

LATE-CAREER ENTREPRENEURSHIP AND EARNINGS: THE ROLE OF FORMAL RETRAINING AND UPSKILLING*

Finn Lattimore^{1†}, Daniel Steinberg² and Anna Zhu³

¹Reserve Bank of Australia

²Gradient Institute

³RMIT University, IZA

May 29, 2023

Abstract

Pursuing formal retraining or an additional educational degree may constitute an important ingredient in the late-career entrepreneurial path, provide a bridge to later retirement and/or increase earnings in later-career. We focus on the causal impacts to economic returns of degrees completed later in life, where motivations and capabilities to acquire additional education may be distinct from education in early years. We revisit this long-studied question with new methodological techniques by adapting machine learning models for causal inference. We find that completing an additional degree increases the probability of being self-employed by 1-2 percentage points and leads to more than \$3000 (AUD, 2019) per year compared to those who do not complete additional study in later-life. For outcomes, treatment and controls we use the extremely rich and nationally representative longitudinal data from the Household Income and Labour Dynamics Australia survey. To take full advantage of the complexity and richness of these data we use a Machine Learning (ML) based methodology to estimate the causal effect. We are also able to use ML to discover sources of heterogeneity in the effects of gaining additional qualifications.

JEL: J12, J18, H53

*Corresponding author: Anna Zhu, RMIT University. Email: anna.zhu@rmit.edu.au.

We thank Tim Robinson, Hayley Fisher, Bruce Bradbury and numerous seminar and conference participants for helpful comments. The authors would like to thank Tessa LoRiggio, Prabath Abeysekara, Michael Duffield, and Yin-King Fok for their excellent research assistance.

Zhu acknowledges the support of the Australian Research Council Linkage Project (LP170100472). This paper uses unit record data from the Household, Income and Labour Dynamics in Australia (HILDA) Survey. The HILDA Project was initiated and is funded by the Australian Government Department of Social Services (DSS), and is managed by the Melbourne Institute of Applied Economic and Social Research (Melbourne Institute). The findings and views reported in this paper, however, are those of the authors and should not be attributed to either DSS or the Melbourne Institute.

[†]This work was performed while the author was working at the Gradient Institute. Views expressed in this paper are those of the author and do not reflect those of the Reserve Bank of Australia.

Keywords: Machine Learning, education, mature-age learners, causal impacts

1 Introduction

Later-career entrepreneurship has been found to extend careers, provide job autonomy and flexibility (Hundley, 2001), improve health (Stephan and Roesler, 2010), increase life satisfaction (Parker and Rougier, 2007, Kautonen, Kibler and Minniti, 2017, Lombard, 2001), and lift earnings (Levine and Rubinstein, 2017, Astebro and Chen, 2014). Such benefits can also improve individuals' retirement security by delaying pension claiming and increasing savings for retirement (Kibler et al., 2012), as well as provide societal benefits via government budget savings in pension payouts and welfare programs (Herneas et al., 2016, Kautonen, Kibler and Minniti, 2017). On the other hand, later-career entrepreneurship may have costly barriers to entry, such as the need to obtain an additional qualification. The transition into entrepreneurship can not only take longer than expected, but can present obstacles (age-discrimination) and compete with other commitments (such as caring responsibilities) (Kibler et al., 2012, Weber and Schaper, 2004). For later-career entrepreneurship to be individually rewarding and socially sustainable it is important to understand which transition pathways elevate entrepreneurial and financial success, and which transition pathways tend to be counter-productive.

In this paper we ask: does formal retraining or upskilling constitute an important ingredient in the late-career entrepreneurial path? And if so, is it economically rewarding for individuals to undertake such formal retraining or upskilling? We also aim to understand which individuals generate higher versus lower returns to formal retraining and upskilling. Specifically, which factors explain the relationship between formal retraining and later-life economic and entrepreneurial success. We depart from the returns to education literature by focusing on education at older ages (25 and over) as opposed to early adulthood. We focus on retraining and upskilling by studying individuals who return to formal education¹ after a break from education acquired in early adulthood. In other words, we focus on degrees that are additional to the first educational degree obtained.

Despite the increasing prevalence of formal retraining or upskilling, this is an understudied topic in both the education and entrepreneurship literature. Until now the literature has focused on the cumulative years of education or early-years educational attainment, such as college or high school degrees. We argue that there is value in focusing on the impact of *retraining* separately from cumulative education years, especially since returns to education are more uncertain for older students who face higher opportunity costs and

¹We focus on formally recognised degrees ranging from short degrees (6 months to 2 years full-time) to postgraduate degrees such as masters and PhDs. Most degrees undertaken in our sample are Certificate 3 or 4 degrees which take an average of 1-2 years to complete. See Figure 5 for the breakdown of the type of formal degrees considered.

complex entry rules. This also allows us to disentangle formal retraining and upskilling from experience in the workforce (or through, for example, volunteering pursuits). This is an important distinction since there is likely a trade-off between these two ingredients in entrepreneurial success (Zacharakis and Meyer, 2000). By doing so we contribute to the strand of literature that examines the human capital impacts on entrepreneurship (Unger et al., 2011).

Despite the significant increase in the uptake of degrees in later life (OECD, 2016), it is unclear, a-priori, if formal retraining and upskilling is an entrepreneurial and economically rewarding endeavour. On the one hand, investing in education in later life may be a promising pathway to starting a business. For example, it may assist individuals to update skills that technological change and automation may demand (Autor, Katz and Kearney, 2008, Acemoglu and Autor, 2011) or help them to retrain in new skills if self-employment involves an occupational pivot. For others, the transition into entrepreneurship from unemployment may require some formal retraining or upskilling due to their disconnection from the workforce. Older individuals who wish to start a business (after the age of 40) may be more likely to face technology-driven erosion of skills, compared to younger individuals, which creates a greater need for retraining. However, older individuals may not reap the full rewards of retraining as they have less time to accrue gains from educational investment. While the literature finds strong positive returns to education for skilled entrepreneurs (Fossen and Büttner, 2013)² and for college degree completion following secondary school (Jacobson, LaLonde and Sullivan, 2005, Chesters, 2015, Zeidenberg, Scott and Belfield, 2015, Polidano and Ryan, 2016, Xu and Trimble, 2016, Belfield and Bailey, 2017, Dynarski, Jacob and Kreisman, 2016, 2018, Mountjoy, 2022), the returns for older individuals are unclear given the trade-offs associated with retraining.

