

Supplementary Material 1

Description of information available in GePaRD

For each person, the GePaRD database contains demographic information as well as information on hospitalizations, outpatient physician visits, and drug dispensations. The hospital data comprise information on the dates of hospitalization, diagnoses, reasons for admission and discharge, and diagnostic and therapeutic procedures. Claims of outpatient physician visits include outpatient treatments, procedures, and diagnoses. Reimbursed drugs are identified based on the Anatomical Therapeutic Chemical (ATC) codes, diagnoses are identified based on the International Classification of Diseases, tenth revision, German modification (ICD-10-GM) codes and procedures and services based on Operation and Procedure classification (OPS) codes and Uniform Assessment Standard (EBM) codes. Dispensation data are available for all reimbursed outpatient dispensations and include the dates of prescription and dispensation, the amount of substance prescribed, and information on the prescribing physician.

The following table lists all exposures, covariates as well as the outcome, the ICD, ATC, OPS or EBM codes used for defining the variable (if applicable) and gives a brief description of the derivation of the variable.

Outcome definition: For the definition of pancreatic cancer we considered inpatient diagnoses (ICD10 GM C25), which are considered to have a high validity. Patients with no inpatient but outpatient diagnosis codes of pancreatic cancer were only classified as pancreatic cancer cases if additional criteria were fulfilled such as coding of diagnostic examinations and death within 6-9 months after diagnosis to avoid misclassification. This algorithm was developed based on case reviewing and subject knowledge.

Exposure	ATC-Code	Description
Metformin	A10BA02, A10BD03, A10BD05, A10BD07, A10BD08, A10BD10, A10BD15, A10BD16, A10BD20, A10BD31	Set to 1 if proportion of days covered (PDC) in respective interval was >0.5 or if PDC was >0.25 (but ≤0.5) for a certain year and in addition >0.5 for the preceding and succeeding year
DPP4-inhibitors	A10BH01-A10BH07, A10BH51, A10BD07, A10BD08, A10BD10	Set to 1 if PDC in respective interval was >0.5

		or if PDC was >0.25 (but ≤0.5) for a certain year and in addition >0.5 for the preceding and succeeding year
Other oral glucose-lowering medication	A10BB01-A10BB12, A10BB31, A10BF01-A10BF03, A10BG01-A10BG03, A10BJ01-A10BJ06, A10BK01-A10BK05, A10BD03-A10BD06, A10BD15, A10BD16, A10BD20, A10BD31	Set to 1 if PDC in respective interval was >0.5
Insulins	A10A	Set to 1 if PDC in respective interval was >0.5
Outcome	ICD-10-GM-Code	Description
Pancreatic cancer	C25	First in- or outpatient diagnosis
Death	Not applicable	Death is documented as a potential reason for leaving the insurance company. For persons dying in hospital, death is additionally documented in the hospital data.
Covariates	Code derivation	Description
Duration of diabetes	ICD: E11, E14 Glucose tests: ATC: V04CA01- V04CA08 EBM: 32025, 32881	Time in days from first diabetes type 2 diagnosis to first dispensation of metformin. Date of diabetes diagnosis was set to date of glucose test if a test was conducted within 91 days preceding the first diabetes diagnosis
Socioeconomic status		0: no information or no graduation 1: secondary school certificate 2: higher education entrance qualification
Indicator variable for smoking/alcohol/drug abuse	Alcohol abuse: ICD: E24.4, F10.0-F10.9, G31.2, G62.1, G72.1, I42.6, K29.2,	One of the following had to be fulfilled: <ul style="list-style-type: none">at least one in- or

	<p>K70.0-K70.4, K70.9, K85.2, K86.0, T51.9, Z50,2</p> <p>ATC: N05CM02, N07BB</p> <p>Smoking:</p> <p>ICD: F17.0-F17.9, T65.2</p> <p>ATC: N07BA</p> <p>Drug abuse:</p> <p>ICD: F11.0-F11.9, F12.0-F12.9, F13.0-F13.9, F14.0-F14.9, F15.0, F16.0-F16.9, F18.0-F18.9, F19.0-F19.9, Z50.3, Z72.0</p>	<p>outpatient diagnosis or dispensation</p> <ul style="list-style-type: none"> at least one in- or outpatient diagnosis or dispensation at least one in- or outpatient diagnosis <p>Maximum of (smoking, alcohol, drug abuse)</p>
Treated hypertension	<p>ICD: I10-I13, I15</p> <p>ATC: C02, C03, C07-C09</p>	<p>At least one in- or outpatient diagnosis <u>and</u> dispensations for hypertension with at least 180 DDDs</p>
Coronary heart disease	<p>ICD: I20-I25</p> <p>OPS: 5360-5363, 12755, 88370-88372, 88375, 88376, 88378, 88399, 8837e, 8837k, 8837m, 8837p, 8837q, 8837t, 8837u, 8837v, 8837w, 8839a, 883d</p> <p>ATC: B01AC, C01DA, C01DX11, C01DX12, C07, C08, C09XA53, C09XA54, C09A-C09D, C10AA-C10AD, C10AX01- C10AX03, C10AX05- C10AX09, C10AX11-C10AX16, C10AX19, C10AX21</p>	<p>One of the following had to be fulfilled:</p> <ul style="list-style-type: none"> at least one inpatient diagnosis at least one inpatient procedure during hospital stay with non-missing main discharge diagnosis at least one outpatient diagnosis in three different quarters <u>and</u> dispensations for treatment of at least 50 DDDs
Congestive heart failure	<p>ICD: I50, I11.0, I13.0, I13.2</p> <p>ATC: C01A, C03, C07, C09A-C09D</p>	<p>One of the following had to be fulfilled:</p> <ul style="list-style-type: none"> at least one inpatient diagnosis at least one outpatient diagnosis in three different quarters <u>and</u> dispensations for treatment of at least 50 DDDs

Lipid lowering drugs	ATC: C10AA	Dispensations with at least 180 DDDs
Chronic obstructive pulmonary disease	ICD: J43.2, J43.8, J43.9, J44 ATC: R03A, R03BB, R03C, R03DX07	One of the following had to be fulfilled: <ul style="list-style-type: none"> at least one inpatient diagnosis at least one outpatient diagnosis and one dispensation at least one day participation in COPD disease management program
Asthma	ICD: J45, J46 ATC: D11AH05, H02AB, R03A, R03BA, R03BC, R03C, R03DA, R03DB, R03DC, R03DX05, R03DX08-R03DX10	One of the following had to be fulfilled: <ul style="list-style-type: none"> at least one dispensation and at least one in- or outpatient diagnosis at least one day participation in asthma disease management program
Dementia	ICD: F00-F03, F05.1, G30	One of the following had to be fulfilled: <ul style="list-style-type: none"> at least one inpatient diagnosis at least one outpatient diagnosis from a neurologist at least two outpatient diagnoses in two consecutive quarters
Hemiplegia	ICD: G04.1, G11.4, G80-G83	At least one in- or outpatient diagnosis
Antidepressants use	ATC: N06AA, N06AB, N06AF, N06AG	Dispensations with at least 180 DDDs
Antipsychotics use	ATC: N05A	Dispensations with at least 180 DDDs
Comorbidity score		Sum of dichotomized

		variables for treated hypertension, coronary heart disease, congestive heart failure, lipid modifying agents, chronic obstructive pulmonary disease, asthma, dementia, hemiplegia, antidepressants use, antipsychotics use
Microvascular complications including retinopathy, nephropathy and neuropathy	ICD: E11.2-E11.4, E14.2-E14.4, G63.2, N08.3, H36.0	At least one in- or outpatient diagnosis
Poor glycaemic control	ICD: E11.0, E11.1, E14.0, E14.1, R73	At least one in- or outpatient diagnosis
Liver disease	ICD: K70-K76, B18.8	One of the following had to be fulfilled: <ul style="list-style-type: none"> at least one inpatient diagnosis at least one outpatient diagnosis in two different quarters)
Severe liver disease	ICD: K70-K77, B18.8	At least one inpatient diagnosis
Chronic kidney disease	ICD: N18, I12, I13, Z49.2 OPS: 8857, 88570, 88571, 885710, 885711, 885712, 885713, 885714, 885715, 885716, 885717, 885718, 885719, 88571a, 88571b, 88571c, 88572, 885720, 885721, 885722, 885723, 885724, 885725, 885726, 885727, 885728, 885729, 88572a, 88572b, 88572c, 8857x, 8857y, 8853, 88530, 88531, 885310, 885311, 885312, 885313, 885314, 885315, 885316, 885317, 885318, 885319, 88531a, 88531b, 88531c, 88531d, 88531e, 88531f, 88532, 885320, 885321, 885322, 885323, 885324, 885325,	One of the following had to be fulfilled: <ul style="list-style-type: none"> at least one in-or outpatient diagnosis one in- or outpatient operation one outpatient treatment

	885326, 885327, 885328, 88533, 88534, 88535, 88536, 88537, 885370, 885371, 885372, 885373, 885374, 885375, 885376, 885377, 885378, 885379, 88537a, 88537b, 88537c, 88538, 885380, 885381, 885382, 885383, 885384, 885385, 885386, 885387, 885388, 885389, 88538a, 88538b, 88538c, 8853x, 8853y, 8854, 88540, 88541, 885410, 885411, 885412, 885413, 885414, 885415, 885416, 885417, 885418, 88542, 88543, 88544, 88545, 88546, 885460, 885461, 885462, 885463, 885464, 885465, 885466, 885467, 885468, 885469, 88546a, 88546b, 88546c, 88547, 885470, 885471, 885472, 885473, 885474, 885475, 885476, 885477, 885478, 885479, 88547a, 88547b, 88547c, 88548, 8854x, 8854y, 8855, 88550, 88551, 885510, 885511, 885512, 885513, 885514, 885515, 885516, 885517, 885518, 885519, 88551a, 88551b, 88551c, 88551d, 88551e, 88551f, 88552, 885520, 885521, 885522, 885523, 885524, 885525, 885526, 885527, 885528, 88553, 88554, 88555, 88556, 88557, 885570, 885571, 885572, 885573, 885574, 885575, 885576, 885577, 885578, 885579, 88557a, 88557b, 88557c, 88558, 885580, 885581, 885582, 885583, 885584,	
--	---	--

	885585, 885586, 885587, 885588, 885589, 88558a, 88558b, 88558c, 8855x, 8855y, 8856 EBM: 13611, 40823, 40824, 40825, 40826, 40827, 40828, 40837, 40838	
End-stage renal disease	OPS: 8853, 8854, 8855, 8857, 885a EBM: 13600, 13602, 13610, 13611	One of the following had to be fulfilled: <ul style="list-style-type: none"> • at least one in- or outpatient procedure • at least one outpatient treatment
Liver disease/ severe liver disease / chronic kidney disease / terminal renal disease		Maximum of (liver disease, severe liver disease, chronic kidney disease, terminal renal disease)
Myocardial infarction or stroke	ICD: I21, I22, I60-I64	At least one inpatient diagnosis
Cancer	ICD-10: C00-C97 (except C25 and C44)	At least one inpatient diagnosis
Number of hospitalisations		Number of distinct hospitalisation dates
Number of visits with diabetologist		Number of distinct diabetologist visit dates
Obesity	ICD: E66.0, E66.2, E66.8, E66.9 OPS: 54343, 54344, 54454, 54455, 5448a-f	One of the following had to be fulfilled: <ul style="list-style-type: none"> • at least one in- or outpatient diagnosis • at least one in- or outpatient procedure

