# Dissecting the causal question underlying the association between cancer and ADRD

L. Paloma Rojas-Saunero

## 1. Introduction

Many observational studies have consistently found that individuals with cancer have a lower risk of developing of Alzheimer´s disease or related dementias (ADRD) when compared to individuals with no history of cancer1–4. These findings have motivated substantial research toward mechanistic explanations, including searching for and hypothesizing that molecular and genetic mechanisms may explain this association5–12. These research inquiries inevitably lead to discussions of repurposing or augmenting current cancer chemotherapeutics for ADRD13.

Nevertheless, inferring any treatment or mechanistic effects from the observed cancer-ADRD inverse association is not straightforward. Researchers have raised concerns related to the competing event of death, unmeasured confounding, and ascertainment error that could explain these results9,14. However, understanding these or other sources of bias first requires that we make explicit the causal question. Moreover, making explicit the causal question is one step toward tying a research study to a question that is relevant to decision-making.

To illustrate the complexities of inferring hypothetical or available treatments' effects on ADRD from the observed cancer-ADRD association, we focus on the Pin1 enzyme. Previous animal studies have shown that the Pin1 enzyme over-expression promotes tumorigenesis, while its down-regulation is attributed to mechanisms that contribute to neurodegeneration and amyloid deposition11,12,15. If we one day could develop a drug that increases Pin1 expression specifically in brain tissue in hopes of preventing dementia, we could pose the question as: *What is the direct effect of this Pin1-targeting drug on the risk of ADRD over time compared to standard treatments?*

To explore how we might learn about this effect using real-world data on cancer and ADRD, we progressively build a causal directed acyclic graph to connect this particular causal question to the observable data and the assumptions we rely on to identify the effect. We exemplify different scenarios with data collected from the Rotterdam Study, a population-based cohort study. We describe the challenges and how they translate into the analytic decisions. Last, we discuss how information on mortality and cause of death can provides insight about the direction of some sources of bias.

## 2. Overview of the causal structure

If this hypothetical Pin1-targeting drug was developed, the best way to understand its effect on dementia risk would be to do a well-conducted randomized trial in which we randomize eligible participants in late midlife (e.g., ages 50-60 years) to receive this drug or not, and closely monitor ADRD over a lengthy follow-up.

Since there is no drug currently available that targets Pin1, at best we can use observational data on Pin 1 expression measurements. For example, suppose that a biomarker test was available to measure Pin1 and we measured this biomarker from (stored) baseline blood samples in a population based-cohort that recruited participants in late midlife. Since the biomarker Pin1 is measured within an observational study, confounding can explain an observed association between it and ADRD. In Figure 1, we show that, Pin1 expression ($P_0$) and ADRD at time $t+1$ ($Y_{t+1}$) may share causes $L$, and that assessing the causal relationship

requires adjusting for these confounders $L$. Previous studies have described age, sex, educational level and race/ethnicity as the minimal adjusting set of covariates (Ospina). However, environmental and behavioral factors such as smoking, which are known to cause microenvironmental changes such as inflammation and changes in tissue remodeling, may translate into Pin1 over-expression and are also related to the development of ADRD.

For simplicity, we treat Pin1 expression as a point intervention and we fix the time-ordering of covariates (that is, we assume $L$ happens prior to $P_0$). In reality, it is possible that Pin1 expression changes over time and is affected by time-varying confounders (like smoking), which produces treatment-confounder feedback loops. Addressing such time-varying confounding would require repeated measurements of $L$ and $P\_0$16.
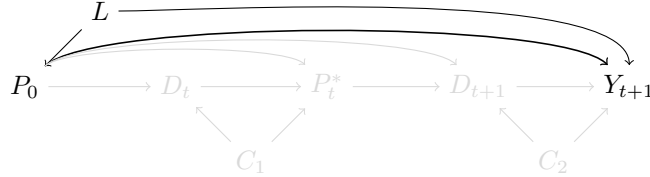


