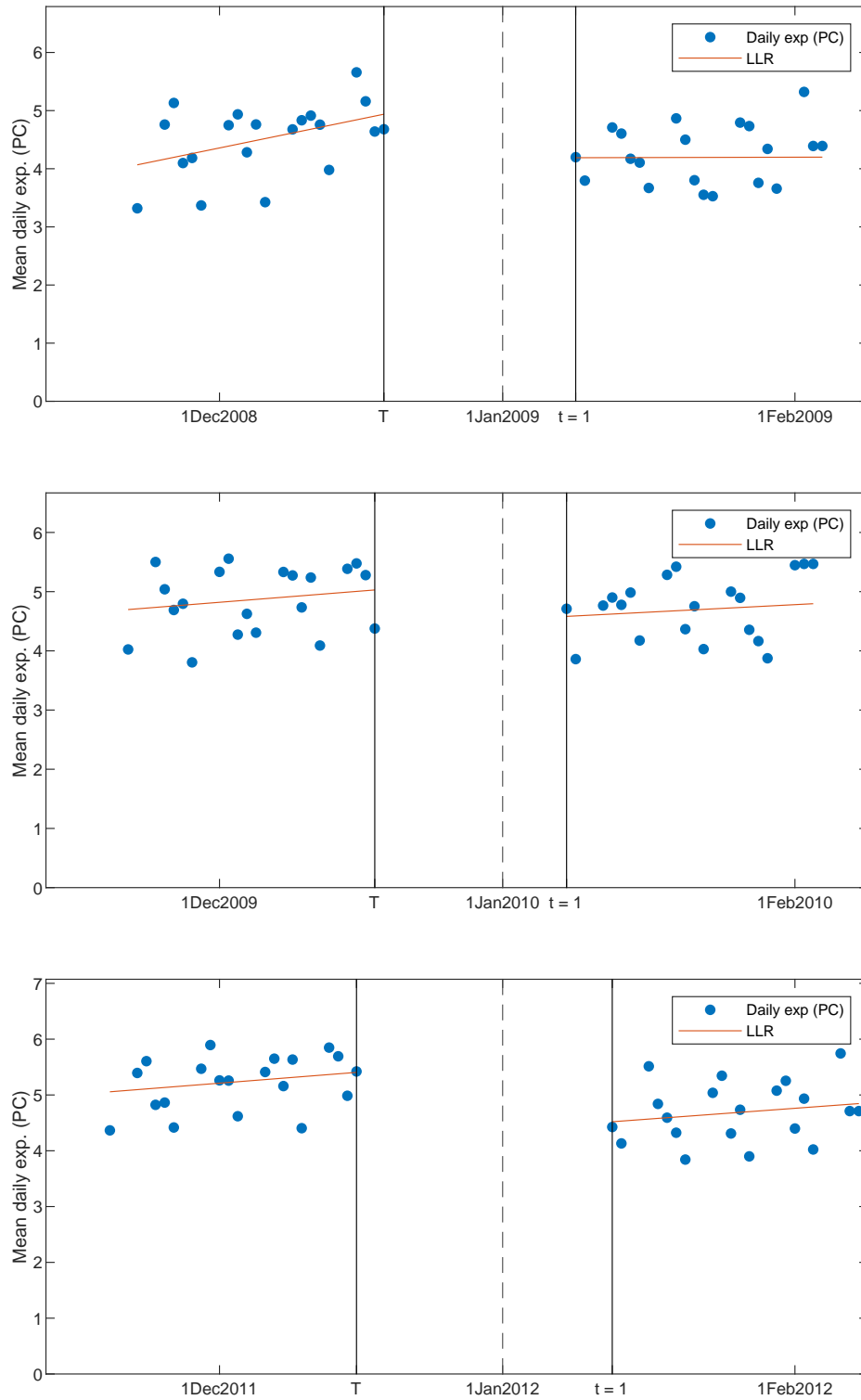
Note: ToY refers to turn of year. The relevant holidays before the turn of the year is the 24th of December, while the relevant holiday after the turn of the year is the 1st of January.

Table A.2: Expenditure deflator values

Year	Deflator
2008	.665
2009	.719
2010	.739
2011	.776
2012	.831
2013	.908
2014	.925
2015	1