Supplementary Material 2

Variable	Purpose	Time period of variable assessment	Type of model when used as dependent variable	Functional form when used as predictor	Modelling type selected in SAS macro ^c for covariate history (covXptype)
Baseline					
Age at cohort entry	Confounder	Pre-baseline	Not predicted	Continuous	
Sex	Confounder	Pre-baseline	Not predicted	Binary	
Duration of diabetes	Confounder; proxy	Pre-baseline	Not predicted	Continuous	
Socioeconomic status	Confounder, proxy	Pre-baseline	Not predicted	Categorical (3 categories)	
Indicator for smoking/alcohol/drug abuse ^a	Confounder	Pre-baseline	Not predicted	Binary	
Comorbidity score	Indicator for general health to reflect ability to adhere to the treatment strategy	Pre-baseline	Not predicted	Continuous	
Outcome					
Pancreatic cancer	Outcome variable	Post-baseline	Pooled across time intervals logistic regression (discrete time hazard)	N/A	N/A
Treatment					
Metformin dispensation	Exposure	Post-baseline	Logistic	Binary	lag1cumavg ^d
Dispensation of DPP-4-inhibitors	Exposure	Post-baseline	Logistic	Binary	lag1cumavg ^d
Time-varying factors					
Other glucose-lowering treatment	Confounder	Post-baseline	Logistic	Binary	lag1bin ^e
Microvascular complications including retinopathy, nephropathy and neuropathy	Factor to intensify treatment; switch to MET+DPP-4i strategy	Post-baseline	Logistic to failure	History-binary ^b	tsswitch1 ^f
Poor glycaemic control	Factor to switch to MET+DPP-4i strategy	Post-baseline	Logistic	Binary	lag1bin ^e
Liver disease/ severe liver disease / chronic kidney disease / terminal renal disease	Contraindications; switch or quit medication	Post-baseline	Logistic to failure	History-binary ^b	tsswitch1 ^f
Myocardial infarction or stroke	Contraindications; switch or quit medication	Post-baseline	Logistic to failure	History-binary ^b	tsswitch1 ^f
Cancer	Confounder	Post-baseline	Logistic to failure	History-binary ^b	tsswitch1 ^f
Number of hospitalisations	Confounder	Post-baseline	Logistic and linear	Continuous	lag1bin ^e
Number of visits with diabetologist	Confounder in sensitivity analysis	Post-baseline	Logistic and linear	Continuous	lag1bin ^e
Obesity ^a	Confounder	Post-baseline	Logistic to failure	History-binary ^b	tsswitch1 ^f
Other important factors that could not be considered	Reason not accounting for the variable				

HbA1c / symptoms of an undetected cancer	Confounder	Unobserved
Chronic pancreatitis	Confounder, contraindication ddp4 inhibitors	Low prevalence; models did not converge

eTable 1: Summary of covariates, type of model and functional form when used as covariate in parametric g-formula

^aThe database contains information on smoking, drug use and obesity diagnosis but no information on changes to these statuses

^bThe term history-binary denotes variables which remain at 1 once they have switched from 0 to 1

^cDetails are given in the user guide of the GFORMULA SAS macro available at <https://causalab.sph.harvard.edu/software/>

^d At time t , the cumulative average of the history of the covariate relative to interval t beginning from time 0 to time $t-2$ is computed. The last term is not included in the average. There are two generated predictors for the covariate at time t , i.e. the lagged covariate at time $t-1$ and the average of the covariate from time 0 to time $t-2$.

^e A lagged variant of the covariate (Cov_{t-1}) is created and used as predictor.

^f The history of the covariate is modelled as a function of the time since the covariate last switched from 0 to 1 at each time t .

Supplementary Material 3

Time-related and other sources of bias

In this section, we will summarize the main sources of time-related and other biases as relevant to the analysis of healthcare claims databases. Some of these are avoidable by suitable design and methods. The section follows Suissa and Dell'Aniello (2020)¹, Suissa and Azoulay (2012)² and Hernán et al. (2016)³ providing a comprehensive overview of time-related biases in pharmacoepidemiology or observational studies.

Immortal time bias

Immortal time bias is a common source of bias in observational studies²⁻⁵ and results from misclassifying unexposed time as exposed time. It occurs when study subjects using specific drugs are compared to those not using drugs but the time interval between time zero and the first drug prescription is counted to the exposure time for the drug-users. The time prior to treatment initiation is actually yet unexposed and called immortal because by definition no outcome event can occur (in the user-group) so that the risk is guaranteed to be zero. Immortal time bias typically results in an overestimation of the outcome of interest in the non-users and hence in an apparent protective effect of the drug.⁶

In our study, we applied a new user design comparing two drugs initiated in persons with similar duration of T2DM. Treatment assignment and eligibility were thus aligned (Hernan et al, 2016). In the analysis stage, both groups were analytically “forced” to adhere to their treatment strategies using the parametric g-formula so that immortal time bias was avoided.

Time-lag bias

Time-lag bias occurs when comparing treatments that are typically prescribed in different disease stages, i.e. first-line treatments are compared with second- or third-line treatments.² This can result in confounding by disease duration (and stage) because an outcome (e.g. cancer incidence) related to the first-line treatment may also be attributed to the second-line treatment if it occurs after a long period of exposure. For instance, metformin is mainly a first-line treatment for type-2 diabetes whereas DPP-4i are usually prescribed as second-line treatment after using metformin for a certain time. As T2DM itself is a risk factor for cancer, an increased cancer risk may be observed for DPP-4i compared to metformin only due to the later disease stage/longer duration of patients receiving DPP-4i.⁷

In our application we enrolled only new drug users showing approximately the same duration of T2DM (mean diabetes duration of 46 days at baseline). We additionally adjusted for diabetes duration. Under strategy B, *everyone* starts DPP-4i *one year after* first metformin use, and the g-formula ensures that the corresponding counterfactual risk is computed correctly (under the stated assumptions).

Latency time bias

According to the Dictionary of Epidemiology – there are two definitions of latency time commonly used in health and life sciences. The first one considers latency time as “the interval from initiation of the disease to clinical emergence or detection of disease”, the second one as “the interval between exposure to the causal agent and appearance or detection of the health process”. In this study we use the first definition of latency time because we mainly use it to get an indication on unmeasured confounding by an undetected disease (see also protopathic bias).

Latency time is a common issue in cancer outcomes because cancers are typically assumed to develop a long time after initial exposure to a certain causal agent.¹ In such cases a sufficiently long follow-up time needs to be considered to avoid a risk underestimation, although it may be difficult in practice to know what ‘sufficiently long’ means when little is known about the latency time. Moreover, estimating the risk as a function over time can help to identify the time point where the risk changes. Our analysis assessed the risk of pancreatic cancer in persons with T2DM over a seven year follow-up period. Latency time bias is particularly important to consider when the absolute risk of new cancers is of interest.

When comparing two treatment strategies, as in our example, both treatment groups would be affected by the latency time in similar manner if treatments were randomised. Thus, in an observational study latency time will possibly bias the overall results if latent pancreatic cancer is related to unobserved confounding factors as described in next.

Protopathic bias

Protopathic bias, also referred to as reverse causality, applies to situations where a medication is initiated in response to an early symptom of the health outcome of interest¹. As a result, exposure to the drug is prior to the time point of the disease diagnosis although the exposure actually occurred after the first manifestation of the outcome. If the early symptom leading to the treatment decision has not been recorded, this is basically a form of unmeasured confounding by an undiagnosed (i.e. latent) disease,⁸ e.g. an occult cancer. The issue is connected with the above latency

time as there needs to be a period of time between onset and diagnosis for protopathic bias to be possible. While the above types of bias (immortal time, time-lag, and possibly latency time) are avoidable, protopathic bias as a form of unmeasured confounding is difficult to deal with. Sensitivity analyses can be carried based on reasonable assumptions (that could be included in the models) about the latency time, or by suitable restrictions.^{8,9}

In our example poor glycaemic control may be an early symptom of (undetected) pancreatic cancer. If then glucose-lowering drugs are prescribed the resulting association will be due to protopathic bias. When estimating the risk function over time this may become apparent as a difference occurring early in time after exposure; this has to be kept in mind when interpreting the results. In our analysis we adjusted for diagnosed poor glycaemic control, but the actual HbA1c levels were not available resulting in possible residual confounding. We conducted a number of sensitivity analyses to check for protopathic bias i) assuming a minimum latency time of one year for the development of pancreatic cancer or excluding persons with T2DM that switched/intensified their medication ii) in the year or iii) 3 to 6 months prior to the pancreatic cancer diagnosis. This exclusion can be seen as an attempt to restricting the study population to those who do not have a latent disease at time-zero mimicking a hypothetical target trial, where participants are screened for a latent disease before treatment assignment.⁸

Prevalent user bias

This bias results from including prevalent users of a drug in study design and analysis.^{10,11} Since prevalent users are survivors of the initial period of drug use, this can e.g. lead to bias if the risk of the outcome changes over time. In addition, users experiencing an outcome early (before start of the study) may stop taking the drug giving an extra survival benefit to the prevalent (sustained) users.

This bias was circumvented by a careful target trial emulation where time zero, eligibility and treatment assignment were aligned; as treatment cannot be assigned retrospectively, this avoids prevalent user bias.³

Immeasurable time bias

This bias occurs if healthcare databases used to study association between drug use and certain health outcomes, especially death, include only prescriptions of medications on an outpatient basis, but not during hospitalizations (inpatient medications).¹ This means patients being hospitalized during the exposure time window will have missing exposure information resulting in bias.

GePaRD includes information on inpatient and outpatient diagnosis, but only on outpatient dispensations. However, it is assumed that patients that are hospitalized due to reasons other than diabetic complications will continue their therapy using their own medications. In addition, we adjusted for the number of hospitalizations.

Confounding by indication

This bias occurs when the risk of an adverse event is related to the indication for drug use but not the use of the drug itself. This means that the clinical condition that determined the prescription of the drug acts as a confounding factor. This situation is also referred to as confounding by indication.

In our example, different drug classes are recommended to control for poor glycaemic control and prevent T2DM complications. T2DM/ poor glycaemic control itself is associated with multiple risk factors making T2DM patients high-risk individuals for e.g. hypertension, CVD and cancer. Thus, poor glycaemic control may confound the association between certain glucose-lowering drugs and cancer.¹²⁻¹⁴ In our example we adjusted for diagnosed poor glycaemic control (being included as time-varying confounder) to account for this source of bias.

Detection bias (also surveillance bias)

This bias occurs when the probability of detecting the study outcome is higher in one exposure group due to increased surveillance, screening or testing of the outcome itself, or an associated symptom. Poor glycaemic control (increased HbA1c) may lead to an addition or switch of medication (e.g. from metformin to combination therapy with DPP-4i) and may further go along with more regular screenings such that detection bias may arise in the exposure group with the metformin/DPP-4i combination therapy. Further, any differential detection may simply reveal cancers that were already latent at the start of DPP-4i. With a long follow-up it is likely that the latent cases will be detected eventually in both groups so that if there is no actual treatment effect, initial differences will level out over time. Moreover, we adjusted for the number of hospitalizations and in a sensitivity analysis additionally for the number of visits at a diabetologist as proxies for increased surveillance.

Competing events

The presence of a competing event, such as death, needs to be taken into account by the choice of causal contrast and analysis to allow for the desired interpretation. Sometimes the term 'competing event bias' is used when analysis and interpretation do not match, even though this is not strictly speaking a bias issue. Young et al.

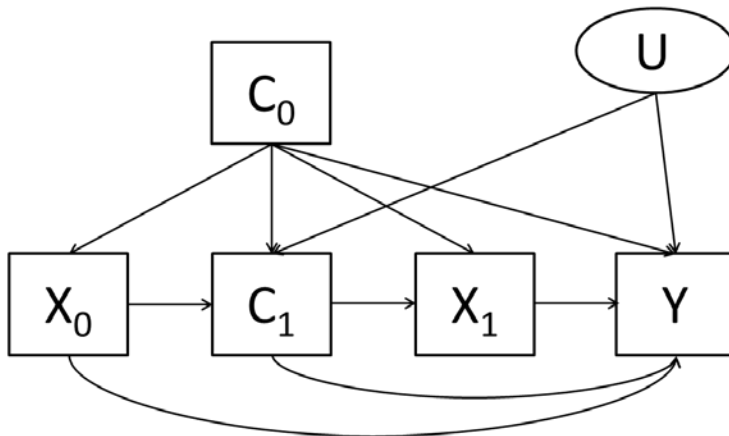
(2020)¹⁵ describe the causal interpretation of different popular statistical approaches for dealing with competing events. In brief, the first option is to target the total causal effect of a treatment on the event of interest by estimating the event-specific cumulative incidence as we do in our analysis. However, this may imply that if treatment, say, increases the risk of death a reduced total risk of the event of interest (pancreas cancer) may result because participants die before the latter is diagnosed. In our application, Figure 1 suggests that this is not the case as the incidences of death are nearly identical for strategies A and B.

The second option is to censor the competing event, which corresponds to a direct effect of treatment on the event of interest under a hypothetical intervention where the competing event can be eliminated. This is typically not a meaningful approach when the competing event is death as in our application. Moreover, for this second option, one needs to additionally adjust for common causes of the two types of events (pancreas cancer diagnosis and death) in order to justify the 'independent censoring' assumption.

Supplementary Material 4

Time-dependent confounding

Time-dependent confounding occurs when there is a time-varying variable affecting both, the disease and the time-varying exposure. It is a common situation that a time-dependent confounder is also *affected by the prior exposure*, *i.e.* the subsequent values of the time-dependent confounder are affected by prior exposure. For instance, obesity (time-dependent confounder) may be causally related to diabetes medications (exposure) as well as cancer risk (outcome) and earlier diabetes medications may affect subsequent obesity risk resulting in feedback as displayed in Figure 1. Another example of a time-varying confounder is poor glycaemic control in our example.