Figure 1: Pin1 as biomarker and risk of dementia

However, unfortunately, in current studies we do not even have a single measurement of Pin1 expression, let alone repeated measurements. Thus, we can only rely on a proxy of this exposure. Since Pin1 over-expression is present in tumors, and tumors are only measured through screening and diagnosis, we considered cancer diagnosis as the proxy for Pin1 over-expression4,17–24. We depict this feature in Figure 2, where $P^*$ represents *incident cancer diagnosis*, the measured proxy of $P_0$. In this causal graph we colored the path $P^* \leftarrow P \leftarrow Y_{t+1}$ because although we are measuring the association between $P^*$ and $Y_{t+1}$ in the observed data, we are assuming that the captured effect is only through $P_0$.
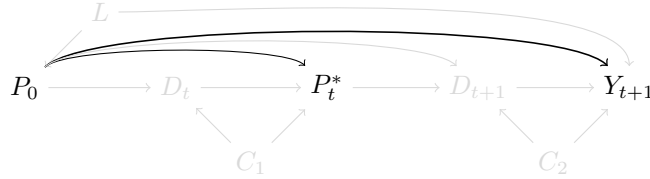


Figure 2: Cancer diagnosis as proxy for Pin1 expression

A major challenge related to assigning cancer diagnosis as the proxy of Pin1 is defining time zero. Not everyone who had Pin1 over-expression will be diagnosed with cancer by late-midlife, but in late-life. For this reason prospective cohort studies have considered cancer diagnosis as a time-varying or time-dependent exposure. This means that a participant within a prospective cohort study contributes to the "regulated Pin1" arm since study entry up to the time of cancer diagnosis and later on to the "Over-expressed Pin1" arm. Other studies have included participants with cancer diagnosis at the time of the diagnosis (for example from cancer registries19,23,25,26) and matched participants by age. In both cases we must remember that the main interest is on the unmeasured $P_0$ thus we should only adjust for covariates prior to $P_0$ and be careful to adjust for post-baseline covariates of $P_0$ or mediators between $P_0$ and $P^*$.

One of the mediators between $P_0$ and $P^*$ is death prior to cancer diagnosis, represented as $D_t$ in Figure 3. Since individuals are at risk of dying from other causes (such as cardiovascular death), we can only measure $P^*$ in the subset of individuals who have survived long enough to have a cancer diagnosis. Several risk factors that increase the risk of cancer might also cause death prior to cancer diagnosis, for example smoking may cause lung cancer and chronic obstructive pulmonary disease (a leading cause of death in this age group27). Therefore, we can only isolate the effect of $P_0$ in $P^*$ if we condition on $C_1$ and block the backdoor pathway that is open by conditioning on $D_t$. In this way we assume a hypothetical scenario in which we could prevent death prior to cancer diagnosis, by conditioning on a rich set of covariates28. This

assumption must hold regardless of whether we use incident cancer as a time-varying exposure, or match cancer patients to participants free of cancer by age.

Furthermore, we note that this is only one of the issues with considering cancer diagnosis as the proxy for Pin1 expression in terms of information bias(ref). We could add more complexity by considering further shared causes of cancer and ADRD diagnosis, or further mediators between Pin1 expression and cancer diagnosis, including but not limited to: screening guidelines, type of healthcare coverage, health-seeking behaviors, and healthcare availability.
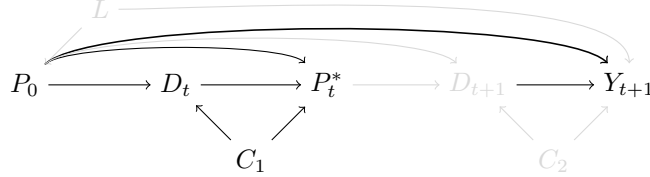


Figure 3: Survival bias when treating cancer diagnosis as the proxy for Pin1 expression

Up to this point we have outlined certain underlying assumptions of using incident cancer diagnosis as a proxy for Pin1 expression. However, as mortality increases steeply in late life, even in the setting of the ideal randomized trial of a hypothetical Pin1-targeting drug, we can only measure ADRD over follow-up in the individuals who survive long enough to have a diagnosis. For this reason, death is a competing event of ADRD because if a participant dies prior to ADRD diagnosis, death prevents observing ADRD at future time-points. As such, to identify the direct effect of $P_0$ in $Y_{t+1}$ we have to first assess if the drug has an any side effects that could increase the risk of death in any way. If the drug has no systemic side-effects, such as that there is no arrow between $P_0$ and $D_{t+1}$, such as observed in Figure 4.
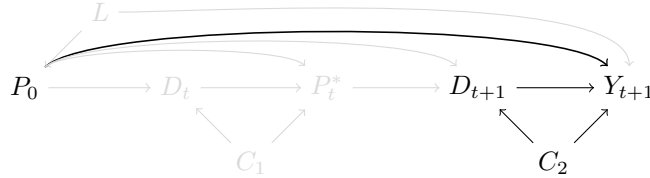