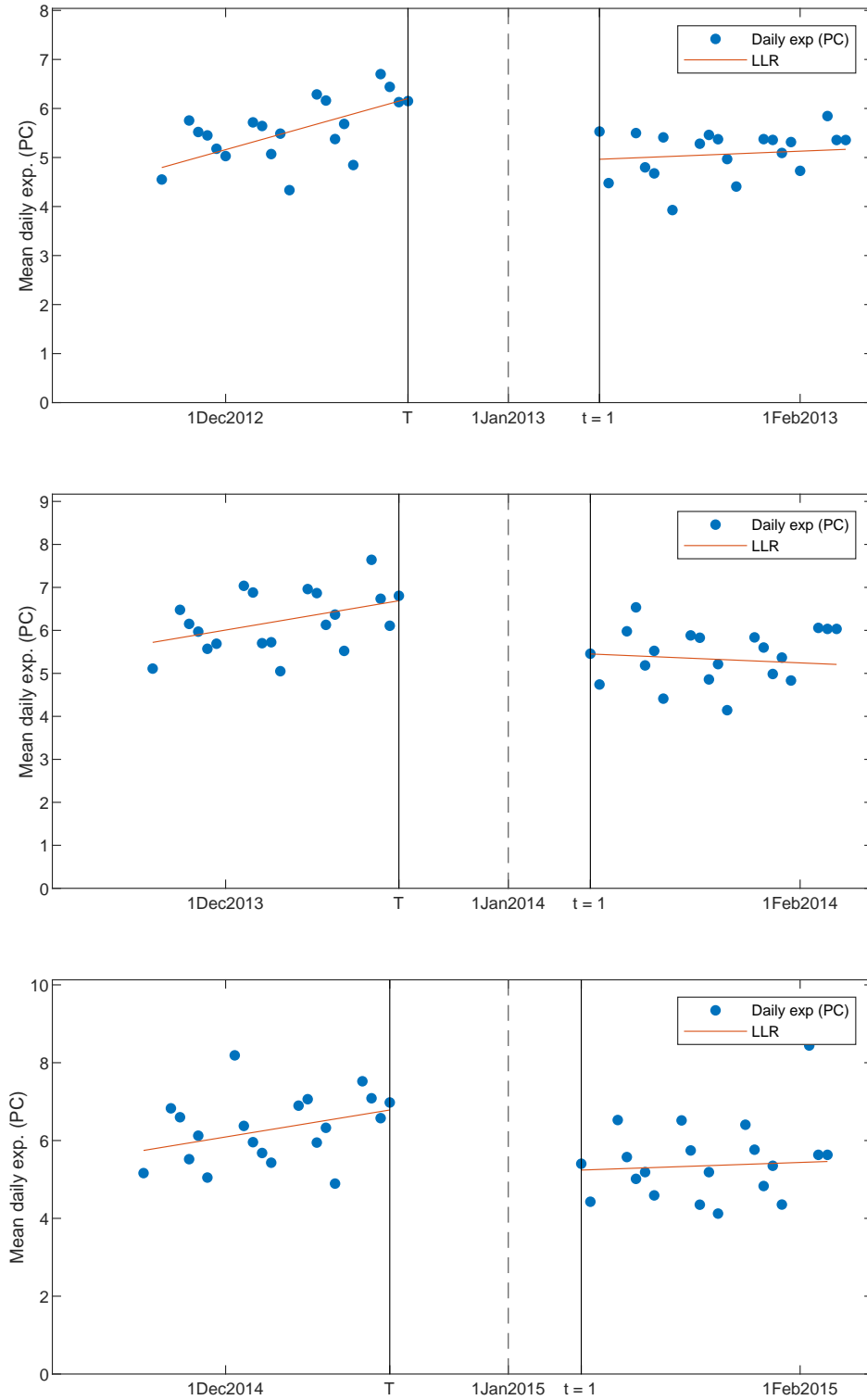
Notes: Deflator values obtained from mean expenditures (only weekdays) from the 37th week of the year to the end of the year, when individuals in our sample face no price of care.

Figure A.2: RD figures for different years



Notes: The figure plots mean daily pseudo-censored healthcare spending (denoted by the blue dots) across days for different years as shown in the x-axis. T denotes the last day of year y and $t = 1$ denotes the first day of year $y + 1$ we use in our estimation of the change in spending due to contracts resetting. The red solid line denotes the local linear regression used for the main analysis. See Section 4.1 for details. Weekends are omitted from our analysis.

Figure A.3: RD figures for different years



Notes: The figure plots mean daily pseudo-censored healthcare spending (denoted by the blue dots) across days for different years as shown in the x-axis. T denotes the last day of year y and $t = 1$ denotes the first day of year $y + 1$ we use in our estimation of the change in spending due to contracts resetting. The red solid line denotes the local linear regression used for the main analysis. See Section 4.1 for details. Weekends are omitted from our analysis.