On the other hand, obtaining an additional education in later life can have large opportunity costs in terms of forgone experience and earnings (Haber and Reichel, 2007, Montgomery and Powell, 2006) and these ingredients may be more important to entrepreneurial success (Zhao et al., 2021, Weber and Schaper, 2004). The financial and time constraints of retraining and upskilling are also likely to be higher in later life than during adolescence³, which can explain why drop-out rates among mature-age learners is relatively high (Cherastidtham and Norton, 2018, Hanson, 2022, Guerra, 2022). The

²The returns are estimated to be higher for opportunity entrepreneurs (those who voluntarily start their own business) versus necessity entrepreneurs (those who lack alternative employment options). We do not make this distinction because the sample size of older entrepreneurs is already small.

³Abel and Deitz (2014) estimate forgone earnings of a 4-year college degree to be around \$23,000 (2013 USD) while Montgomery and Powell (2006) estimate forgone earnings of an MBA, a graduate management program, to be higher at around \$28,000 (2013 USD).

costs of financing education, such as higher debt, can itself hinder entrepreneurship (Krishnan and Wang, 2019). The type of retraining is also relevant. Entrepreneurial specific educational programs may have little impact on students’ own self-assessed proficiency of their entrepreneurial skills and can even have a negative impact on their intention to be an entrepreneur (Oosterbeek, Van Praag and Ijsselstein, 2010). Thus, for later-life educational attainment to be a worthwhile endeavour, a key unanswered question is whether it is economically rewarding for individuals.

This paper contributes to addressing this gap in knowledge by estimating the causal effect of formal retraining and upskilling on earnings and entrepreneurship. We define entrepreneurship as ‘becoming’ self-employed in later life. This includes individuals who transition from being an employee or non-employed to being self-employed after the age of 40 versus those who (1) remain employees, (2) remain self-employed, or (3) remain non-employed. We compare the entrepreneurial success and earnings of individuals who undertook formal retraining after the age of 25 (our treated group) with those who did not obtain a further degree after the age of 25 (our control group). Our retraining variable is defined before entry into entrepreneurship, which overcomes possible reverse causation (Cassar, 2006).

Our causal and heterogeneity analysis utilises, at the outset, every piece of background information that is asked about an individual⁴ in a rich, nationally representative survey dataset. Effectively, we begin with over 3,400 characteristics or features about the individual, which enables a stronger like-for-like comparison of treated and control groups, as well as deeper insights into how the estimated returns to formal retraining and upskilling may vary by background or circumstance. Specifically, we use 19 waves of nationally representative panel data (from 2001-2019) from the Household Income and Labour Dynamics Australia (HILDA) survey. Our empirical strategy is to adapt Machine Learning (ML) techniques for the aim of causal analysis. We need to adapt these tools because they are traditionally used for predictive purposes. The three key benefits of the ML approach over traditional methods is (1) it effectively acts as a hyper-propensity-matching technique that utilises as much information as possible to construct similar groups of individuals in the treated and control groups; (2) it systematically reduces information redundancy between variables that are highly correlated, so it addresses overfitting issues; and (3) it identifies new sub-populations for which the treatment effects are different, which means we avoid missing out on important heterogeneities that may arise from pre-specifying sub-groups (Athey and Imbens, 2017, Knaus, Lechner and Strittmat-

⁴Variables included span from information about their family and neighbourhood to the individual’s intention to change job or industry, predicted retirement age, caring responsibilities, financial debt, economic circumstances, demographic characteristics and the earnings trajectories prior to retraining.

ter, 2021). We use the ML approach combined with the full power of an incredibly rich survey dataset to overcome the main challenge in estimating the causal returns to formal retraining and upskilling: selection. The factors that enable or motivate mature-age learners to pursue and complete retraining or upskilling may also be independent precursors to later-life earnings or entrepreneurial success. For example, those who intend to enter entrepreneurship in later life may invest more in retraining themselves for the role. Alternatively, previous studies have found that individuals who seek a further community college degree tend to have slower growing earnings in the period before their study starts compared to similar individuals who do not seek further study (Jacobson, LaLonde and Sullivan, 2005, Dynarski, Jacob and Kreisman, 2016, 2018).

A key empirical contribution of this paper is that we adapt ML tools for causal inference purposes. We recognise that, as with all statistical models, we make assumptions when we use ML techniques for causal inference, and these need to be tested. We test the assumption that the additional controls included in the ML models work to address selection-into-treatment issues described above. First, we replicate the results based on published works by Chesters (2015). Then we contrast the ‘systematically’ selected control variables in the ML model with those that were manually selected in the Ordinary Least Squares (OLS) model by Chesters (2015), and comment on the potential biases from manual variable selection. We have chosen this published work because it uses the same data (HILDA) and examines the same topic. Second, we compare the results of the ML model with models traditionally used to address unobserved heterogeneity such as a Fixed Effects Difference-in-Difference approach. Across all results, we find that the ML models tend to estimate lower returns to retraining compared to these baseline models.

Human capital theory, employment choice theory, and self-determination theory provide valuable frameworks for examining the impact of retraining and upskilling on earnings and entrepreneurship. Human capital theory (Becker, 1964) recognises that acquiring additional education can improve human capital – an individuals’ skills and knowledge – which can enhance productivity and earnings potential. Individuals who engage in retraining and upskilling may acquire new competencies and knowledge that positively impact their suitability for entrepreneurship. Employment choice theory (Heckman, 1993) suggests that retraining and upskilling can influence an individuals’ choice to pursue entrepreneurship through the acquisition of new perspectives, specialized knowledge, and expanded networks. It emphasises that the choice of employment type is influenced by a range of factors, such as personal preferences, skills, motivation, financial considerations, and opportunities for growth and fulfilment. Self-determination theory (Deci and Ryan, 1985) reflects on the psychological needs that promote well-being. Retraining and

upskilling that supports autonomy, aligns with personal goals, and provides a sense of competence can provide individuals with the confidence, skills, and motivation necessary to pursue higher-paying positions or start their own business. These theories highlight the complex dynamics involved in later-life education and its potential outcomes for individuals' economic prospects and entrepreneurial pursuits, such as the interplay between human capital investment, individual motivation and personal preferences.