X_0, X_1 : Exposure at $t=0$ and $t=1$ (e.g. metformin mono-/combination therapy)
 C_0 : Confounders at baseline
 C_1 : Time-varying confounders at $t=1$ (e.g. obesity)
 Y : Outcome variable (e.g. pancreatic cancer)
 U : Unobserved factor

eFigure 1: Time-varying confounding affected by prior exposure

Whereas standard regression models do not correctly adjust for time-varying confounding, the g-methods suggested by Robins and Hernán^{16,17} provide valid ways of adjustment (see Daniel, et al.¹⁸ for an accessible introduction). Time-dependent confounding cannot be accounted for at the design stage, unless treatment / exposure can be sequentially randomised. Hence it typically needs to be considered at the analysis stage. Here we focussed on the parametric g-formula.

Supplementary Material 5

The g-formula is valid under the following assumptions:

Conditional sequential exchangeability

This assumption demands that the potential outcomes under certain fixed exposure levels are independent of the observed exposures. It is made within levels of past observed covariate values (conditional) and at each time point (sequentially). In Figure 1, this is expressed by the absence of edges from unobservables U into exposures X_0 and X_1 while observed covariates C_1 can affect X_1 . Informally, we speak of “no unmeasured confounding”.

Treatment version irrelevance (counterfactual consistency)

It is assumed that the effect of the exposure is the same whether it is set by the considered hypothetical intervention or whether it occurs naturally. This assumption guarantees that the strategies under consideration are well-defined.

Positivity

This assumption demands that for each strategy and within each confounder and treatment history observed in the data through time t , it is possible to also observe in the data a value of treatment at t consistent with the strategy for all t . Positivity is met when there are exposed and unexposed subjects within all confounder and prior exposure levels. This can be verified empirically.

Correct model specification

The g-formula requires correct specification of the conditional (on the past) probabilities of the outcome and time-varying covariates in all follow-up intervals. Due to the use of multiple models, the parametric g-formula is especially vulnerable to the assumption of correct model specification. Informal checking is possible by comparison of the observed data to the data simulated under the natural course. The g-formula should preferably be applied in situations with good knowledge on the causal relationships among the variables of interest.

Moreover, it is important to use flexible models to avoid the g-null paradox.¹⁹

Supplementary Material 6

In the following, the algorithm to apply the parametric g-formula is outlined:

Step 1: Probability modelling

- a) Separate regression models for the treatment and each covariate in year t are fitted as a function of t and past treatment and covariate history, restricted to those who survived and remained uncensored until t . The covariate history may e.g. be summarized by the baseline covariates and the two most recent values of time-varying covariates. The estimated conditional distributions of the covariates are used to construct an estimate of the joint distribution of the covariates.
- b) A discrete-time hazard regression model for the occurrence of a pancreatic cancer diagnosis in year t is fitted as a function of t and past treatment and covariate history at each time t (a logistic regression pooled over time intervals). A logistic model is also fitted to estimate the conditional discrete-time hazard of the competing event (death) at each time t , which is necessary to compute the event-specific cumulative incidence function.

Step 2: Monte Carlo Sampling

In this step, a large number of covariate histories consistent with the intervention is generated. For each treatment strategy, do the following n times (with n being as large as possible to reduce simulation error):

- a) For each year t , the treatment and covariates using the model coefficients estimated in Step 1a are simulated based on previously simulated treatment and covariates through $t-1$ (values at baseline are sampled from the observed data).

- b) The simulated treatment value at time t is replaced with the value of treatment that should be assigned according to the specified treatment strategy.
- c) The discrete-time hazard of pancreatic cancer at t is estimated using the estimated regression coefficients from Step 1b for each of the n simulated histories through t consistent with the specific strategy. The discrete hazard of the competing risk event is also estimated at each t from the estimated model coefficients.

Step 3: Estimation of the risks for the two treatment strategies

- a) For each of the n histories, the estimated hazards from Step 2c are used to compute the n history-specific 7-year risks by the end of follow-up under each strategy
- b) The 7-year risk under treatment strategy A/B in the study population is obtained by averaging the n history-specific risks.

95%- confidence intervals can be obtained by nonparametric bootstrapping, i.e., by repeating these three steps in 100 bootstrap samples.

After concatenating the datasets from step 2, the hazard ratio can be estimated by comparing the hazards in treatment strategy A dataset with those in the treatment strategy B dataset.

Comparison of the distribution of the simulated variables with the ones observed in the actual population can provide an indication on gross model misspecification.

Supplementary Material 7

Implementation of the parametric g-formula

For modelling the outcome, the pancreatic cancer incidence, a logistic regression was fitted by pooling across all time points; this estimated the conditional discrete-time hazard at each time t given past covariates. Also, a logistic model to estimate the conditional discrete-time hazard of the competing event (death) at each time t was fitted.

At each time point $t > 0$ parametric models are fit to estimate the joint distribution of the p covariates given past covariate history through $t-1$, $f_t(\text{cov}_1, \dots, \text{cov}_p \mid \text{past})$. This joint distribution is estimated via a product of conditional distributions:

$f_t(\text{cov}_1 \mid \text{past})f_t(\text{cov}_2 \mid \text{cov}_1, \text{past}) \dots f_t(\text{cov}_p \mid \text{cov}_{p-1}, \dots, \text{cov}_1, \text{past})$ (cmp. Step 1a in Supplementary Material 6).

A separate parametric model is fitted for each conditional distribution of this product. For the covariates $\text{cov}_1, \dots, \text{cov}_p$, an arbitrary ordering was chosen. Under correct model specification, the g-formula would be insensitive to this ordering of the

covariates. As we cannot use saturated models, we conducted a sensitivity analysis changing the arbitrary order of covariates to check for model misspecification (cf. previous examples²⁰⁻²²). The type of model used for each covariate when considered as dependent variable as well as the functional form when used as predictor is summarized in Supplementary Material 2.

In brief, each time-varying predictor was classified as binary, history-binary or continuous. Binary-dependent variables, like poor glycaemic control, and history-binary-dependent variables (indicators that move only from zero to one like myocardial infarction or stroke) were modelled using logistic regression. The models for history-binary variables were limited to those with no history at the beginning of the 365-day-interval. To account for the high number of zero values, continuous covariates like the number of hospitalizations were modelled based on both, a logistic model (using an indicator whether the covariate is >0) to estimate the probability of the covariate being zero and a linear regression model for the natural log of the covariate restricted to records with the covariate being >0 for the estimation of non-zero values.

All models included, as predictors, all baseline covariates, the current and previous value of all binary and continuous covariates as main terms; history-binary variables were modelled as a function of the time since the covariate last switched from 0 to 1 at each time t . No interaction terms were added.

Plausibility of identifying assumptions

Conditional and sequential exchangeability: As a common issue in observational data, unmeasured confounding cannot be fully excluded e.g. due to the lack of HbA1c and more precise information on lifestyle factors. However, proxies were included for all known important confounders except for chronic pancreatitis (prevalence was too low).

Causal consistency: This assumption is fulfilled when the treatment strategies being assessed are well-defined and correspond to the treatment strategies observed in the data, e.g. the outcome for a patient who happens to receive a metformin monotherapy is the same as if he/she had been assigned to the metformin monotherapy in the target trial, which is plausible.

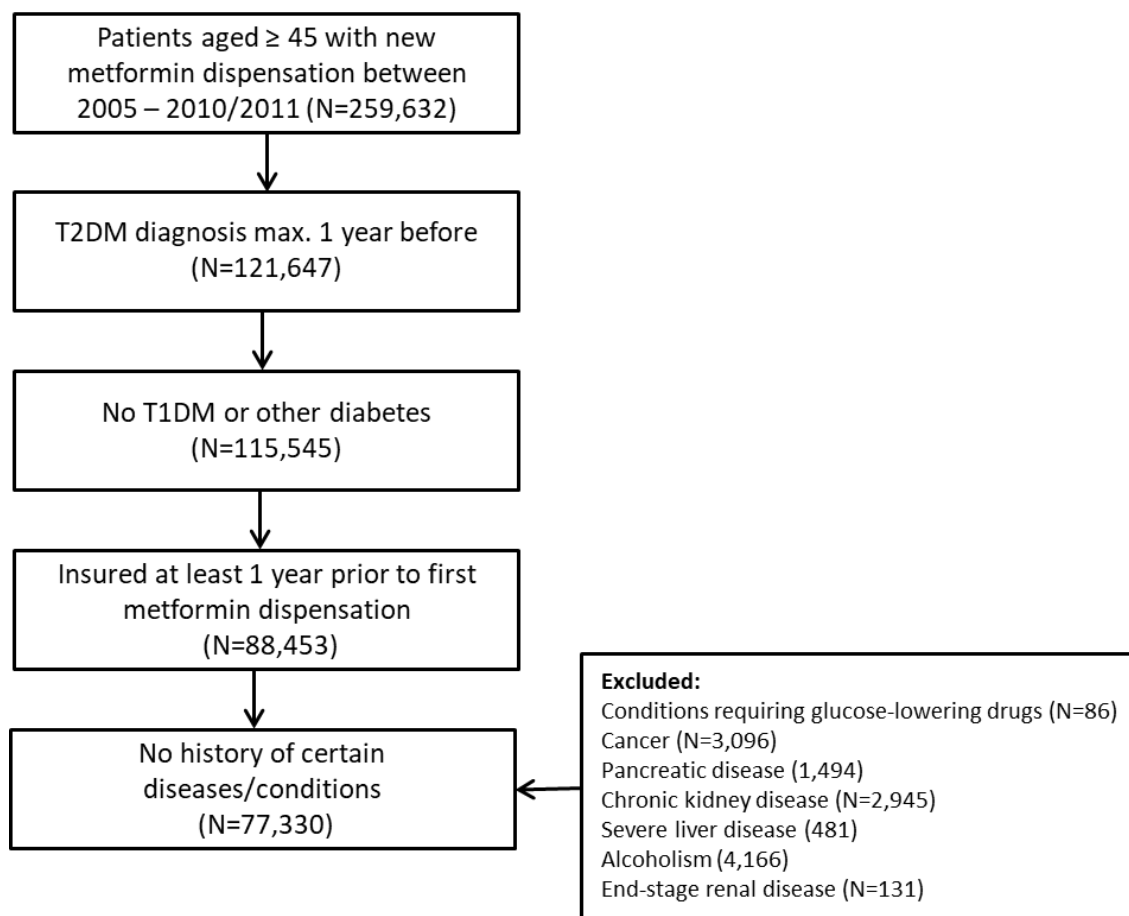
Positivity: The positivity assumption was checked empirically. It was slightly violated with regard to both treatment groups. This is mainly due to the high number of binary/categorical variables (1043 possible combinations) with partly low prevalence such that e.g. combinations like cancer + myocardial infarction/stroke + (severe) liver disease/renal disease/chronic kidney disease are not observed in all strata of the

other covariates. Essentially this means that even this huge dataset does not provide enough data and that the models need to extrapolate information. However, we do not believe that the problem is structural so that it will not strongly affect our results. Nevertheless, instead of adjusting for contraindications, future studies using the parametric g-formula could define dynamic treatment strategies based on contraindications to mitigate this problem. Especially in older patients suffering from multiple diseases a dynamic strategy may better reflect “real-life” treatment practice.

To assess the validity of our parametric assumptions, we compared the observed means of the outcome and time-varying covariates with those predicted by our models. The parametric g-formula closely replicated the observed risk (0.86%) and the mean covariates under the natural course.