Figure A.4: Results by gender

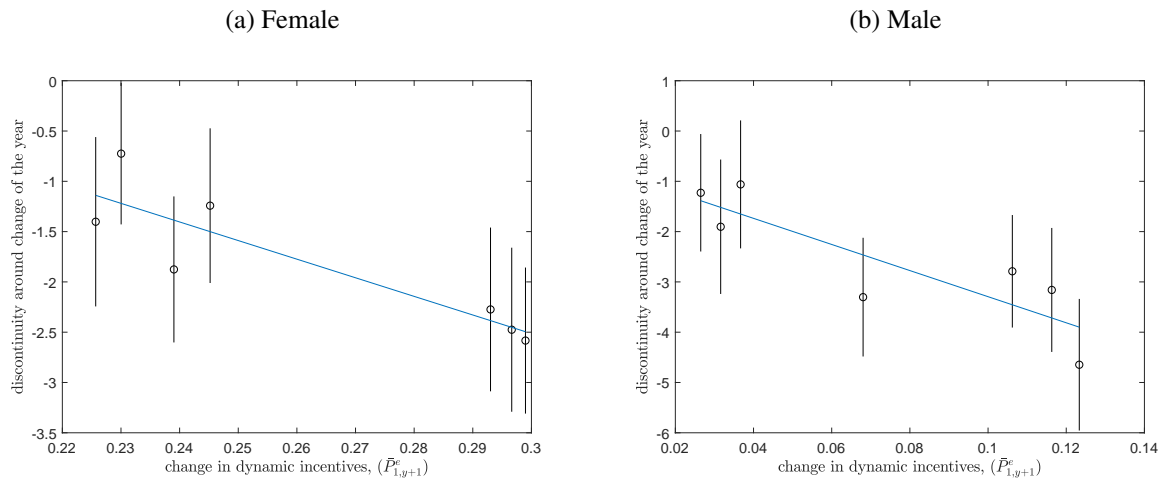


Figure A.5: Results by income group

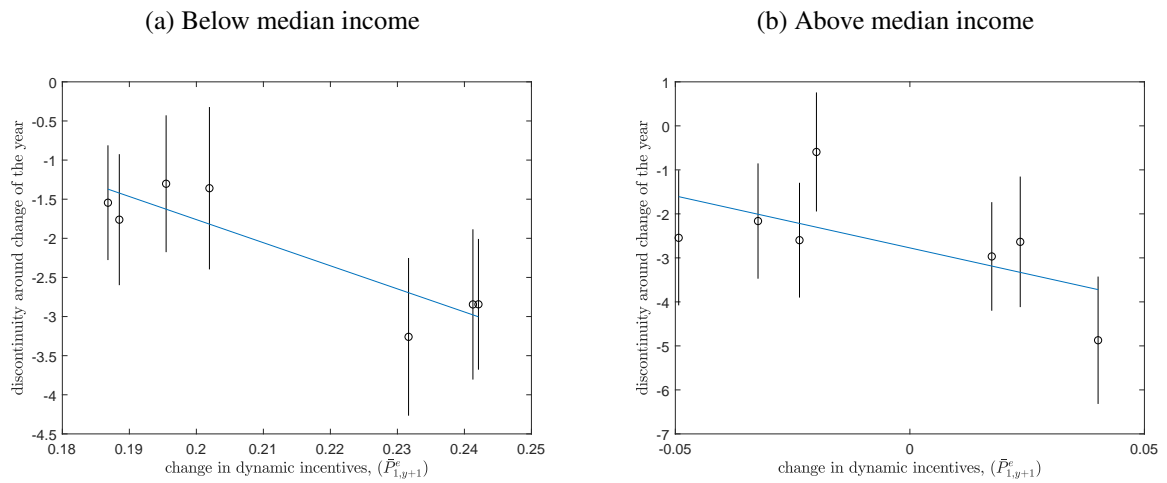
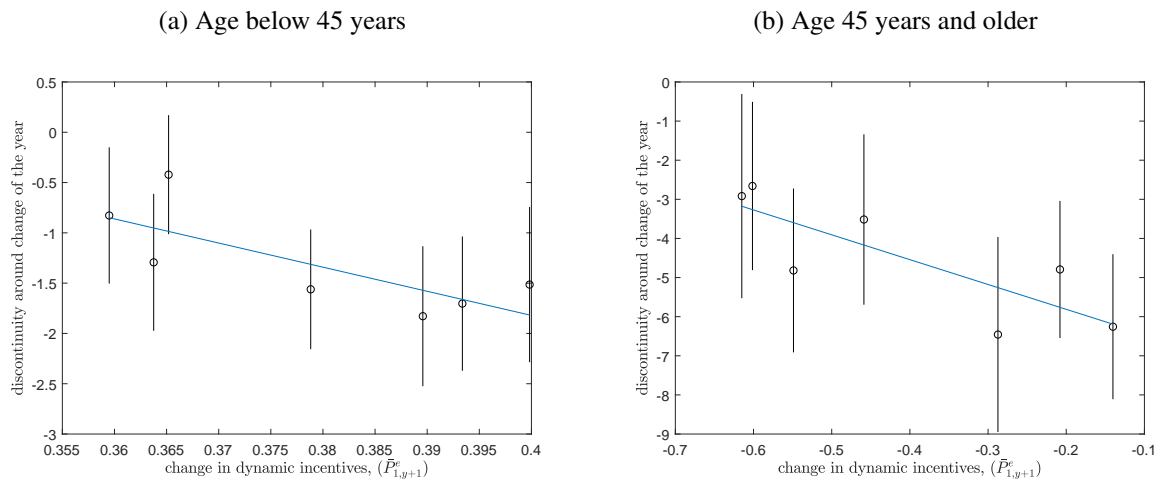