Our results show that an additional degree in later life increases the probability of being self-employed (at age 40 or older) by 1-2 percentage points and increases total future earnings by more than \$3,000 (AUD) per year on average, compared to those who do not complete further study. We consistently estimate this causal effect using a selection-on-observables strategy based on T-learner, Doubly Robust and Bayesian models. As entrepreneurship is a binary variable, we also use Logit and Probit transformations to model the non-linearity in parameters.

Our ML approach also identifies new sub-populations for which the treatment effects are different. For entrepreneurship, we document that behavioural factors, such as the age at which a person left home and the division of labour in the household, are the most significant factors related to benefits from retraining. Household income is also hugely important: those with above median household income tend to benefit more from retraining than those with below median household income. This aligns with the entrepreneur literature that finds higher household income prior to becoming self-employed is associated with more ambitious revenue targets ([Cassar, 2006](#)). For earnings, the most significant factors affecting returns to retraining are the starting levels of and pre-study trends in personal and household income. Age and wealth-related variables also account for variation in estimated effects. All of these variables are consistently selected as being significant for prediction out of the 3,400 features within the HILDA data. This selection is consistent across different ML models (which includes linear and non-linear model classes) and across numerous bootstrap draws of the original sample.

2 Context: Higher education and Vocational Study in Australia

Mature-age education in Australia is among the highest in the world. In 2014, Australia's participation in vocational education by those aged 25-64 was the highest among OECD countries. The tertiary education rate for those aged 30-64 was the second highest ([Perales and Chesters, 2017](#)). Mature-age Australians are increasingly enrolling in university or college to change employers, change careers, gain extra skills, improve their promotion

prospects and earning capability or search for better work/life balance. Redundancy and unemployment have also been driving forces for individuals to return to education later in life (Coelli, Tabasso and Zakirova, 2012).

The increase in mature-age learners accessing higher education has in part been driven by government policy. In 2009, the Australian government adopted a national target of at least 40% of 25-34-year-olds having attained a qualification at bachelor level or above by 2025 (O'Shea, May and Stone, 2015). This was part of a policy that transitioned Australia to a demand-driven system (Universities Australia, 2020). The policy had a large effect on access to higher education, as it removed the cap on the number of university student places. By 2017, 39% of 25-34-year-olds had a bachelor's degree or higher (Caruso, 2018). The demand-driven system effectively came to an end in 2018 when the government capped funding at 2017 levels. While the initial uptake of higher education in the demand-driven system was strong, especially among mature-age students (Universities Australia, 2019), the cap coincided with a decline in enrolments for both undergraduate and vocational courses (Universities Australia, 2020, Atkinson and Stanwick, 2016, NCVET DataBuilder, 2021).

The cost of a bachelor's degree for domestic students in Australia is the sixth highest among OECD countries (Universities Australia, 2020). In 2018, the average annual cost of a bachelor's degree was around \$5,000. VET and TAFE/college courses in Australia cost a minimum of \$4,000 per year on average while post-graduate courses cost a minimum of \$20,000 per year on average (Studies in Australia, 2018). Mature-age students can cover the cost of further study themselves or they can receive support from the government. Many undergraduate students can access the Commonwealth Supported Place (CSP) scheme which subsidises tuition fees for those studying at public universities and some private higher education providers. Postgraduate students are generally not covered by the CSP. Students at university, approved higher education providers or in VET can access financial support from the Higher Education Loan Program (HELP) scheme, which provides income-contingent loans. This allows students to defer their tuition fees until their earnings reach the compulsory repayment threshold, upon which repayments are deducted from their pay throughout the year at a set rate. CSPs and HELP loans, however, are withdrawn from students who fail half of their subjects.

3 Data

We use data from the Household Income and Labour Dynamics Australia (HILDA) survey. These data are rich, and we exploit the full set of background information on

individuals (beginning with more than 3,400 variables per observation). HILDA covers a long time span of 19 years, starting in 2001. We use the 2019 release. This means we observe respondents annually from 2001 to 2019.

3.1 Sample exclusions

Our main analysis sample contains respondents who were 25 years or above in 2001. This allows us to focus on individuals who obtain a further education beyond that acquired in their previous degree. Our main analysis focuses on measuring the impact of further education using wave 19 outcomes. Here, the feature inputs to the models are taken from the individuals in 2001. We drop any individuals who were ‘currently studying’ in 2001. This also ensures that our features, which are defined in 2001 are not contaminated by the impacts of studying but clearly precede the study spell of interest. These sample exclusions result in 7,359 respondents being dropped because they are below the age of 25 in 2001 and a further 1,387 respondents being dropped because they were studying in 2001. We then restrict the sample to those who are present in both 2001 and 2019. This ensures that we observe base characteristics and outcomes for every person in our analysis sample. This results in a further 5,727 respondents being dropped from the sample. Our analysis sample has 5,441 observations. More details of our main analysis sample and data can be found in Online Appendix [F](#).

3.2 Outcomes

We code the values of the outcome variables based on the survey panel data from Wave 17 onwards. This is to allow us to measure the long-term impact of formal re-training. It also ensures the outcome is measured ‘after’ the individual has started their formal retraining. This minimises the chance that our analysis results are subject to reverse causation issues.

For the outcome of entrepreneurship, we include individuals who ‘became’ self-employed from Wave 17 onwards. This means they were observed to not be self-employed in the previous wave and then subsequently transitioned into self-employment. For the outcome of earnings, we only use the 2019 wave to measure this outcome. This means that we measure earnings at least 2 years after an individual completes their formal re-training. We use annual earnings to measure the economic returns to education. Last, we also analyse outcomes related to the labour market such as employment, changes in earnings, changes in occupation, industry, and jobs. These are also measured after 2018.

3.3 Treatment

We define further education as an individual who obtains a further degree in a formal, structured re-training or educational program. For short, we refer to this as ‘re-training’ throughout the paper since it characterises the second (or subsequent) degree individuals obtain after their first degree. These programs must be delivered by a certified training, teaching or research institution. Thus, we do not analyse informal on-line degrees (such as Coursera degrees). We also do not consider on-the-job training as obtaining further education.

Our treatment variable is a binary variable that takes the value of 1 if an individual has obtained an additional degree anytime between wave 2 (2002) and wave 17 (2017). We drop any respondent who obtained a qualification after wave 17.