Supplementary Material 8

eFigure 2: Flow-chart depicting the selection process leading to the final analysis sample



Supplementary Material 9

Time-varying variables	Baseline/ cohort	Year 1 (N=77,323)	Year 2 (N=74,929)	Year 3 (N=72,758)	Year 4 (N=70,691)	Year 5 (N=68,747)	Year 6 (N=66,772)	Year 7 (N=64,821)
	entry (day 1) (N=77,330)							
Metformin	77,330 (100%)	43,537 (56.3%)	38,042 (50.8%)	38,471 (52.9%)	38,736 (54.8%)	38,535 (56.1%)	37,889 (56.7%)	37,003 (57.1%)
Metformin monotherapy	77,330 (100%)	39,953 (51.7%)	32,222 (43.0%)	30,273 (41.6%)	28,110 (39.8%)	25,703 (37.4%)	23,169 (34.7%)	20,866 (32.2%)
DPP-4-inhibitors	0 (0.0%)	1,487 (1.9%)	3,239 (4.3%)	5,099 (7.0%)	7,237 (10.2%)	9,433 (13.7%)	11,500 (17.2%)	12,954 (20.0%)
Metformin and DPP-4i	0 (0.0%)	1,162 (1.5%)	2,532 (3.4%)	4,049 (5.6%)	5,877 (8.3%)	7,762 (11.3%)	9,489 (14.2%)	10,633 (16.4%)
Metformin and DPP-4i but no other antidiabetic treatments	0 (0.0%)	1,116 (1.4%)	2,361 (3.2%)	3,635 (5.0%)	5,118 (7.2%)	6,501 (9.5%)	7,667 (11.5%)	8,271 (12.8%)
Switchers from monotherapy to combination therapy	0 (0.0%)	1,116 (1.4%)	1,193 (1.6%)	1,260 (1.7%)	1,554 (2.2%)	1,650 (2.4%)	1,673 (2.5%)	1,454 (2.2%)
Other oral antidiabetic treatment	0 (0.0%)	3,346 (4.3%)	5,268 (7.0%)	6,361 (8.7%)	6,988 (9.9%)	7,362 (10.7%)	7,559 (11.3%)	7,948 (12.3%)
Insulin	0 (0.0%)	759 (1.0%)	1,274 (1.7%)	1,914 (2.6%)	2,686 (3.8%)	3,537 (5.1%)	4,538 (6.8%)	5,474 (8.4%)
Pancreatic cancer	0 (0.0%)	214 (0.3%)	102 (0.1%)	84 (0.1%)	61 (0.1%)	74 (0.1%)	51 (0.1%)	66 (0.1%)
Death	0 (0.0%)	802 (1.0%)	828 (1.1%)	951 (1.3%)	1,032 (1.5%)	1,088 (1.6%)	1,156 (1.7%)	1,279 (2.0%)
Censored	7 (0.0%)	1,378 (1.8%)	1,241 (1.7%)	1,032 (1.4%)	851 (1.2%)	813 (1.2%)	744 (1.1%)	63,476 ^a (97.9%)
Chronic Pancreatitis	3 (0.0%)	48 (0.1%)	88 (0.1%)	137 (0.2%)	162 (0.2%)	181 (0.3%)	210 (0.3%)	236 (0.4%)
Microvascular diabetic complications ^b	4,505 (5.8%)	13,049 (16.9%)	17,494 (23.3%)	21,128 (29.0%)	24,452 (34.6%)	27,396 (39.9%)	30,098 (45.1%)	32,249 (49.8%)
Hypoglycaemia	142 (0.2%)	225 (0.3%)	194 (0.3%)	177 (0.2%)	189 (0.3%)	195 (0.3%)	189 (0.3%)	180 (0.3%)
Poor glycaemic control	6,315 (8.2%)	4,183 (5.4%)	3,427 (4.6%)	3,270 (4.5%)	3,019 (4.3%)	2,749 (4.0%)	2,654 (4.0%)	2,465 (3.8%)
Any liver or kidney disease	8,887 (11.5%)	13,044 (16.9%)	14,854 (19.8%)	16,426 (22.6%)	17,825 (25.2%)	19,335 (28.1%)	20,660 (30.9%)	21,917 (33.8%)
Liver disease	8,838 (11.4%)	12,086 (15.6%)	13,129 (17.5%)	13,917 (19.1%)	14,497 (20.5%)	15,039 (21.9%)	15,417 (23.1%)	15,736 (24.3%)
Severe liver disease	0 (0.0%)	91 (0.1%)	175 (0.2%)	240 (0.3%)	296 (0.4%)	372 (0.5%)	436 (0.7%)	500 (0.8%)
Chronic kidney disease	53 (0.1%)	1,156 (1.5%)	2,176 (2.9%)	3,240 (4.5%)	4,360 (6.2%)	5,766 (8.4%)	7,219 (10.8%)	8,676 (13.4%)
Terminal renal disease	0 (0.0%)	71 (0.1%)	131 (0.2%)	216 (0.3%)	317 (0.4%)	465 (0.7%)	650 (1.0%)	886 (1.4%)

Time-varying variables	Baseline/ cohort entry (day 1) (N=77,330)	Year 1 (N=77,323)	Year 2 (N=74,929)	Year 3 (N=72,758)	Year 4 (N=70,691)	Year 5 (N=68,747)	Year 6 (N=66,772)	Year 7 (N=64,821)
Cancer	0 (0.0%)	1,114 (1.4%)	904 (1.2%)	884 (1.2%)	859 (1.2%)	818 (1.2%)	798 (1.2%)	833 (1.3%)
Cardiovascular events	1,621 (2.1%)	2,332 (3.0%)	2,888 (3.9%)	3,431 (4.7%)	3,962 (5.6%)	4,455 (6.5%)	4,939 (7.4%)	5,434 (8.4%)
Myocardial infarction	733 (0.9%)	1,061 (1.4%)	1,320 (1.8%)	1,574 (2.2%)	1,849 (2.6%)	2,085 (3.0%)	2,330 (3.5%)	2,603 (4.0%)
Stroke	892 (1.2%)	1,294 (1.7%)	1,601 (2.1%)	1,902 (2.6%)	2,186 (3.1%)	2,471 (3.6%)	2,748 (4.1%)	3,009 (4.6%)
Alcohol abuse	31 (0.0%)	522 (0.7%)	916 (1.2%)	1,276 (1.8%)	1,577 (2.2%)	1,783 (2.6%)	2,015 (3.0%)	2,235 (3.4%)
Smoking	4,422 (5.7%)	5,911 (7.6%)	6,736 (9.0%)	7,399 (10.2%)	7,849 (11.1%)	8,318 (12.1%)	8,719 (13.1%)	9,012 (13.9%)
Drug abuse	426 (0.6%)	628 (0.8%)	794 (1.1%)	988 (1.4%)	1,160 (1.6%)	1,310 (1.9%)	1,414 (2.1%)	1,522 (2.3%)
Obesity	26,685 (34.5%)	31,443 (40.7%)	33,407 (44.6%)	34,567 (47.5%)	35,396 (50.1%)	35,933 (52.3%)	36,296 (54.4%)	36,306 (56.0%)
Number of hospitalizations, mean (sd)	0.3 (0.68)	0.4 (0.85)	0.3 (0.86)	0.4 (0.91)	0.4 (0.93)	0.4 (0.95)	0.4 (0.97)	0.5 (1.00)
Number of visits at a diabetologist, mean (sd)	2.9 (6.79)	3.3 (7.40)	2.2 (6.10)	1.4 (4.50)	0.7 (2.62)	0.5 (2.33)	0.5 (2.39)	0.5 (2.42)

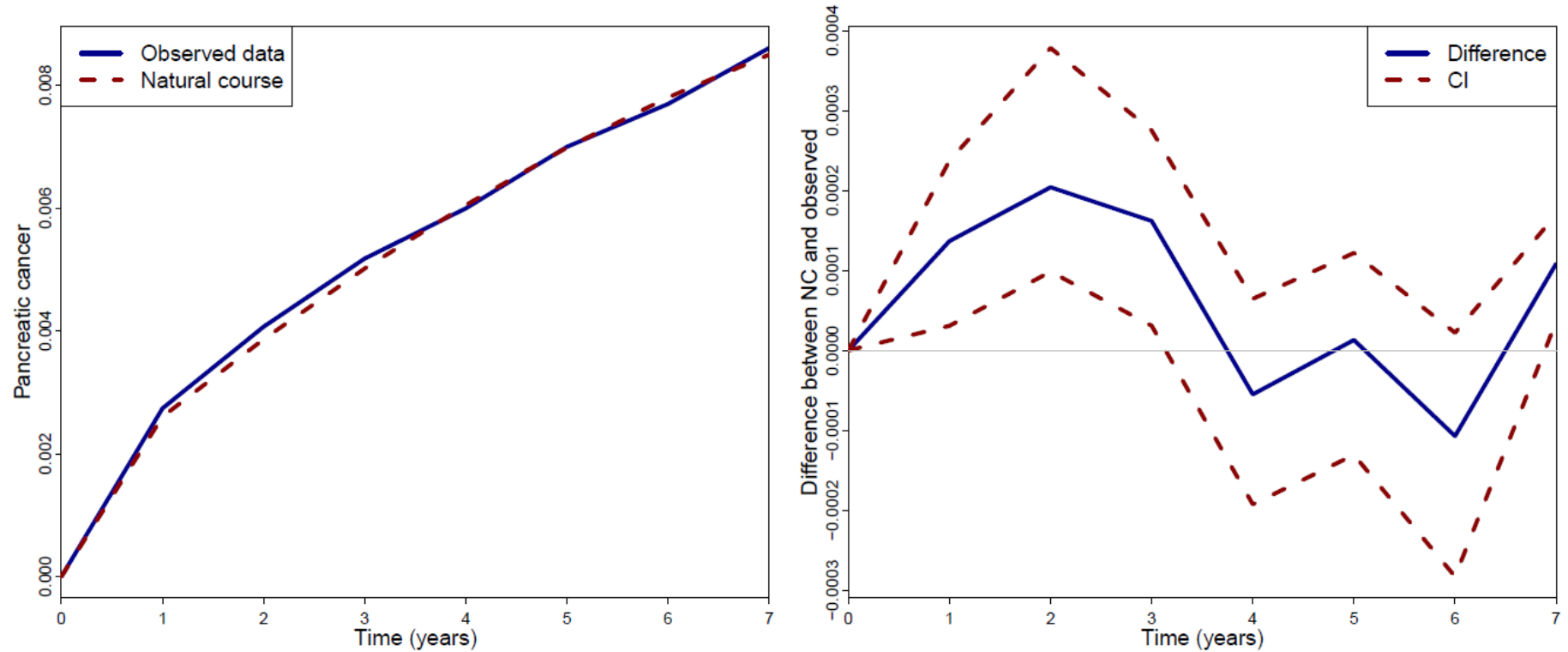
eTable 2: Distribution of exposures and time-varying covariates (number and percentage) over the 7-year follow-up period in the observed data

^aAfter 7 years, all remaining patients are censored (end of follow-up)

^bComplications due to poor glycaemic control such as retinopathy, nephropathy and neuropathy

Supplementary Material 10

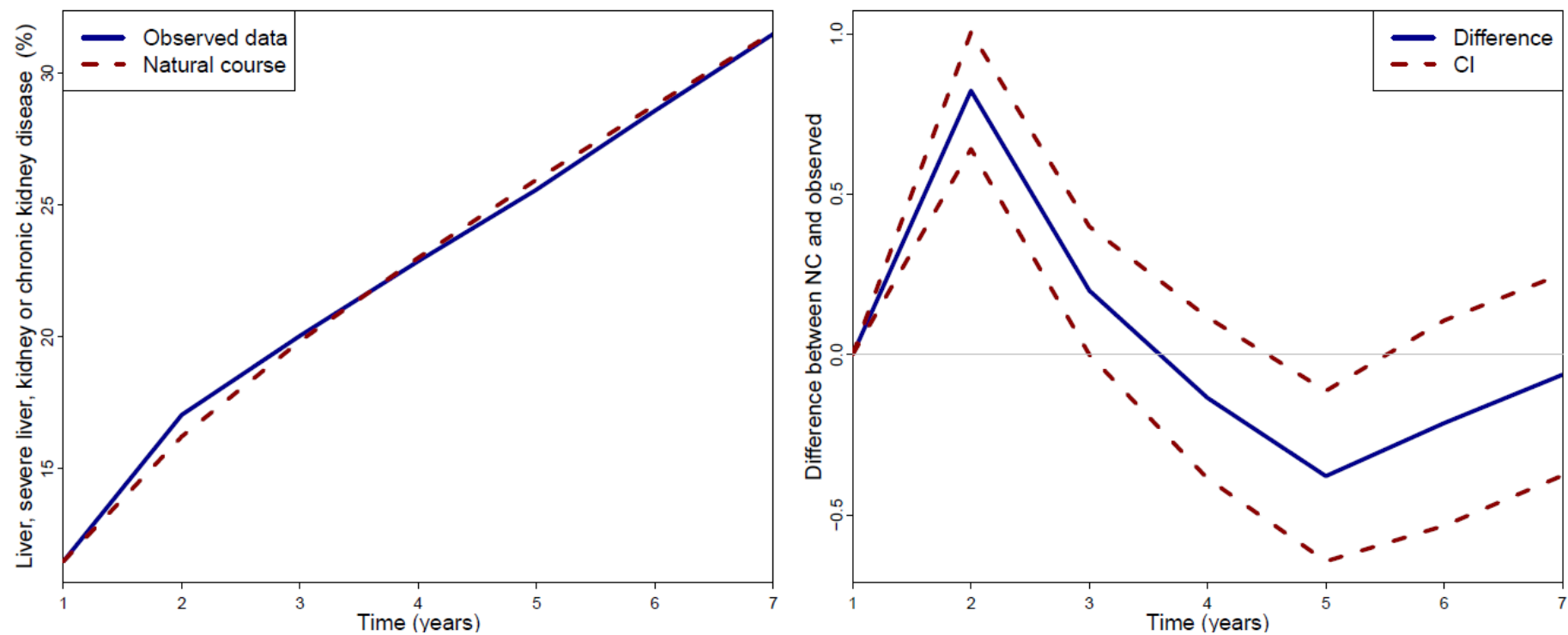
Observed means and predicted means under the natural course for the outcome (pancreatic cancer)