Figure A.6: Results by age



C Inference

C.1 Relating discontinuity sizes to changes in dynamic incentives

Our parametric test proceeds in two steps: estimating the discontinuities around the turn of the year and then relating them to the expected end-of-year price. The measurement error in the first step needs to be taken into account when calculating the standard errors in the second step. Following [Hanushek \(1973\)](#), the regression equation for the second step is given by:

$$\hat{\Delta}_y = \alpha + \beta \bar{P}_{1,y}^e + (v_y + u_y) \quad (20)$$

where $\hat{\Delta}_y$ is the estimated change in care consumption for year-pair $\{y-1, y\}$, u_y is the sampling error that arises from using the estimates $\hat{\Delta}_y$ as the dependent variable and v_y is the error term for the regression equation that uses Δ_y as the dependent variable. We assume v_y and u_y are independent, v_y is homoskedastic with $Var(v_y) = \sigma^2$, $Var(u_y) = e_y^2$ and $Cov(u_y, u_{y'}) = 0 \forall y \neq y'$. $\sqrt{e_y^2}$ corresponds to the standard error of the estimated discontinuity $\hat{\Delta}_y$. Let $\varepsilon_y = v_y + u_y$. It can be shown that the variance-covariance matrix of ε_y is given by:

$$\mathbb{E}[\varepsilon \varepsilon'] = \sigma^2 \Omega = \begin{bmatrix} 1 + \frac{e_1^2}{\sigma^2} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & 1 + \frac{e_Y^2}{\sigma^2} \end{bmatrix} \sigma^2$$

In order to efficiently estimate the parameters β , we first run an OLS regression of $\hat{\Delta}_y$ on $\bar{P}_{1,y}^e$ and then use the residuals from that regression to obtain an estimate for $\hat{\sigma}$. [Hanushek \(1973, p. 67\)](#) shows that $\hat{\sigma}$ can be estimated by the following equation:

$$\hat{\sigma}^2 = \frac{s^2(n-1) - \sum_i e_i^2 - tr(X'X)^{-1}X'GX}{n-1}$$

where s^2 is the variance of the residuals from the first OLS regression, n is the number of observations, X is the vector of $\bar{P}_{1,y}^e$ and G is defined as $\sigma^2(I - \Omega)$. tr indicates the trace of the resulting matrix.

We then use this estimate of $\hat{\sigma}^2$ to run a weighted least-squares regression with weights, w_y , given by:

$$w_y = \frac{1}{\sqrt{\hat{\sigma}^2 + e_y^2}}$$

C.2 Non-parametric monotonicity tests

We use the MR and Up/Down tests that have been proposed in [Patton and Timmermann \(2010\)](#) as non-parametric tests for monotonicity. [Patton and Timmermann \(2010\)](#) argue that these test statistics can provide strong evidence for an economic theory. It was also the only non-

parametric test of monotonicity we could find that made use of the estimated coefficients rather than just their ranking.

Consider our 7 estimates of changes in expenditures around the turn of the year, $\hat{\Delta}_y$, ordered with respect to the expected end-of-year price, with the smallest expected end-of-year price first and denote them by $\tilde{\Delta}_1, \dots, \tilde{\Delta}_7$. Economic theory posits that years with a higher expected end-of-year price should exhibit a larger decrease in percentage spending. This implies:

$$\tilde{\Delta}_1 > \tilde{\Delta}_2 > \dots > \tilde{\Delta}_7. \quad (21)$$

Let d_i denote the difference between consecutive ranks i and $i + 1$ in (21). Consider the minimum of all d_i 's. If this minimum is larger than 0, it must be the case that every other consecutive difference is larger than 0. This provides the basis for the MR test statistic:

$$J_T = \min_i \{d_i\} \quad (22)$$

with the corresponding hypothesis as:

$$H_0 : d = 0$$

$$H_1 : d \geq 0$$

where d is a vector containing all the d_i 's from above.³⁰

Since the asymptotic distribution of J_T does not have a closed form solution, we have to obtain p -values via bootstrapping. This is done, for each bootstrap repetition, by drawing from the distribution of the estimated coefficients, ordering them with respect to the expected end-of-year price, and then computing J_T^b , where J_T^b is the value of the MR test statistic for bootstrap repetition b . To impose the null hypothesis, we then subtract the quantity in (22), J_T , from J_T^b . The p -value is given by the following:

$$\frac{1}{B} \sum_{b=1}^B 1 \left\{ \left(J_T^b - J_T \right) > J_T \right\}$$

A problem with this MR test is that it can lack power, as shown in Monte Carlo simulations done in [Patton and Timmermann \(2010\)](#). For our specific application, this lack of power stems from 2 reasons:

1. There is a cardinality in the ranking that is not taken into account. For example, the change in expected end-of-year price from 2012 to 2013 is rather large, but the test statistic treats all of these changes the same.

³⁰Although we use only the difference between consecutive ranks in our application, this approach is valid for any difference implied by theory. [Patton and Timmermann \(2010\)](#) report very little gains from including all possible combinations of differences implied by theory in the testing procedure.

2. The magnitude of the change in coefficients is not taken into account. For example, the estimated coefficients changes by quite a large amount from 2012 to 2013.

To diagnose whether the test statistic fails to reject the null of no monotonic relationship due to a lack of power, [Patton and Timmermann \(2010\)](#) propose the Up and Down test statistic, given by the following:

$$J_{Up} = \sum_{i=1}^N |d_i| \{d_i > 0\} \quad (23)$$

$$J_{Down} = \sum_{i=1}^N |d_i| \{d_i < 0\} \quad (24)$$

Intuitively, the test statistic in (22) checks if there are significant changes in line with the theory, taking into account the magnitude of these changes, whereas (24) checks the opposite. This addresses the 2nd power concern for the MR test we mentioned above. We obtain *p*-values for both (22) and (24) using the same bootstrap procedure mentioned for the MR test.