HILDA documents formal degree attainment in two ways. The first is to ask respondents, in every wave, what is their highest level of education. The second way is to ask respondents, in every wave, if they have acquired an additional educational degree since the last time they were interviewed. We utilise both these questions to construct our measure of further education. Using the first question, we compare if the highest level of education in 2017 differs from that in 2001. If there has been an upgrade in educational qualification between these two years, we set the treatment indicator to be one and zero otherwise. This question, however, only captures upgrades in education; it fails to capture additional qualifications that are at the same level or below as the degree acquired previously by the respondent. We rely on the second survey question to fill this gap.

These two survey questions thus capture any additional qualification obtained from 2002 to 2017, inclusive. Additional qualifications refer to the following types of degrees: Trade certificates or apprenticeships, Teaching or nursing qualifications, Certificate I to IV, Associate degrees, Diplomas (2-year and 3-year fulltime), Graduate certificates, Bachelor, Honours, Masters and Doctorate degrees.

3.4 Covariates/features

We define our covariates, or features as they are known in machine learning parlance, using 2001 as the base year. Since we drop any respondents who were currently studying in 2001, we ensure that all features were defined before a respondent begins further study.

A unique approach to our feature selection strategy is that we use all the information available to us from the HILDA survey in 2001. This means that we have more than 3,400

raw variables per observation. Before using the features in a ML model, we delete any features that are identifiers or otherwise deemed irrelevant for explaining the outcome.

In order to reduce redundancy in this vast amount of information, we next apply a supervised Machine learning model to predict outcomes 5 years ahead of 2001 i.e. in 2006. We then select the top 100 variables that are most predictive of the outcome in 2006.⁵ These variables are listed in Table 1.

4 Descriptive Figures and Tables

We calculate the average returns to degree completion for mature-age students who completed degrees between 2002 and 2017. The window in which study and degree-completion took place is noticeably large. However, sample size limitations with our survey data mean that it is not feasible to run an ML analysis, disaggregated by the timing-of-completion.

In order to obtain some insights into the potential heterogeneity over time, we present a series of descriptive graphs in this section. Here, our aim is not to present any causal analysis but to describe which groups studied earlier in the time period (and thus had more time to accumulate returns). These graphs can also point to the potential different factors driving study across the time period, and different effects on earnings depending on how much time has elapsed since completion.

Figure 3 presents the distribution of degree completion over time. There is a steep decline in degree-completion proportions over time. This is likely to reflect the aging profile of HILDA survey respondents and that further study is disproportionately higher among the younger cohorts (25-44 year olds) (see Figure 4).

Over time, Figure 5 shows that the composition of degrees completed has shifted. Among those who completed a degree in later years, compared to those who completed a degree in the earlier period, a higher percentage completed a Certificate III or IV, Diploma or Advanced Diploma as opposed to a lower-level degree (Certificate I or II or below). In all years, the most frequently completed degrees are Cert 3 or 4, Associate degrees, Diplomas and Advanced Diplomas.

⁵Confounders are features that both have an impact on the outcome and on the treatment. [Chernozhukov et al. \(2018\)](#) suggest including the union of features kept in the two structural equations (outcome on features and treatment on features). Here, we only include the features that predict the outcome equation because including features that are only predictive of the treatment can erroneously pick up instrumental variables (see [Pearl \(2012\)](#) for a discussion of this issue).

The predominance of Cert 3 or 4 degrees is common across gender. Although, Figure 6 shows the distribution of degrees is more heavily skewed towards these degrees for men than they are for women.

Figure 7 shows an increase in average earnings of around \$380 between 2002 and 2017. The proportion entering entrepreneurship decreased slightly by 2 percentage points over the period but has been volatile since 2009. For example, in 2016 entrepreneurship reached a high of 20% before falling to 5% in 2017. Earnings also show high volatility from 2008. This likely reflects the smaller sample sizes in the later years of the survey. In our main analysis we average the returns over time as the samples within each year are inadequate to draw inference about heterogeneity across time.

5 Method

We aim to estimate the causal impact of obtaining a new qualification. Our empirical challenge is a missing data one in the sense that we do not observe the counterfactual outcome for each person – what would have their income been if they had/had not obtained a new qualification?

We use capitalisation to denote random variables, where $Y \in \mathbb{R}^+$ is the outcome variable for income and $Y \in \{0, 1\}$ for entrepreneurship, $T \in \{0, 1\}$ is the binary treatment indicator, and $X \in \mathcal{X}$ are the conditioning variables (which can be a mix of continuous or categorical in type). Small case is used to denote realisations of these random variables, e.g. y , t and x , and we may use a subscript for an individual realisation, e.g. y_i for individual i from a sample of size n .

Under the potential outcomes framework of Imbens and Rubin (2015), $Y(0)$ and $Y(1)$ denote the outcomes we would have observed if treatment were set to zero ($T = 0$) or one ($T = 1$), respectively. In reality, we only observe the potential outcome that corresponds to the realised treatment,

$$Y = T \cdot Y(1) + (1 - T) \cdot Y(0). \quad (1)$$

The missing data problem (or the lack of counterfactuals) is especially problematic when the treated group is different from the control group in ways that also affect outcomes. Such selection issues mean that we cannot simply take the difference in the average of the non-missing values of $Y(0)$ and $Y(1)$.

To address the missing data problem, we turn to a range of ML-based techniques. Standard ML tools are purposed to predict, but our aim is to estimate the causal parameter. These are different aims, and so we have to adapt the ML tools. We may potentially bias our causal parameter of interest if we were to use the off-the-shelf tools. For example, if we were to select the important confounders using an ML model to predict the outcome Y , then we may undervalue the importance of variables that are highly correlated to the treatment T but only weakly predictive of Y (Chernozhukov et al., 2018).

We approach filling the missing data indirectly with three types of ML models that have been specially adapted to causal inference. They are: the T-Learner, Doubly Robust and Bayesian models.

Identification assumptions

To interpret the estimated parameter as a causal relationship, we require the following identification assumptions: Conditional independence ($Y(0)$ and $Y(1)$ are independent of T conditional on X); SUTVA: Stable Unit Treatment Value Assumption ($Y = Y(0) + T \cdot (Y(1) - Y(0))$); Overlap Assumption (no subpopulation defined by $X = x$ is entirely located in the treatment or control group); Exogeneity of features (the features included in the conditioning set are not affected by the treatment). With the strong ignorability and overlap assumptions in place, treatment effect estimation reduces to estimating two response surfaces, one for treatment and one for control.

5.1 T-Learner model

The first adaptation of ML models for causal estimation is the T-learner approach. We aim to measure the amount by which the response Y would differ between hypothetical worlds in which the treatment was set to $T = 1$ versus $T = 0$, and to estimate this across subpopulations defined by attributes X .