Figure 4: Death as a competing event for dementia over follow-up

However, if the drug has any side-effects we have to decide if either we want to estimate a total effect (that is, the effect on ADRD even if partly mediated by the effect on death) or a direct effect (that is, the effect directly on ADRD not mediated by death) (ref). We face this challenge in our observational data setting, since we are measuring cancer diagnosis as the proxy for Pin1 over-expression, and cancer is a leading cause of death29. We can visualize this data feature in the causal diagram in Figure 5. In this causal diagram we include $D_{t+1}$ as a representation of death over follow-up, an arrow between $P_0$ and $D_{t+1}$ since $P_0$ may increase the mortality risk. We also include an arrow between $P_t^*$ and $D_{t+1}$ since cancer diagnosis, and subsequent treatment (or lack of treatment) may have an effect on death. Last, the arrow from $D_{t+1}$ and $Y_{t+1}$ represents the key feature of a competing events data structure: an individual who dies over follow-up (prior to dementia diagnosis) cannot subsequently develop ADRD, and since $D_{t+1}$ and $Y_{t+1}$ are events related to aging, $C_2$ represent the shared causes of both events.

In this causal diagram we observe that, had we measured and adjusted for $L$, we could estimate the total effect of $P_0$ in $Y_{t+1}$ without further assumptions. However the total effect includes all pathways between $P_0$ and $Y_{20}$28. This means that if Pin1 has an effect on mortality through cancer or other mechanisms, we will observe a protective effect of Pin1 over-expression in ADRD, partially or fully mediated through death. Since the question of interest is focused on the direct effect of $P_0$ in $Y_{t+1}$ as in Figure 4 (violet arrows) we need to conceptualize the different mechanisms through which $P_0$ could affect $D_{t+1}$ and $Y_{20}$. With this in mind we can conceive different causal questions (estimands) to represent this direct effect , such

as the controlled direct effect (CDE) and the natural separable direct effect. In this section we discuss the controlled direct effect as the causal question of interest since it translates to frequently used methods in this literature (such as Kaplan-Meier estimator and Cox-proportional hazard model) and leave the separable direct effects question for discussion.
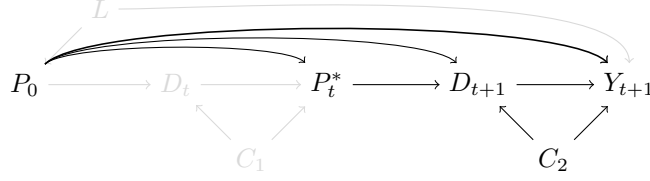


Figure 5: Direct effect

The CDE represents the effect of Pin1 over-expression in ADRD in a setting where we could have prevented death over the study period. It relies on the assumption that we have measured all $C$ to block the pathway $Y20 \leftarrow C \rightarrow D_{15} \rightarrow P \rightarrow P*$. This assumption is defined as the independent censoring assumption conditional on covariates. In this setting death is treated as a censoring event and it can be interpreted as those who died would have the same risk of developing dementia if prevented from dying than those who remained alive and free of dementia at a given time point, conditional on shared causes of dementia and death.

Therefore, if we combine the challenges in section two related to cancer diagnosis as a proxy for Pin1 over-expression, and having death as a competing event of ADRD we observe the complexity of the DAG in Figure 5. Though this is yet a simplified version since we omitted time varying $P*$, $L$ and $C's$.
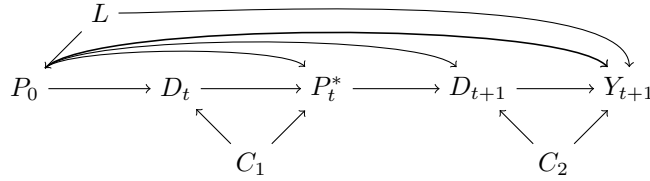


Figure 6: Direct effect