D Additional robustness checks

D.1 Placebo test

In the Netherlands only individuals above the age of 18 (inclusive) are subject to cost-sharing in the form of deductibles, while the price of care for individuals below the age of 18 is always zero. Individuals below the age of 18, thus, provide a simple placebo test for our empirical methodology—since they do not face any price of care for all the years in consideration, their changes in healthcare utilization should not exhibit an increasing relationship with $\bar{P}_{1,y+1}^e$.³¹ It is important to mention that dental care was removed from estimating the changes from resetting contracts since it is not covered in the basic package for our main sample.

Figure A.7 depicts the relationship between the estimated coefficients and $\bar{P}_{1,y+1}^e$, while Table A.3 reports the tests for monotonicity.

All four tests, for both measures of healthcare utilization, suggest a lack of forward-looking behavior in this particular subgroup.

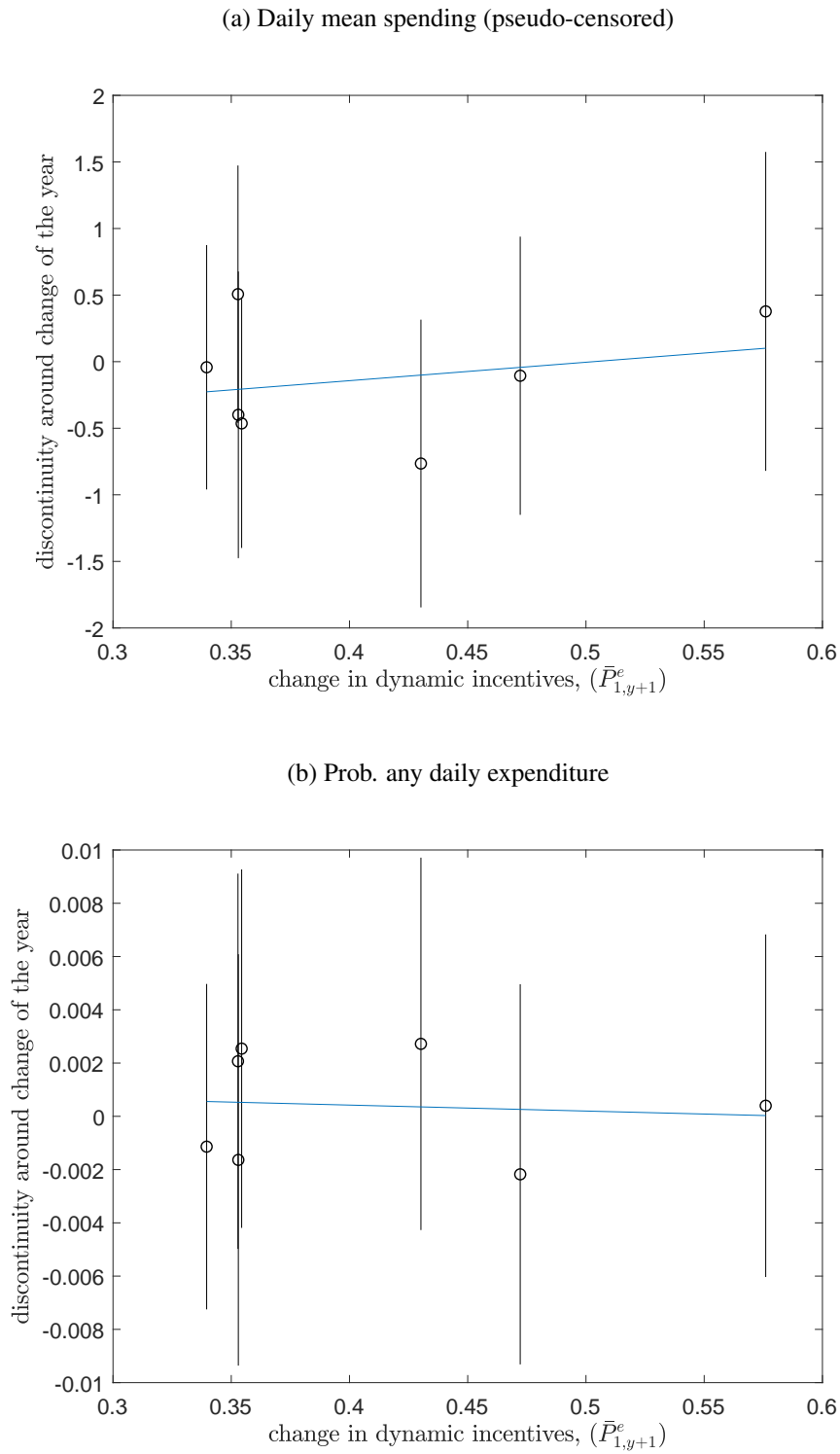
D.2 Weekly data

We also run our empirical approach on weekly level data. We create this weekly level data set through the following steps:

1. We create the daily level final data set, with all the relevant dates for the donut hole.

³¹Since these individuals do not face the deductible, we compute their $\bar{P}_{1,y+1}^e$ as the 1 minus the proportion of total individuals that would have crossed the actual deductible in year $y + 1$.