The T-learner is a two-step approach where the conditional mean functions defined in Equations (2) and (3) are estimated separately with any generic machine learning algorithm.

$$\mathbb{E}[Y|X=x, T=1] \approx \mu_1(x) \quad \text{and} \quad (2)$$

$$\mathbb{E}[Y|X=x, T=0] \approx \mu_0(x) \quad (3)$$

Machine learning methods are well suited to find generalizable predictive patterns, and we employ a range of model classes including linear (LASSO and Ridge) and non-linear

(Gradient Boosted Regression ([Friedman, 2001](#))). Once we obtain the two conditional mean functions, for each observation, we can predict the outcome under treatment and control by plugging each observation into both functions. Taking the difference between the two outcomes results in the Conditional Average Treatment Effect (CATE), and averaging yields the Average Treatment Effect (ATE),

$$\hat{ATE} = \frac{1}{n} \sum_{i=1}^n \mu_1(x_i) - \mu_0(x_i). \quad (4)$$

In practice, however, this indirect way of minimising the mean squared error for each separate function to proxy for the minimum mean squared error of the treatment effect can be problematic. See, for example, [Künzel et al. \(2019\)](#), [Kennedy \(2020\)](#) for settings when the T-learner is not the optimal choice. One potential estimation problem arises when there are fewer treated individuals than control individuals and the individual regression functions are non-smooth. In this instance the response surfaces can be difficult to estimate in isolation, and the T-learner does not exploit the shared information between treatment and control observations. For example, if X relates to Y in the same fashion for treated and control observations the T-learner cannot utilise this information. As a result, the estimate μ_1 tends to over smooth the function; in contrast, the estimate μ_0 regularises to a lesser degree because there are more control observations. This means a naïve plug-in estimator of the CATE that simply takes the difference between $\mu_1 - \mu_0$ will be a poor and overly complex estimator of the true difference. It will tend to overstate the presence of heterogeneous treatment effects. We turn to other ML models to address this potential problem.

5.2 Doubly Robust model

The second approach is the Doubly Robust learner (DR-learner). It is similar to the T-learner in that it separately models the treatment and control surfaces, but it uses additional information from a propensity score model. In this case the propensity score model is a machine learning classifier that attempts to estimate the treatment assignment process,

$$\mathbb{E}[T=1|X=x] = \mathbb{P}(T=1|X=x) \approx \rho(x), \quad (5)$$

where $\rho(x)$ as a probabilistic machine learning classifier. This allows information about the students' background, and the nature and complexity of their situation that may have led them to pursue further education to be incorporated into the model. Thus, the doubly robust approach can improve upon the T-learner approach because it can

reduce misspecification error either through a correctly specified propensity score model or through correctly specified outcome equations. Another feature of the Doubly Robust approach is that it places a higher weight on observations in the area where the relative count of treatment and control observations is more balanced (i.e. the area of overlap). This may allow better extrapolations of the predicted outcomes within the region of overlap. The ATE is estimated from three separate estimators,

$$\hat{ATE} = \frac{1}{n} \sum_{i=1}^n \left[\frac{t_i(y_i - \mu_1(x_i))}{\rho(x_i)} + \mu_1(x_i) \right] - \frac{1}{n} \sum_{i=1}^n \left[\frac{(1 - t_i)(y_i - \mu_0(x_i))}{1 - \rho(x_i)} + \mu_0(x_i) \right] \quad (6)$$

Previously, with the T-learner, we were just estimating $\mu_0(x)$ and $\mu_1(x)$. With the DR-learner, we augment $\mu_0(x)$ and $\mu_1(x)$. For example, for the treated observations, we augment $\mu_1(x)$ by multiplying the prediction error by the inverse propensity scores. This up-weights those who get treated but who are statistically similar to the control observations. We then apply this same augmentation to the $\mu_0(x)$ for the control observations.

5.3 Bayesian Models

The third approach is to use Bayesian models. We follow the general formulation presented by [Hahn, Murray and Carvalho \(2020\)](#) that suggests a predictive model of the following form,

$$\mathbb{E}[Y|X=x_i, T=t_i] \approx \mu_0(x_i, \rho(x_i)) + \tau(x_i) \cdot t_i, \quad (7)$$

where $\mathbb{E}[T = 1|X=x_i] \approx \rho(x_i)$ is the propensity score of individual i for the treatment. The component $\mu_0(x_i, \rho(x_i))$ is known as the ‘prognostic’ effect, and is the impact of the control variates, X , on the outcome without the treatment. Then we are left with $\tau(x_i)$, which is the individual treatment effect,

$$\begin{aligned} \mathbb{E}[Y|X=x_i, T=1] - \mathbb{E}[Y|X=x_i, T=0] &\approx [\mu_0(x_i, \rho(x_i)) + \tau(x_i)] - \mu_0(x_i, \rho(x_i)), \\ &= \tau(x_i). \end{aligned}$$

Average treatment effect is then just simply estimated as,

$$\hat{ATE} = \frac{1}{n} \sum_{i=1}^n \tau(x_i).$$

This approach allows us to place explicit and separate priors on the prognostic and treatment components of the models. This minimises bias in the form of regularisation

induced confounding (RIC) which is discussed in more detail in [Hahn et al. \(2018\)](#), [Hahn, Murray and Carvalho \(2020\)](#). It is also a very natural way to estimate heterogeneous treatment effects, since we can parameterise $\tau(x_i)$ directly as an additive effect on μ_0 , rather than having to separately parameterise control and treatment surfaces.

However, since the entrepreneurial outcome is binary we had to modify Equation (7) to have a sigmoid curve to keep the model likelihood in the range $[0, 1]$. This also ensures the model’s posterior distribution over parameters reflects the true range of the outcome. So our Bayesian model becomes,

$$\mathbb{E}[Y|X=x_i, T=t_i] \approx \sigma(\mu_0(x_i, \rho(x_i)) + \tau(x_i) \cdot t_i), \quad (8)$$

where $\sigma(\cdot)$ is a logistic sigmoid. Unfortunately this nonlinear sigmoid breaks the separability of the prognostic and treatment models when computing the ATE, which we now compute as,

$$\hat{ATE} = \frac{1}{n} \sum_{i=1}^n \sigma(\mu_0(x_i, \rho(x_i)) + \tau(x_i)) - \sigma(\mu_0(x_i, \rho(x_i))).$$

We explore three different model classes for μ_0 and τ . The first is a linear model for both prognostic and treatment models. The next uses a Gaussian process (GP) for the earnings outcome, and an approximate GP for the entrepreneurship outcome. Lastly we use Bayesian additive regression trees (BART) for the earnings outcome (only). We detail these models in the following sections.