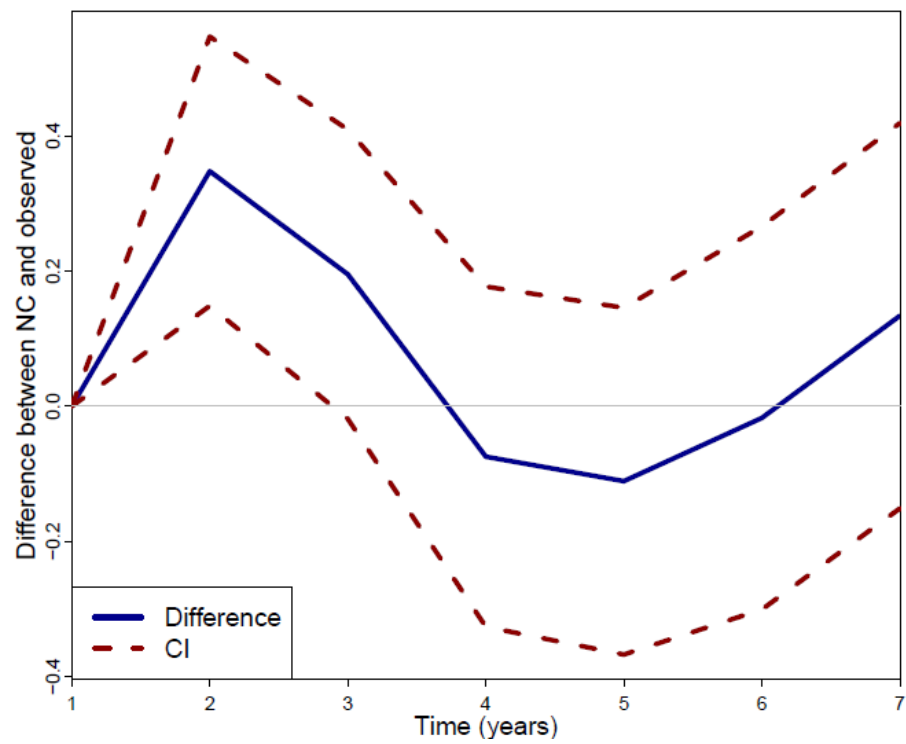
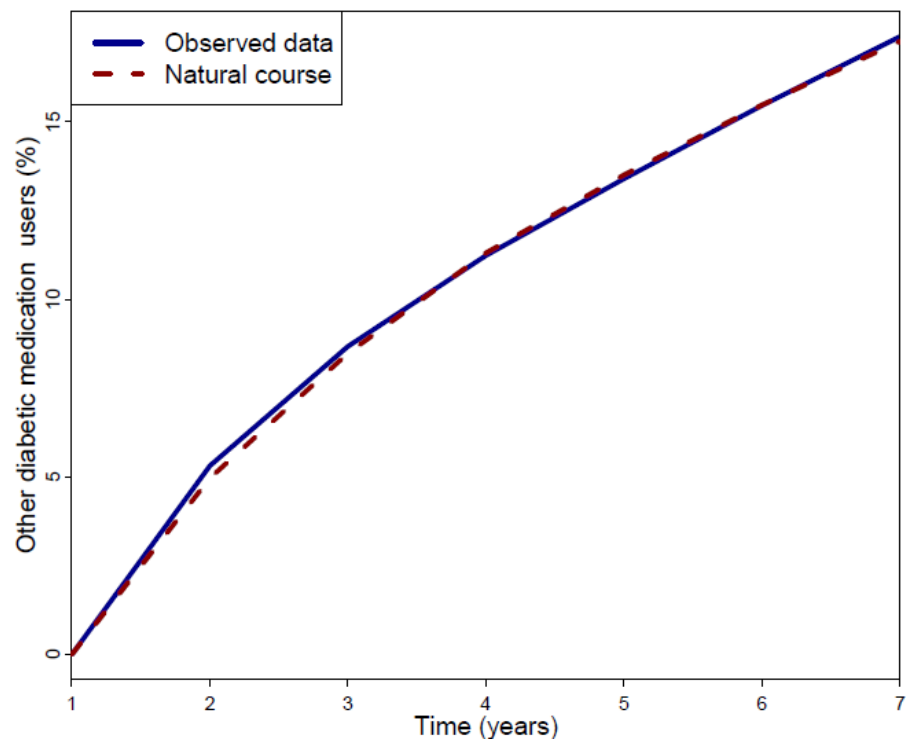


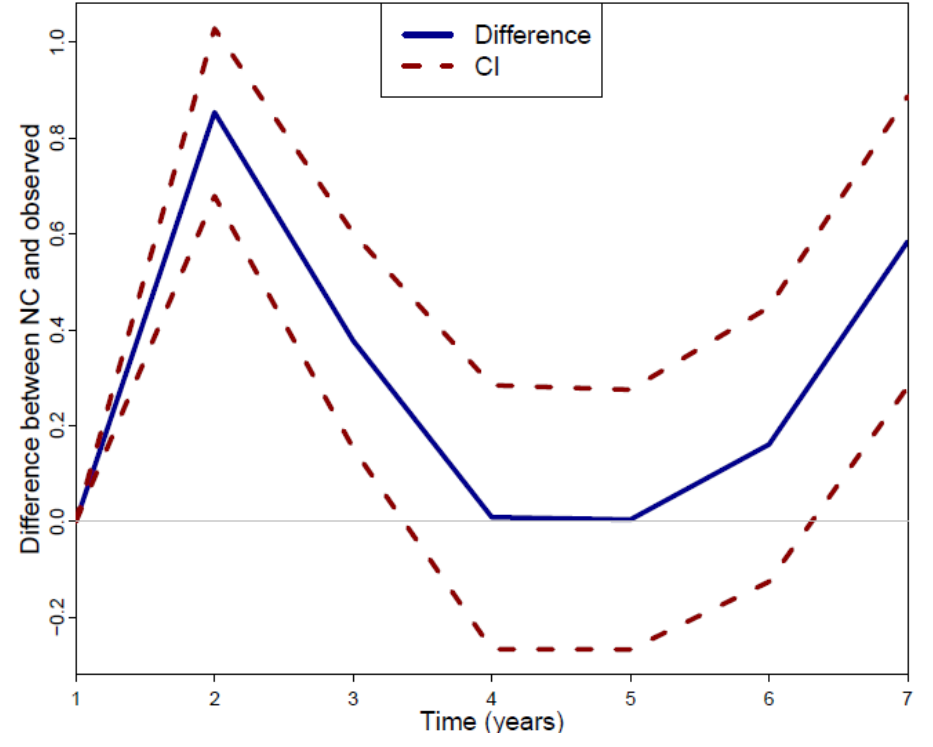
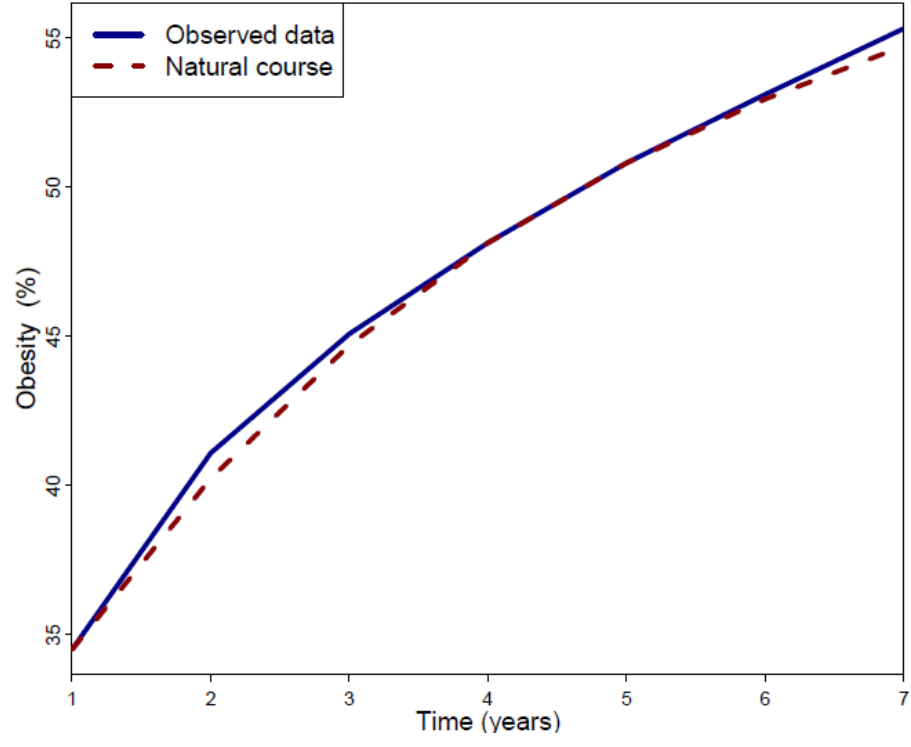
Left column: observed cumulative incidence of pancreatic cancer (solid line), natural course (dotted line) estimates by follow-up year.

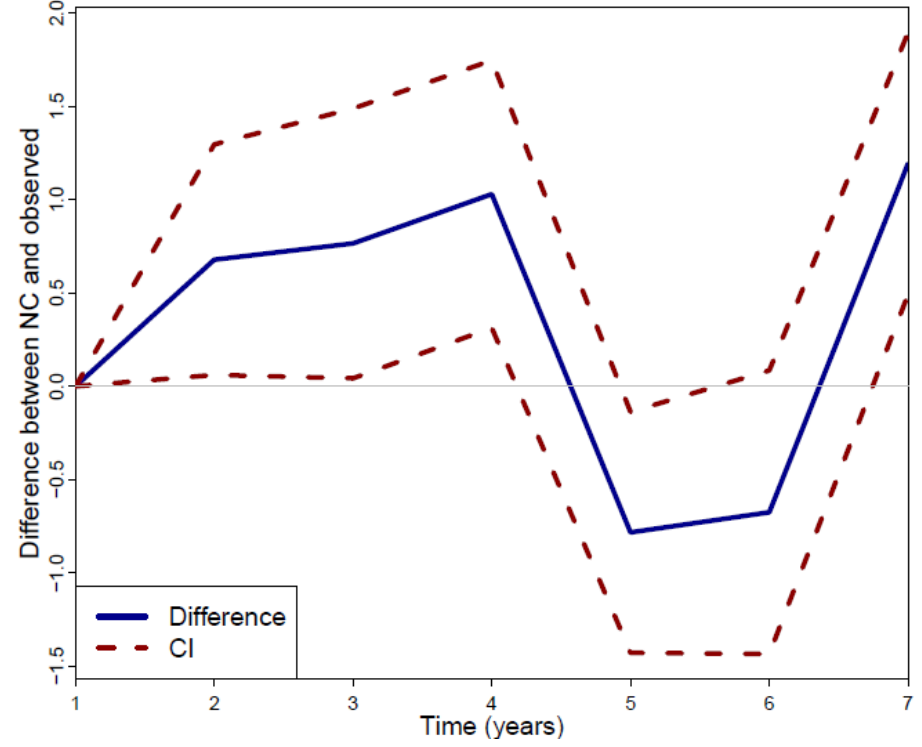
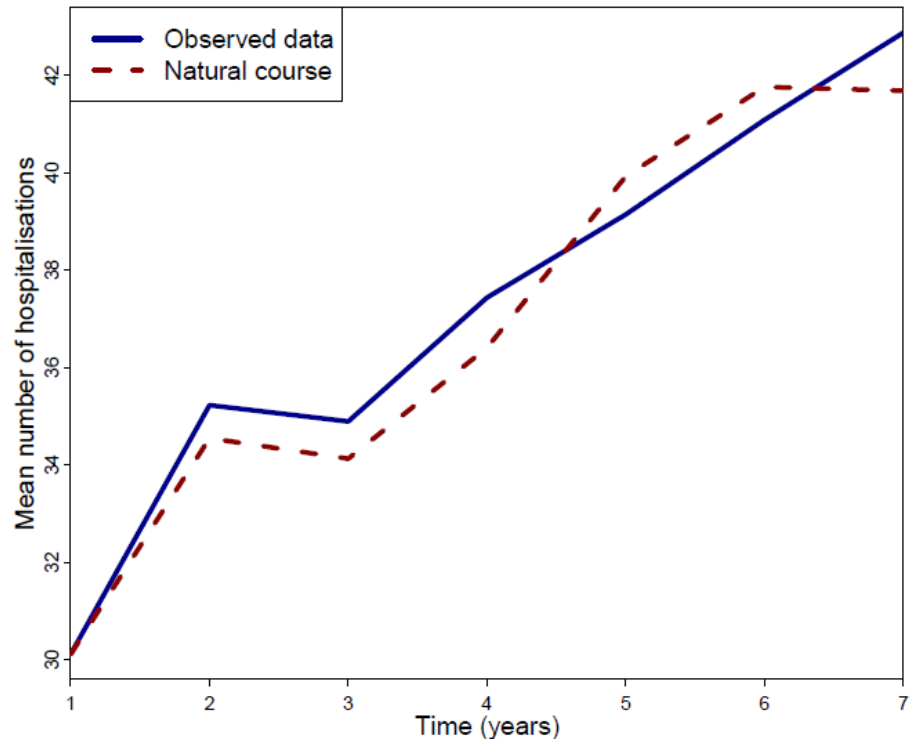
Right column: differences between observed and natural course estimates (solid lines) and 95% pointwise confidence intervals (dotted lines)