Figure A.7: Relationship between discontinuity sizes and dynamic incentives (placebo)



Notes: These figures plot the relationship between the estimated changes due to resetting contracts and the computed $\bar{P}_{1,y+1}^e$ for the placebo sample (individuals aged below 18). The solid line is the OLS regression line.

Table A.3: Dependence of discontinuity sizes on dynamic incentives (placebo)

	daily expenditure	any daily expenditure
effect of dynamic incentives (γ^e)	1.215 (2.3187)	-0.002 (0.0105)
p -value nonparametric MR test	0.219	0.180
p -value nonparametric Up test	0.737	0.739
p -value nonparametric Down test	0.511	0.620

Notes: The effect of dynamic incentives was estimated by regressing the estimated discontinuities on $\bar{P}_{1,y+1}^e$. The “Slope” test refers to a two-sided t-test on the slope coefficient of a weighted regression of the estimated changes reported in Table 2 and the computed $\bar{P}_{1,y+1}^e$. The MR, Up and Down tests were conducted using 10,000 bootstrap repetitions. See Appendix C for details.

Table A.4: Dependence of discontinuity sizes on dynamic incentives (weekly)

	daily expenditure	any daily expenditure
effect of dynamic incentives (γ^e)	-52.027 (7.8074)	-0.250 (0.1084)
p -value nonparametric MR test	0.099	0.987
p -value nonparametric Up test	0.004	0.000
p -value nonparametric Down test	0.981	0.029

Notes: The effect of dynamic incentives was estimated by regressing the estimated discontinuities on $\bar{P}_{1,y+1}^e$. The “Slope” test refers to a two-sided t-test on the slope coefficient of a weighted regression of the estimated changes reported in Table 2 and the computed $\bar{P}_{1,y+1}^e$. The MR, Up and Down tests were conducted using 10,000 bootstrap repetitions. See Appendix C for details.

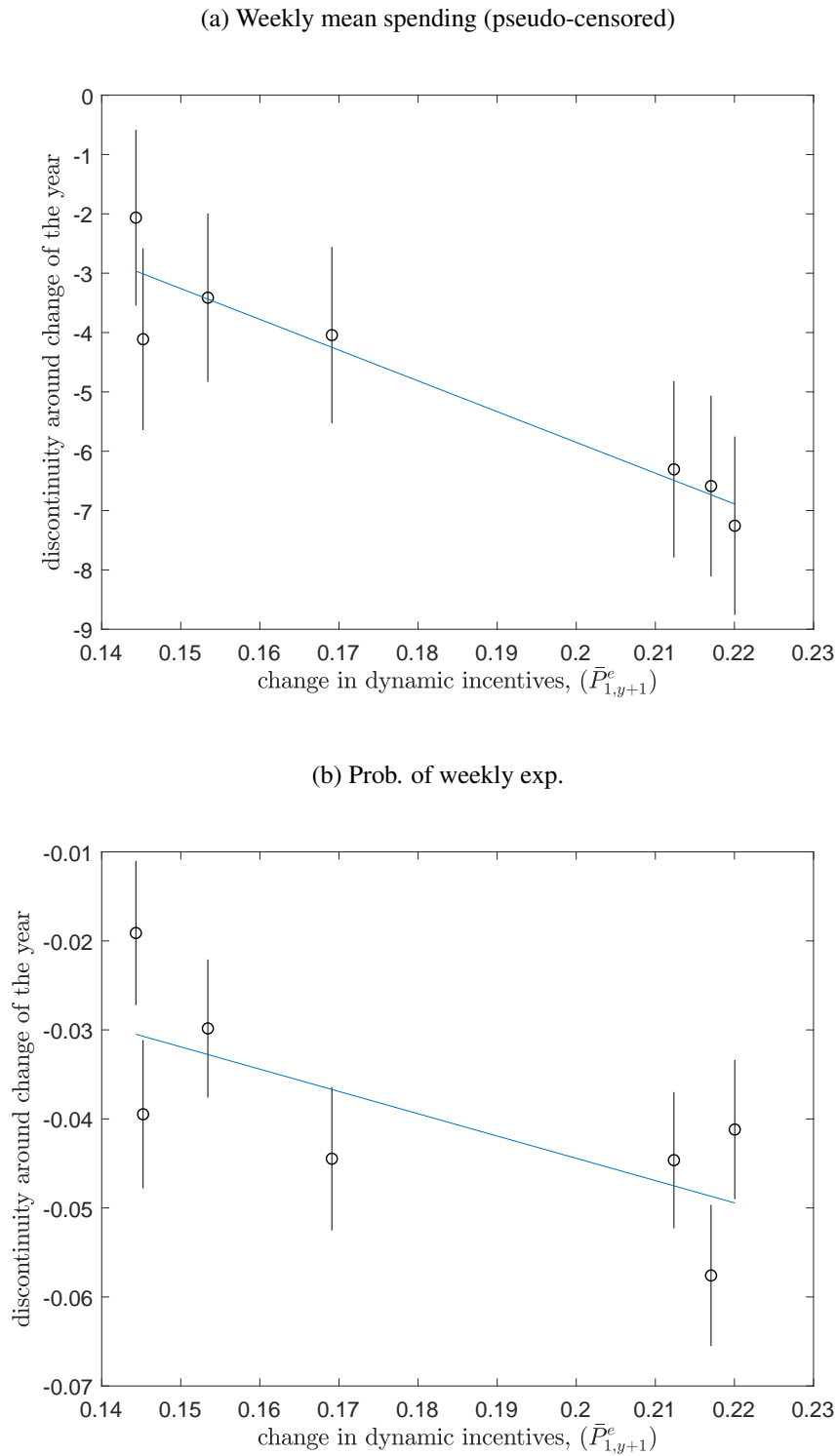
2. We start counting 7 days from the day before the turn of the year (our T in the RD). This is one week for days before the turn of the year.
3. We start counting 7 days from the day after the turn of the year (our $T + 1$ in the RD). This is one week for days after the turn of the year.
4. We sum expenditures within these 7 days, NOT taking into account weekend expenditures.