Hierarchical Linear/Logistic Models

The first Bayesian model uses linear prognostic and treatment components,

$$\begin{aligned} \mu_0(x_i, \rho(x_i)) &= w_0 + w_x^\top x_i + w_\rho \rho(x_i), \\ \tau(x_i) &= w_t + w_{tx}^\top x_i. \end{aligned}$$

The propensity score, $\rho(x_i)$, is obtained from a logistic regression model. We also tested a gradient boosted classifier ([Friedman, 2001](#)) for this using five-fold nested cross validation. It did not seem to be significantly more performant than the logistic model on held-out log-loss score.

For the earnings outcome we use a Normal likelihood, and for the entrepreneurship outcome we use a Bernoulli likelihood with a logistic sigmoid, as in Equation (8). We place

hierarchical priors over all model weights, the forms of which we have detailed in Appendix E.

For model inference, we use the no U-turn MCMC sampler (Hoffman and Gelman, 2014) in the numpyro software package (Bingham et al., 2019, Phan, Pradhan and Jankowiak, 2019). We first burn in the Markov chain for 2,000 samples, then draw 1000 samples from the posterior parameters to approximate the ATE,

$$A\hat{T}E_{\text{earnings}} = \frac{1}{Sn} \sum_{s=1}^S \sum_{i=1}^n \tau^{(s)}(x_i), \quad (9)$$

$$A\hat{T}E_{\text{entrep.}} = \frac{1}{Sn} \sum_{s=1}^S \sum_{i=1}^n \sigma \left(\mu_0^{(s)}(x_i, \rho(x_i)) + \tau^{(s)}(x_i) \right) - \sigma \left(\mu_0^{(s)}(x_i, \rho(x_i)) \right). \quad (10)$$

where (s) denotes a sample from the posterior parameters has been used to construct a random realisation of the treatment model component, and $S = 1000$.

Gaussian Process Regression

Gaussian process (GP) regression can be viewed as a non-linear generalisation of Bayesian linear regression that makes use of the kernel trick (Williams and Rasmussen, 2006, Bishop, 2006). Another way of understanding a GP is that it parameterises a distribution over functions (response surfaces) directly using kernels, rather than model weights as is the case with Bayesian linear regression.

To implement the functional relationship in Equation (7) in a Gaussian process, we create the kernel function over $\langle x, t \rangle$ pairs,

$$k(\langle x_i, t_i \rangle, \langle x_j, t_j \rangle) = \nu_{\mu_0}^2 k_{\mu_0}(\langle x_i, \rho(x_i) \rangle, \langle x_j, \rho(x_j) \rangle) + t_i t_j \cdot [\nu_{\tau}^2 k_{\tau}(x_i, x_j) + \tau_0].$$

Here k_{μ_0} and k_{τ} are the prognostic and treatment kernels respectively, ν_{μ_0} and ν_{τ} allow us to scale the contribution of these kernels to the functional relationships learned, and τ_0 permits a constant treatment effect. We chose a Gaussian process with a Matérn kernel which can learn non-linear and interaction-style relationships between input features and the outcome. The kernel parameters are given in Appendix E.

The ATE is approximated as,

$$A\hat{T}E = \frac{1}{Sn} \sum_{s=1}^S \sum_{i=1}^n f_*^{(s)}(x_i, 1) - f_*^{(s)}(x_i, 0),$$

where $f_*^{(s)}(x_i, t)$ are samples from the Gaussian process posterior predictive distribution⁶ with kernel inputs $k_*(\langle x_i, t \rangle, \langle x_i, t \rangle)$, which is equivalent to sampling from the distribution over $\tau(\cdot)$. We use $S = 100$ samples.

Unfortunately inference is more difficult in a Gaussian process that models the binary entrepreneurship outcome in Equation (8), and so we must resort to approximation. There are numerous ways we can approximate a Gaussian process classifier, the most simple of which is to use the hierarchical Bernoulli-logistic model from the previous section but to transform the inputs, x_i , with a Nyström basis expansion (Williams and Seeger, 2000). We can then use the same Monte-Carlo inference procedure and ATE approximation in Equation (10), but with $x_i \rightarrow \phi(x_i)$ where ϕ is the Nyström transformation. We give more details of this transformation in Appendix E.

Bayesian Causal Forests

The last Bayesian model we use is the Bayesian causal forest introduced in Hahn, Murray and Carvalho (2020). Broadly it models the prognostic and treatment components as Bayesian additive regression trees (BART). We use the accelerated BART (XBART) implementation of this algorithm detailed in Krantsevich, He and Hahn (2022). BART (Chipman, George and McCulloch, 2010) has been shown to be an effective and easily applicable non-parametric regression technique that requires few assumptions in order to capture complex relationships that can otherwise confound effect estimation. ATE is estimated in the same way as for the linear model in Equation (9), but where the BART posterior is used for the treatment effect distribution. We give more information about this model in Appendix E. Unfortunately there exists no implementation of the Bayesian causal forest for binary outcomes, and so we only use this model for the earnings outcome.

5.4 Model selection and model evaluation

For the non-Bayesian models we separate the evaluation of the model class and estimation of the ATE and CATE parameters in two procedures. We evaluate the predictive capacity of each model class using nested cross-validation. The procedure is represented in Figure 1. Here, our aim is to compare the predictive performance of three model classes: LASSO, Ridge and Gradient Boosted Regression (GBR). Our second procedure is to estimate the ATE and CATE parameters. The procedure is represented in Figure 2. We use bootstrap sampling (with replacement) to generate uncertainty estimates for the parameters, which we obtain over several draws of the same model class, but with model parameter re-fitting.

⁶See Equations (2.22)-(2.24) of Williams and Rasmussen (2006).

Focusing on the first procedure, we apply nested cross-validation to evaluate which model class performs best. In a first step, as Figure 1 shows, we pre-process the full dataset (containing 3,400 variables) to generate a dataset with a smaller set of highly predictive features (containing 91 variables). We apply a supervised machine learning approach with a LASSO model to select our top 91 predictors of the outcome of interest using outcomes measured in 2006. Note that in our later estimations of the treatment effect, the outcome is measured in 2019. We implement this intermediary step in order to reduce the correlation between variables and eliminate redundant information.