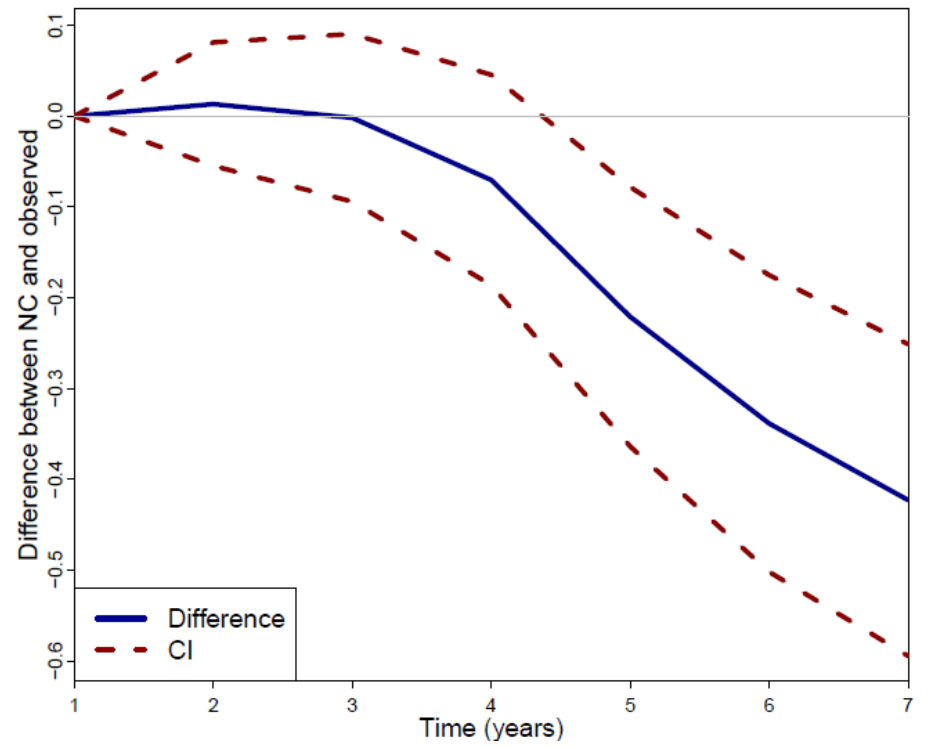
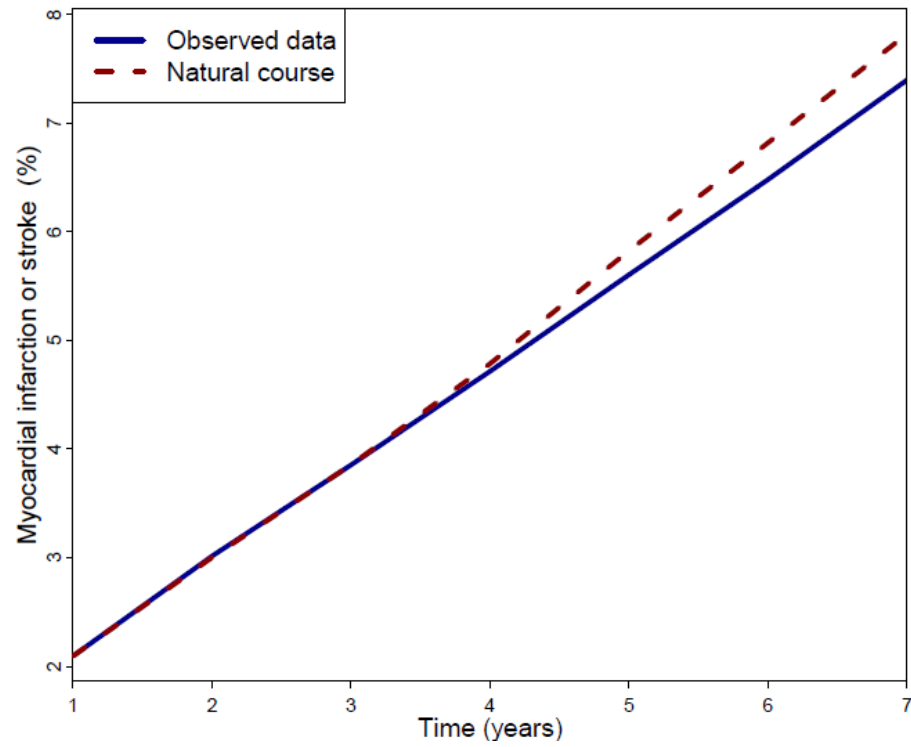
Left: Observed means (solid lines) and predicted means under the natural course (dotted lines) for time-varying covariates; right: differences between observed and natural course estimates (solid lines) and 95% pointwise confidence intervals (dotted lines)

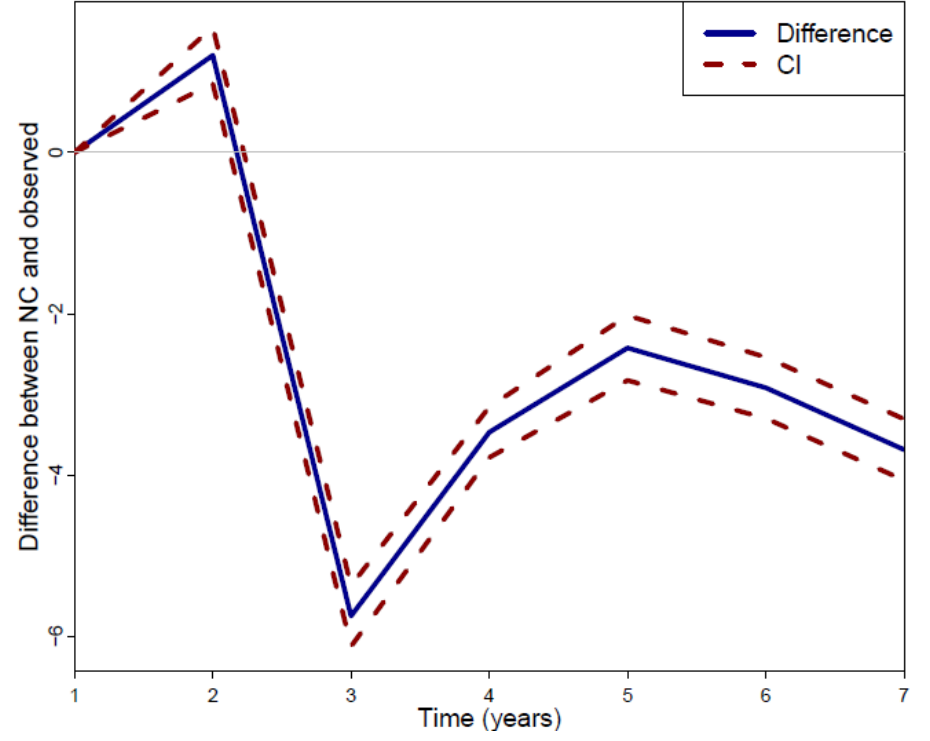
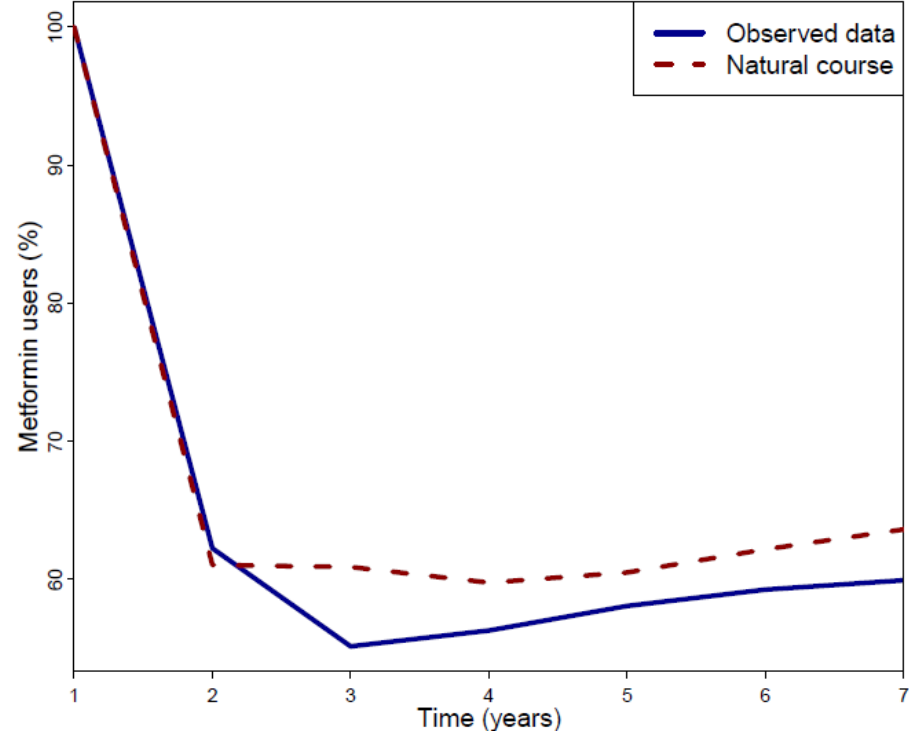