The results are reported in Figure A.8 and Table A.4.

D.3 Pseudo-censored cutoff

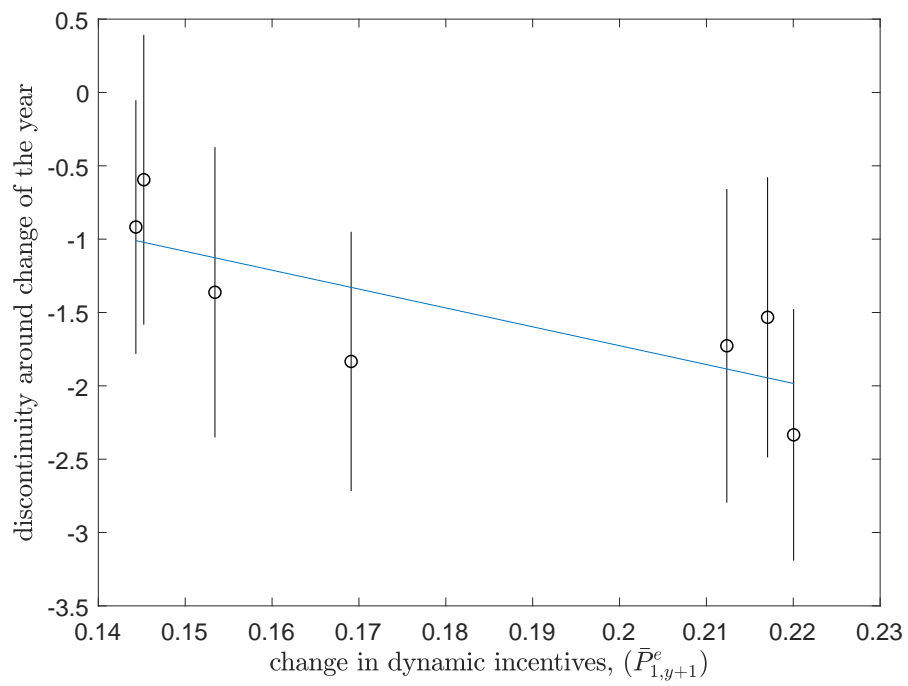
In our baseline specification, we pseudo-censored expenditures at €500. In this section, we report results from pseudo-censoring expenditures at €5000. The results are reported in Figure A.9 and Table 5 respectively.

Figure A.8: Relationship between discontinuity sizes and $\bar{P}_{1,y+1}^e$ (weekly data)



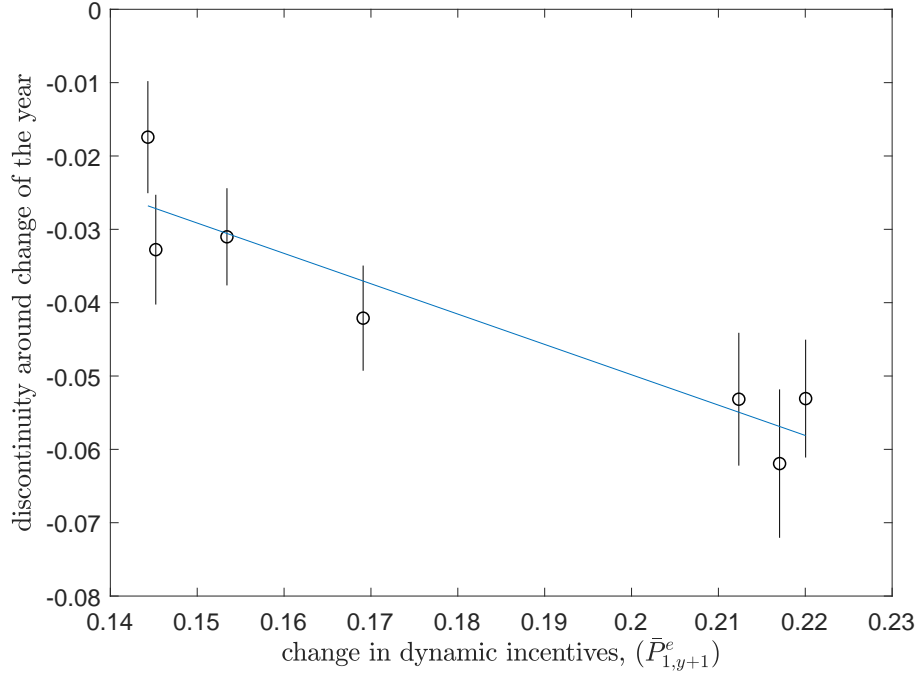
Notes: These figures plot the relationship between the estimated changes due to resetting contracts and the computed $\bar{P}_{1,y+1}^e$ for the baseline sample, where claims were aggregated at the weekly level. The solid line is the OLS regression line.