We assume that the top 91⁷ features that are most predictive of the outcome in 2006 correlate with the features that would be most predictive of the outcome in 2019. By choosing to apply this pseudo-supervised ML approach on the same outcome variable, but measured at a different time point, we obtain a good indication of the features that are useful for a model to perform well. Improved model performance here will also mean that the selected features are likely to represent the important confounders. We have chosen 2006 to ensure there is no overlap with 2019 outcomes to avoid overfitting issues with subsequent models.⁸

Using the top 91 predictors, we then apply nested cross validation to evaluate the predictive capacity of each model class (LASSO, Ridge, GBR). First, we split the data into train and test folds with an 80-20 split. Within the 80 percent train fold we perform 5-fold cross-validation in order to train and evaluate the performance of each configuration of hyperparameters. We do this separately for the outcome surface using the treated observations and the outcome surface using the control observations. From this, we select the models with the best mean predictive scores. We then evaluate the predictive performance of the selected model on the holdout test.

We repeat this process ten times (10-outer scores) for each model class. This allows us to evaluate the performance based on the mean and standard deviation of these scores. Note that thus far, we have not evaluated any particular configuration of the model, rather the performance of the model class on random (without replacement) subsets of data. The

⁷We were aiming for approximately 100 features, and 91 was the closest we could get the LASSO estimator to select by changing the value of its regularisation strength.

⁸We do not compromise predictive performance when we use the selected subset of features as opposed to the full set of features. For example, the predictive performance from a Gradient Boosted Tree model that predicts earnings in 2006, using 5-fold nested cross-validation, is statistically similar between models that use the 91 feature set and the full, 3,400 feature set (with Root Mean-Squared Errors (RMSEs) of 484.251 and 482.286, respectively). This is a negligible loss in predictive performance. There is a slightly larger associated loss between the restricted and full feature sets from models predicting earnings in 2019 (RMSEs of 843.548 and 831.931, respectively), but this is still not statistically significant.

nested cross validation procedure protects us against overfitting when reporting predictive performance, as the model selection and validation happens on different data.

Table 5 shows that the GBR is the best performing model class. It yields the highest out-of-sample R-squared and the lowest MSE. This is true for both the outcome surfaces separately. As the DR-learner model relies on the same treatment and control outcome surfaces estimated in the T-learner, we do not repeat Table 5 for the DR results. A further component of the DR model, however, is the propensity score. Here, we implement a regularised logistic regression to predict the likelihood of being treated (to obtain a further degree). Specifically, we use cross validation to fit a Logistic regression and obtain the predictions from the original sample. The holdout performance of the fitted Logistic regression model yields an area under the ROC curve of 0.71.

Inference via bootstrapping

Once we have selected the best performing model class, we turn to the estimation of the parameters and their associated uncertainty. We use a bootstrapped validation procedure to capture the uncertainty arising from model hyperparameter selection in addition to that from estimating parameters of a fixed model from noisy, finite data.

A common approach to inference in the causal machine learning literature is to use cross-fitting (Chernozhukov et al., 2018) or sample splitting (Athey and Wager, 2019). However, sample sizes of survey-based data are often not large enough to split the dataset into separate train and test datasets for each model. A suitable alternate procedure is to use bootstrapping. Bootstrap resampling allows us to estimate variation in the point model parameter estimates. In this way, we side-step the need to rely on the assumption of asymptotic normality, and it is more efficient than sample splitting to generate standard errors. In our bootstrapping procedure, we ensure that the standard errors reflect the sources of uncertainty stemming from both the selection of the model and the estimation of the model. As a result, we generate standard errors that avoid any potential pre-test issues. Appendix C describes the bootstrapping procedure in detail.

Inference for the Bayesian models

The inference process for the Bayesian models is a little different since the hyper-parameters of the models are either fixed or selected automatically by the learning algorithm (maximum likelihood type-II or MCMC). Bayesian inference procedures tend to afford some protection against over-fitting since they are parsimonious when choosing posterior distributions over model parameters that vary from their prior distributions, which induces

a natural model complexity penalty⁹. As such, we use all the available data to learn the model posterior distributions, which we then sample from to form empirical estimates of the (C)ATE as outlined in the previous section.

6 Results

There are clear entrepreneurial and economic benefits to gaining an additional qualification in later life (25 years or older). The effects remain strong up to a decade-and-a-half after course completion. Table 2 shows that re-training lifts the chance an individual becomes self-employed by approximately 1-2 percentage points. This is consistently estimated across all the models. In Table 2, we start by comparing regression model results with a linear probability model (OLS) and that of generalised linear models (GLMs) that have a Logit and Probit link that specifically account for the binary nature of entrepreneurship. The linear probability model and the GLMs are S-learners, in that they simply take the treatment variable as a single input feature. We find the results based on the marginal effects (the difference in the predicted probabilities evaluated at $t_i = 1$ and 0 for all i with all other x_i unchanged) are consistent across these three models. This suggests the linearity-in-parameters assumption is valid. We use regression models for the T-learner and Doubly Robust analyses, and logistic transformations for the Bayesian models. The T-learner and Doubly Robust models are relatively similar in terms of magnitude (a 2 percentage point increase in the probability of entrepreneurship) and statistical significance. The S-Learners and Bayesian models estimate a slightly smaller effect size, although their confidence intervals include the values estimated in the former approaches.

For the outcome of earnings, Table 3 displays a gain of approximately \$55-110 per week in gross earnings across the approaches. In 2019, this was roughly 5-8 percent of the average gross weekly earnings of \$1256.20 for all Australian employees (ABS, 2019; 6345.0 Wage Price Index, Australia). The effect sizes from the GBR T-Learner model are smaller than that of the two linear models. GBR better captures non-linearities. For example, age is likely to exhibit a highly non-linear relationship with earnings in 2019. Those who were aged 46 or above in 2001 will be aged 65 or above in 2019. This means they are more likely to have retired by 2019 compared to those who were aged below 46 in 2001. As a result, we may expect a shift down in earnings at age 46.¹⁰ The Doubly Robust (DR)

⁹This point can be understood more thoroughly by examining the evidence lower bound in variational Bayesian inference, see Chapter 10 of Bishop (2006).