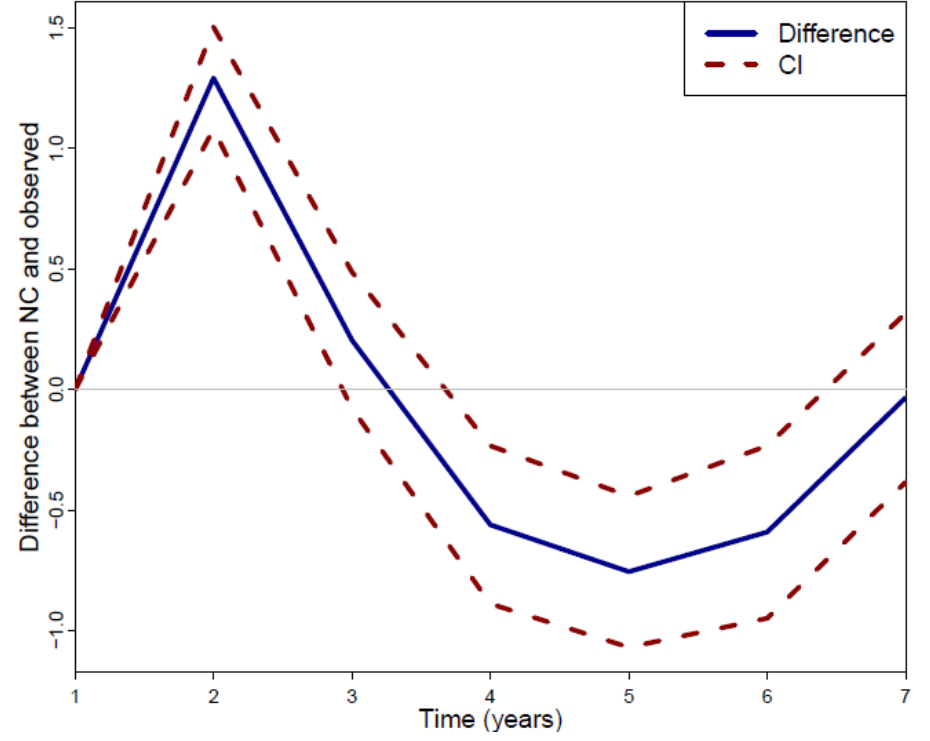
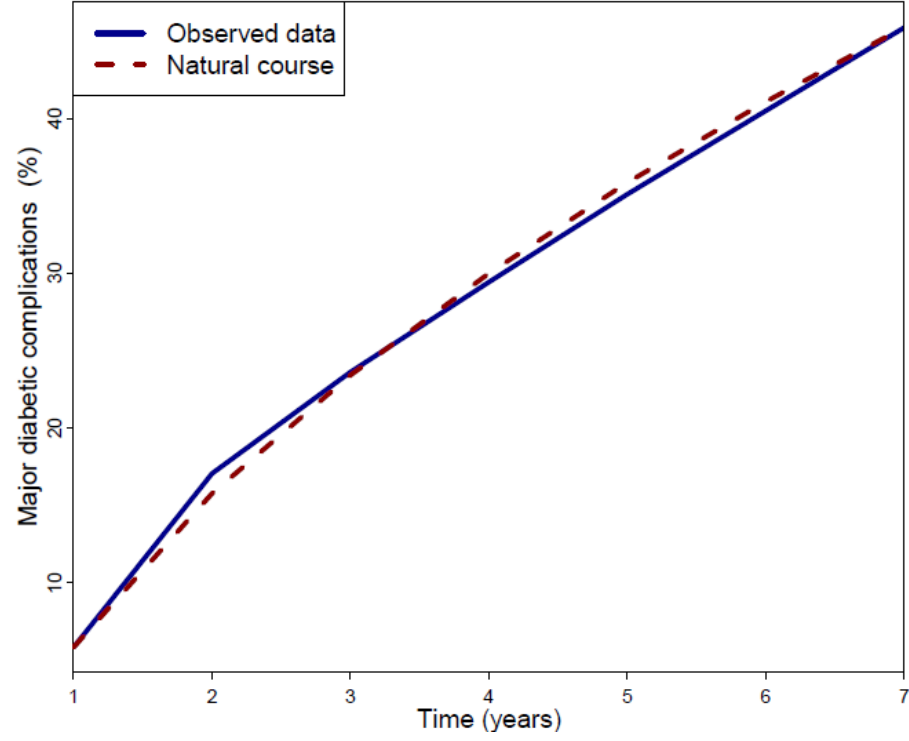


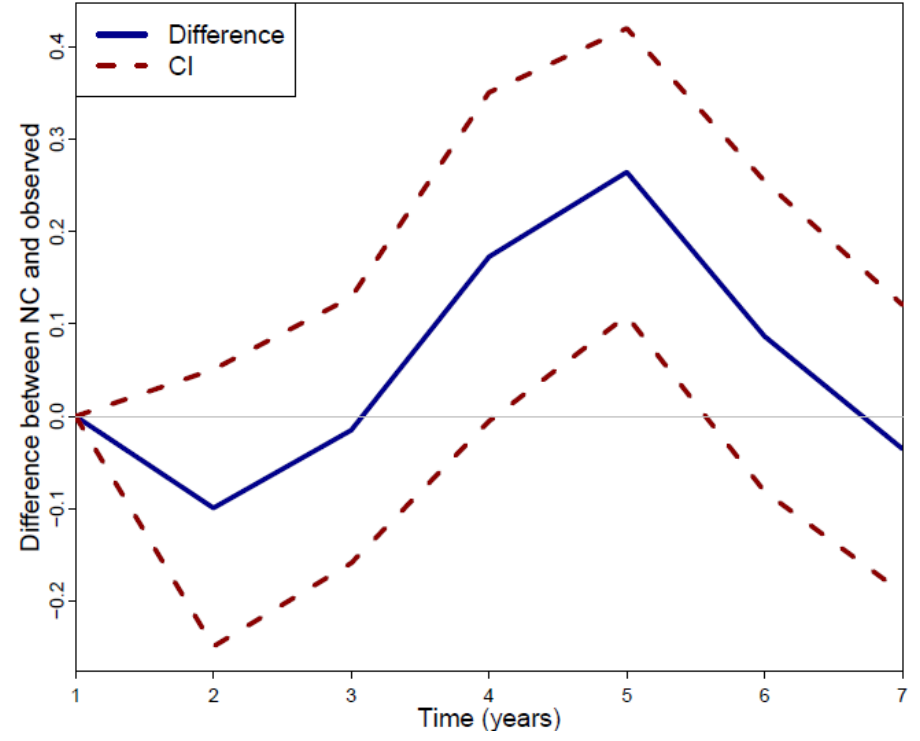
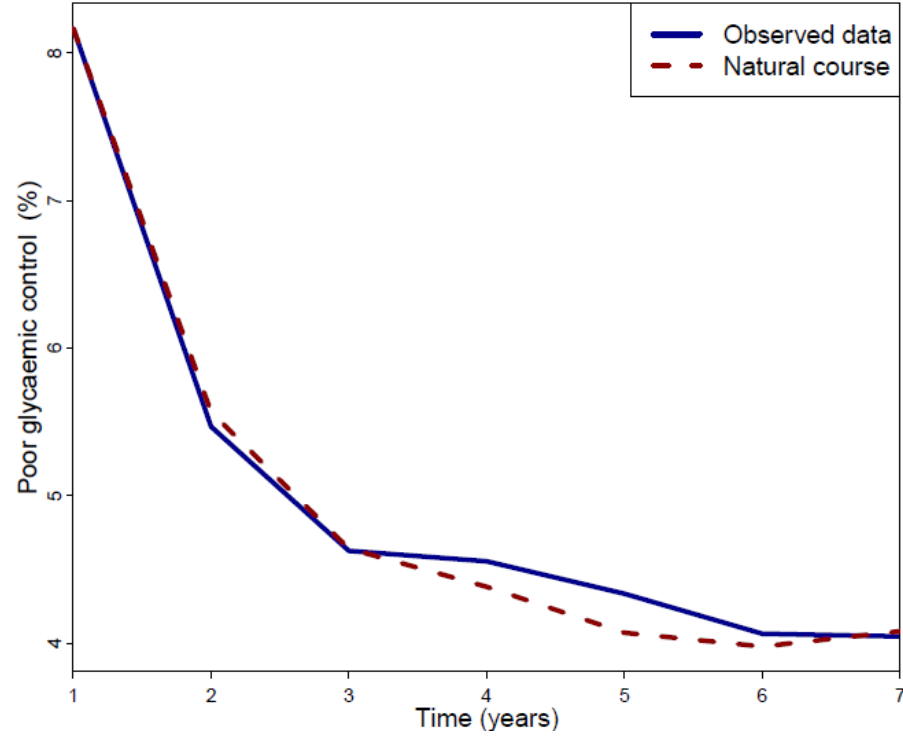


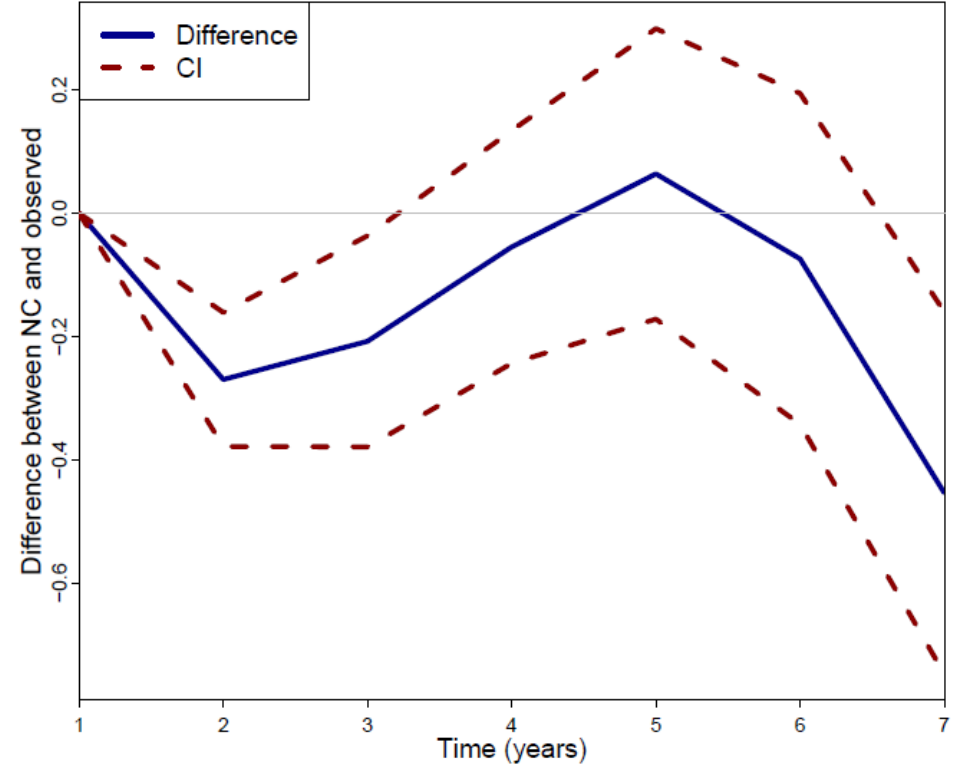
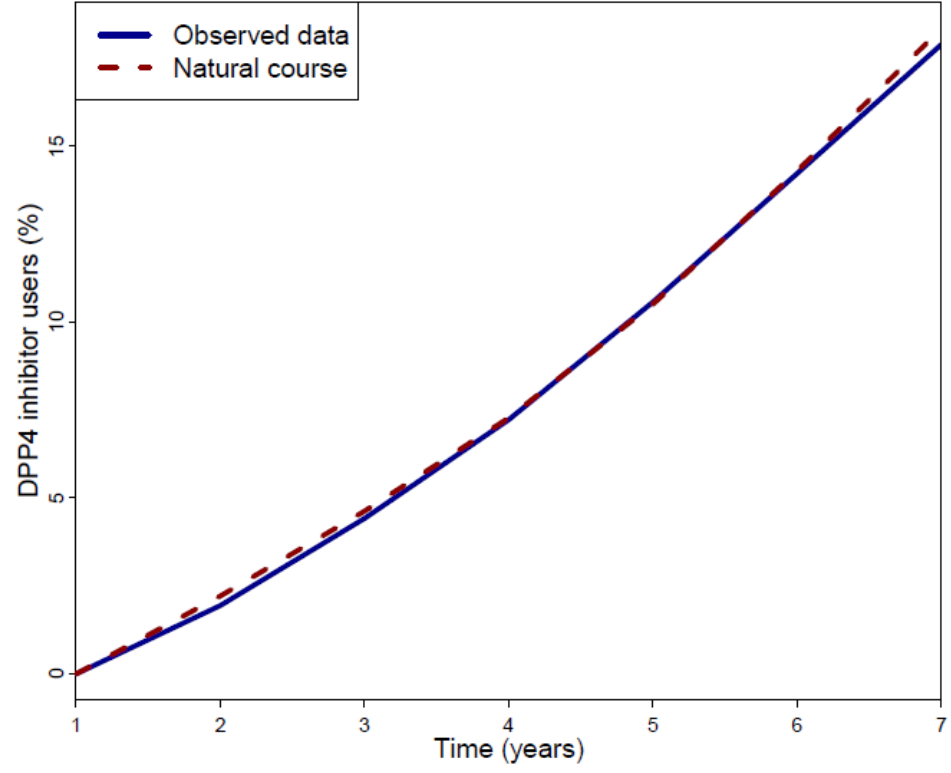