Figure A.9: Relationship between jump and $\bar{P}_{1,y+1}^e$ (PC at 5000)



Notes: These figures plot the relationship between the estimated changes due to resetting contracts and the computed $\bar{P}_{1,y+1}^e$ for the baseline sample as defined in Section ???. The dependent variable here is pseudo-censored at 5000, i.e., any expenditure above 5000 was coded to be equal to 5000. The solid line is the OLS regression line.

Figure A.10: Relationship between discontinuity size and $\bar{P}_{1,y+1}^e (\ln(\exp + 1))$



Notes: These figures plot the relationship between the estimated changes due to resetting contracts and the computed $\bar{P}_{1,y+1}^e$ for the baseline sample as defined in Section ?? . The dependent variable here is logged expenditure plus 1. The solid line is the OLS regression line.

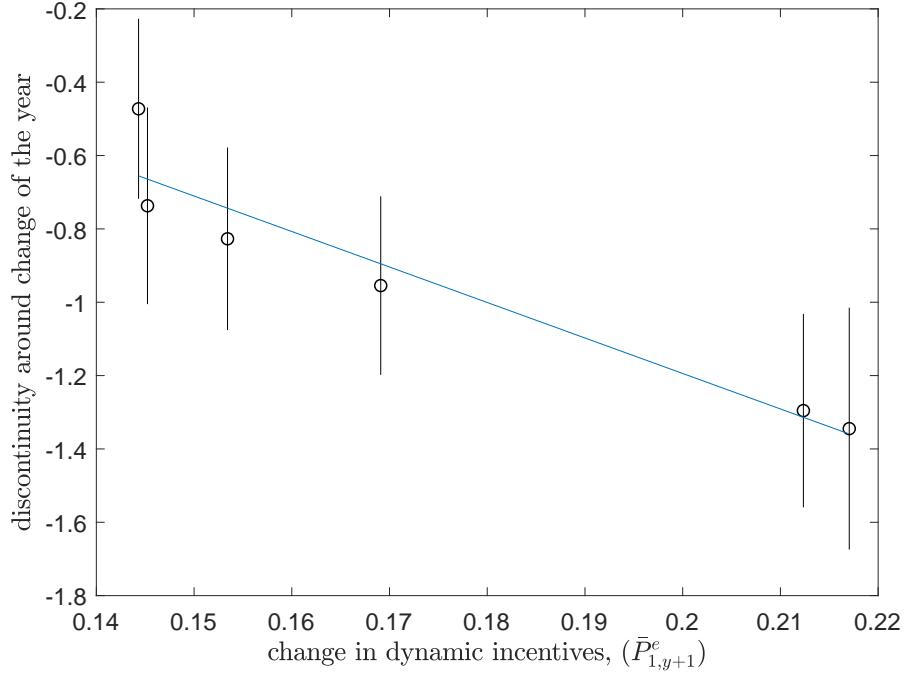
D.4 Log specification

In this specification, we use the logged expenditure (plus 1) dependent variable for healthcare utilization in place of pseudo-censored expenditures. These results are reported in Figure A.10 and in the third row of Table 5. Our approach is robust to such alternative specifications.

D.5 Removing last year-pair

In footnote 8, we mentioned that the data were collected for a project in which a new payment model for GPs was evaluated. This pilot began in July 2014 and could thus affect our estimated discontinuity size for the year-pair $\{2014, 2015\}$. To ensure that this one year-pair is not the sole driver of our findings, we apply our empirical approach on all year-pairs other than the year-pair $\{2014, 2015\}$. The results are shown in Figure A.11 and in the last row of Table 5. We see that our conclusions barely change.

Figure A.11: Removing last year-pair



Notes: These figures plot the relationship between the estimated changes due to resetting contracts and the computed $\bar{P}_{1,y+1}^e$ for the baseline sample as defined in Section 3.2. We omit the estimated discontinuity size from the last year-pair. The solid line is the OLS regression line.

Table A.5: Dependence of discontinuity sizes on dynamic incentives

specification	γ^e	MR test	Up test	Down test
log dependent variable	-0.414 (0.0737)	0.658	0.000	0.688
PC at 5000	-13.072 (4.7448)	0.032	0.264	0.958
last year-pair removed	-9.698 (1.5325)	0.000	0.009	0.998

Notes: The table reports our estimates for γ^e and our results from non-parametric tests of monotonicity across different specifications. See Section 6.1 to 6.4 for details. The effect of dynamic incentives, γ^e , was estimated by regressing the estimated discontinuities on $\bar{P}_{1,y+1}^e$. The MR, Up and Down tests were conducted using 10,000 bootstrap repetitions. See Appendix C for details.