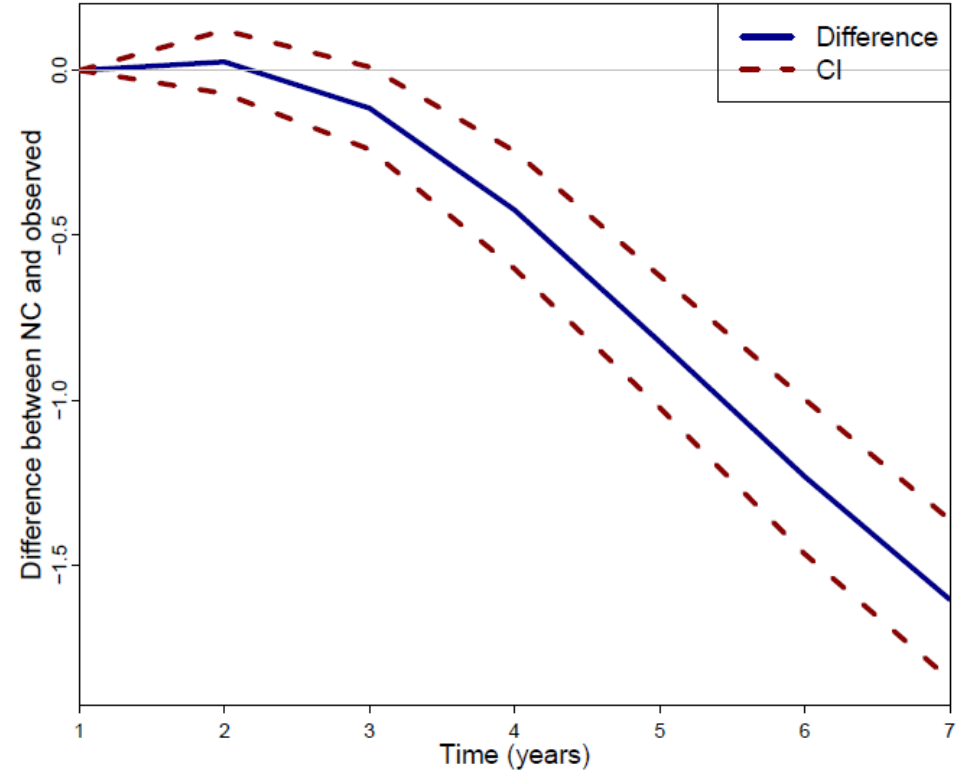
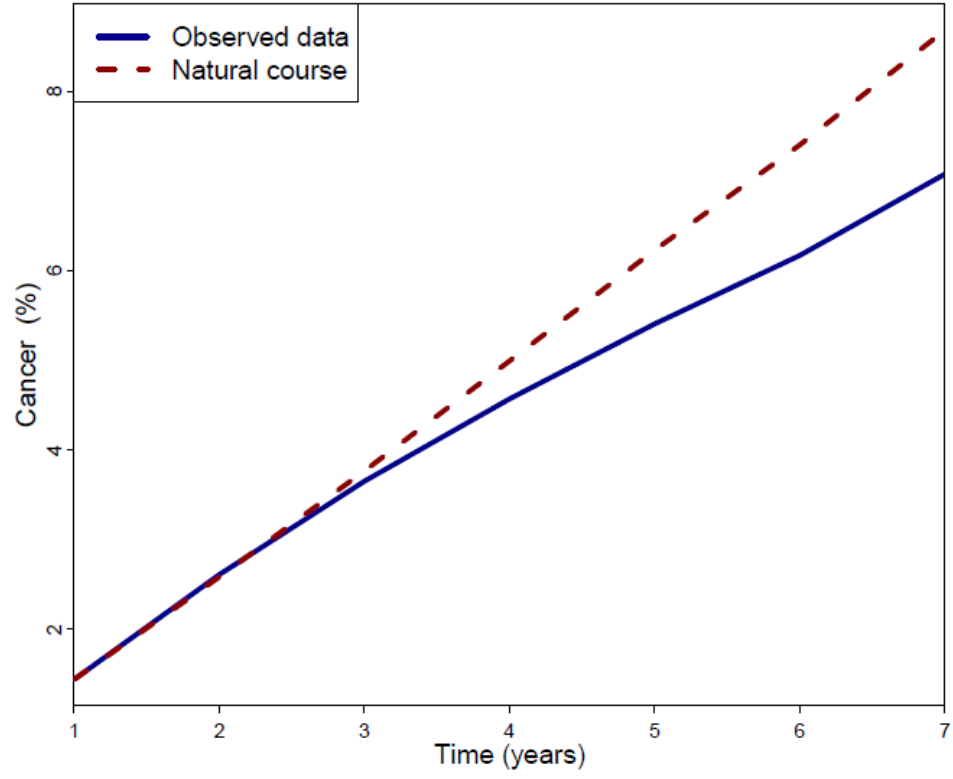












Supplementary Material 11

In the following, the results of several sensitivity analyses are summarized.

Corresponding risk estimates are presented in SCD11:

- 1) Using two instead of one year prior to baseline to apply the eligibility criteria as well as change of the arbitrary order of the covariates in the fitting of the models in step 1) of the parametric g-formula algorithm did not markedly alter the results.
- 2) Choosing smaller sample sizes for the simulations (10 000, 35 000 instead of actual sample size of 77,330) lead to problems of quasi-separation in some of the bootstrap samples; confidence intervals for risk ratios and risk differences were much wider as compared to the model with a sample size of 77,330.
- 3) When using a 91-day (quarter) instead of a 365-day interval, the run time increased markedly, convergence problems arose, observed and predicted values of certain variables showed large deviations and certain important covariates could not be considered (e.g. poor glycaemic control) as the models became too complex.
- 4) When using death as censoring event (instead of competing event) the 7-year risk estimates increased slightly as compared to the main model.
- 5) Including further covariates also led to convergence problems which was the reason to finally combine certain variables into one variable (e.g. comorbidity score).
- 6) Changing the cut-off for the proportion of day covered by the medication to check for exposure misclassification (use a PDR of >0.8 instead of >0.5 as a stricter criterion for treatment classification) resulted in slightly higher estimates for the risk ratio and risk difference with confidence intervals being wider. When choosing the more stringent cut-off of 0.8, the number of observed patients being treated with metformin and/or DPP-4i is smaller by definition. This potentially leads to a selection effect (patients with lower doses are left out) and further increases the problem of violating the positivity assumption. This corroborated the decision to use a cut-off of 0.5.

Supplementary Material 12

Model	Cases	Treatment strategy	7-year risk (%)	95% CI	Risk ratio	95% CI	Risk difference (%)	95% CI
Main	652	Metformin	0.86	0.79 – 0.96	Ref		Ref	
		Met / DPP-4i	1.26	0.94 – 1.69	1.47	1.07–1.94	0.40	0.07 – 0.82
Additional adjustment for visits at diabetologist	652	Metformin	0.87	0.79 – 0.94	Ref		Ref	
		Met / DPP-4i	1.27	0.95 – 1.72	1.47	1.08–1.94	0.41	0.08 – 0.82
2-y exclusion criteria	558	Metformin	0.90	0.80 – 0.99	Ref		Ref	
		Met / DPP-4i	1.26	0.89 – 1.59	1.40	0.97 – 1.83	0.36	-0.03 – 0.73
Reordering of covariates	652	Metformin	0.88	0.80 – 0.98	Ref		Ref	
		Met / DPP-4i	1.33	1.01 – 1.77	1.51	1.09 – 1.98	0.45	0.09 – 0.89
Simulation based on 35,000 ^a	652	Metformin	0.86	0.72 – 0.99	Ref		Ref	
		Met / DPP-4i	1.27	0.75 – 1.70	1.46	0.88 – 2.13	0.40	-0.11 – 0.86
Simulation based on 10,000 ^b	652	Metformin	0.87	0.61 – 1.14	Ref		Ref	
		Met / DPP-4i	1.27	0.44 – 2.86	1.47	0.52 – 3.65	0.40	-0.40 – 2.23
91-day intervals (instead of 365 days) ^c	652	Metformin	0.89	0.81 – 1.00	Ref		Ref	
		Met / DPP-4i	1.34	1.02 – 1.63	1.49	1.17 - 1.83	0.44	0.15 – 0.73
No competing risk (death is censored)	652	Metformin	0.93	0.84 – 1.03	Ref		Ref	
		Met / DPP-4i	1.36	1.02 – 1.85	1.47	1.07 – 1.97	0.44	0.07 – 0.91
Change of cut-off for PDC	652	Metformin	0.89	0.68 – 1.14	Ref		Ref	
		Met / DPP-4i	1.38	0.87 – 2.17	1.55	1.02 – 2.33	0.49	0.02 – 1.08

eTable 3: Risks of development of pancreatic cancer under a sustained metformin monotherapy and under a combination therapy with DPP-4-inhibitors in our main model and certain sensitivity analyses; GePaRD data 2007 – 2017

^aQuasi-separation occurred in 1 out of 100 bootstrap samples; results may be interpreted with caution

^bQuasi-separation occurred in 60 out of 100 bootstrap samples; results should be interpreted with great caution

^cPoor glycaemic control could not be included in the 91-day interval model due to convergence problems; observed and predicted values of certain covariates (e.g. metformin) showed large discrepancies

References

1. Suissa S, Dell'Aniello S. Time-related biases in pharmacoepidemiology. *Pharmacoepidemiol Drug Saf* 2020;**29**(9):1101-1110.
2. Suissa S, Azoulay L. Metformin and the risk of cancer: time-related biases in observational studies. *Diabetes Care* 2012;**35**(12):2665-73.
3. Hernán MA, Sauer BC, Hernández-Díaz S, Platt R, Shrier I. Specifying a target trial prevents immortal time bias and other self-inflicted injuries in observational analyses. *J Clin Epidemiol* 2016;**79**:70-75.
4. Suissa S. Immortal time bias in observational studies of drug effects. *Pharmacoepidemiol Drug Saf* 2007;**16**(3):241-9.
5. Suissa S. Immortal time bias in pharmaco-epidemiology. *Am J Epidemiol* 2008;**167**(4):492-9.
6. Suissa S, Azoulay L. Metformin and cancer: mounting evidence against an association. *Diabetes Care* 2014;**37**(7):1786-8.
7. Yang X, Kong AP, Luk AO, et al. Validation of methods to control for immortal time bias in a pharmacoepidemiologic analysis of renin-angiotensin system inhibitors in type 2 diabetes. *J Epidemiol* 2014;**24**(4):267-73.
8. Danaei G, Robins JM, Young JG, et al. Weight Loss and Coronary Heart Disease: Sensitivity Analysis for Unmeasured Confounding by Undiagnosed Disease. *Epidemiology* 2016;**27**(2):302-310.
9. Robins JM. Causal models for estimating the effects of weight gain on mortality. *Int J Obes (Lond)* 2008;**32 Suppl 3**:S15-41.
10. Yang XL, Ma RC, So WY, et al. Addressing different biases in analysing drug use on cancer risk in diabetes in non-clinical trial settings--what, why and how? *Diabetes Obes Metab* 2012;**14**(7):579-85.
11. Danaei G, Tavakkoli M, Hernán MA. Bias in Observational Studies of Prevalent Users: Lessons for Comparative Effectiveness Research From a Meta-Analysis of Statins. *American Journal of Epidemiology* 2012;**175**(4):250-262.
12. Stattin P, Bjor O, Ferrari P, et al. Prospective study of hyperglycemia and cancer risk. *Diabetes Care* 2007;**30**(3):561-7.
13. Wuermli L, Joerger M, Henz S, et al. Hypertriglyceridemia as a possible risk factor for prostate cancer. *Prostate Cancer Prostatic Dis* 2005;**8**(4):316-20.
14. Yang X, Ko GT, So WY, et al. Associations of hyperglycemia and insulin usage with the risk of cancer in type 2 diabetes: the Hong Kong diabetes registry. *Diabetes* 2010;**59**(5):1254-60.
15. Young JG, Stensrud MJ, Tchetgen Tchetgen EJ, Hernán MA. A causal framework for classical statistical estimands in failure-time settings with competing events. *Statistics in Medicine* 2020;**39**(8):1199-1236.
16. Robins J, Hernan M. Estimation of the causal effects of time-varying exposures. In: Fitzmaurice G, Davidian M, Verbeke G, Molenberghs G, eds. *Advances in Longitudinal Data Analysis*. New York: Chapman & Hall/CRC Press, 2009;553-599.
17. Robins J. The control of confounding by intermediate variables. *Statistics in Medicine* 1989;**8**(6):679-701.
18. Daniel RM, Cousens SN, De Stavola BL, Kenward MG, Sterne JA. Methods for dealing with time-dependent confounding. *Stat Med* 2013;**32**(9):1584-618.
19. McGrath S, Young J, Hernán M. Revisiting the g-null paradox. *arXiv* 2021;**eprint 2103.03857**.
20. Zhang Y, Young JG, Thamer M, Hernan MA. Comparing the Effectiveness of Dynamic Treatment Strategies Using Electronic Health Records: An Application of the Parametric g-Formula to Anemia Management Strategies. *Health Serv Res* 2018;**53**(3):1900-1918.
21. Westreich D, Cole SR, Young JG, et al. The parametric g-formula to estimate the effect of highly active antiretroviral therapy on incident AIDS or death. *Stat Med* 2012;**31**(18):2000-9.
22. Young JG, Cain LE, Robins JM, O'Reilly EJ, Hernan MA. Comparative effectiveness of dynamic treatment regimes: an application of the parametric g-formula. *Stat Biosci* 2011;**3**(1):119-143.