¹⁰Age fixed-effects alone are unlikely to capture the differential age effects across other variables such as across different occupations, or by gender, and earnings. The linear ML models include age fixed effects. However, they do not include interactions between age and other variables whereas GBR does

models estimate smaller effects compared to the T-learner models. Table 3 displays a gain of approximately \$62-69 per week in gross earnings across the DR approaches. The estimated effect sizes are statistically different from zero. The confidence intervals for the DR estimates also exclude the point estimates from the T-Learner approach.

The Bayesian models estimate similar sized effects to the DR models for the most part. However, they tend to have more uncertainty associated with their estimates. They all remain significant with the 95% confidence intervals remaining above \$0. The hierarchical linear model and the Gaussian process both estimate a gain of approximately \$61-\$63 per week in gross earnings, with the Gaussian process being more certain in its estimate. Interestingly, the Gaussian process prefers a much smoother and smaller treatment effect component compared to its prognostic component¹¹. The Bayesian causal forest estimates a slightly higher gain of \$84.50 per week in gross earnings, which is more inline with the GBR T-learner. This suggests that the tree ensemble methods may be able to more easily capture non-linear relationships than the other models.

Proportionate changes in earnings can be measured by taking the log of the earnings measures. In Appendix Figure 15, we see that the proportionate change in earnings was large at 50 percent. This is likely to be because of people entering the labour market as a result of the new qualification. We find that a new qualification increases the likelihood of employment by approximately 8 percent (see Figure 8).

Acquiring an additional qualification may increase earnings through a number of potential mechanisms. We find evidence that, in Figure 8 for example, it increases the chance that individuals move from being unemployed or out of the labour force to being employed. The increase in employment is approximately 8 percentage points and is statistically significant. We also find evidence pointing to workers switching occupations or industries. This suggests that further education in later life can support the economic goals of a larger workforce as well as a more mobile one.

7 Sub-group Analysis

Qualification advancements may not benefit individuals in the same way. In this section we analyse if there is heterogeneity in the treatment impacts. We use a data-driven

include them. To illustrate how GBR adequately captures non-linearities we re-estimated our results focusing on those who were aged 25-45 in 2001. This is the same as interacting a binary variable (for age 25-45) with every other feature in the model. In Appendix Figure 14, we see that the results across the models are now more similar than when we use the full sample.

¹¹The treatment kernel length scale is long, and the kernel has a small amplitude and offset ($l_\tau = 243$, $\sigma_\tau^2 = 0.0517^2$, and $\tau_0 = 0.0312^2$) whereas the prognostic kernel parameters stay relatively close to their initial settings ($l_{\mu_0} = 16$, and $\sigma_{\mu_0}^2 = 1.42^2$).

approach to select the sub-groups. Specifically, we identify the important variables for which we expect to see the largest changes in the treatment effects. This involves using a permutation importance procedure.

7.1 Permutation Importance Feature Selection Method

We use a permutation importance selection method ([Breiman, 2001](#), [Molnar, 2020](#)) to evaluate the relative importance of individual features. Our aim here is to understand where the heterogeneous treatment effects are most pronounced. In other words, we aim to identify the sub-groups for which the treatment effects differ most significantly. Figures 9 and 10 display the top ten features and a residual category for all the other features. In selecting the important features our objective is to understand how to partition the data by the treatment effects as opposed to predicting the outcomes themselves. Appendix D describes this procedure in detail.

For entry into entrepreneurship, interestingly, the two most important features are: Age first left home and Do fair share of child care. Together, these two factors explain 27% of the importance of all the variables. This points to behavioural and childhood factors – rather than current economic conditions – being the top determinants of entrepreneurial returns from retraining. Although, the other categories point to more income- and financial-wealth based factors as also being important.

For earnings, the features that are most important are: weekly gross wages on the main job and income- or wealth-related variables. Together, this class of income/wealth variables accounts for 40% of the importance of all variables. We focus on these selected features since our Nested CV approach pointed to the better predictive performance of the GBR model over the linear models. Other important features include those related to employment, including occupational status, employment expectations, and employment history. The demographic background of the individual, namely their age, is also important.

Figures 16 and 17 display the distribution of the MSE importance values across the bootstrap samples for the GBR model. They display the distributions for the top 3 features for the respective outcomes. This suggests that in some of the bootstrap samples, where the MSE is larger, the individual treatment effects from the permuted data differ greatly from the original individual treatment effects. We present the results for the DR (GBR) model, however, similar features are chosen from the T-learner models. Results are available upon request.

I think more explanation may be needed here? Add some of the permutation importance description from the appendix back in

Figures 12 and 11 show that there is heterogeneity in the treatment impacts. We display heterogeneity by features that are both commonly used in such analysis and which are considered most important according to the permutation procedure. For each feature, we divide the sample into two groups. For continuous variables, we take the median value and divide the sample into those who are above and below this median value.

Larger entrepreneurial gains are achieved when qualifications are acquired at older ages. Younger learners may be more likely to seek re-training in order to become employees. These age-heterogeneity results contrast with the outcome of earnings, where mature-age learners are less likely to gain.¹² This may point to the different aspirations of those mature-age learners who do and do not become entrepreneurs in the late-career stage.

Weekly personal income has a large impact on the effect size for the outcome of earnings but little impact on whether or not retraining lifts the chance of being an entrepreneur. For earnings, those with below median income in 2001 derive more benefits than those with above median income, possibly because high income earners hit an earnings ceiling. Weekly personal income and age are likely to be highly correlated, with older individuals tending to earn a higher personal income. We cannot say which variable is the main driver of the heterogeneous treatment effects and there may also be interaction effects between them.

We also investigate if there are heterogeneous treatment effects according to commonly used variables in Figures 12 and 11. Males are far more likely to reap entrepreneurial gains compared to females when they retrain; whereas females are slightly more likely to reap earnings gains. Those with resident children are more likely to reap entrepreneurship-rewards when they retrain, compared to those without resident children. For earnings, however, similar treatment effects apply to those with and without a resident children. The gaps in these effects widen when we consider parents who have younger children in the household, although the gap widens in their favour when we consider earnings gains and tends to work against them for entrepreneurship. This suggests that entrepreneurship in the late-career stage involves time-commitments that prevent relatively time-poor individuals (with higher caring responsibilities) from entry – and that retraining or an additional degree provide less buffer against barriers to entry for them – compared to those with less caring responsibilities.

¹²Younger people have had more time to accumulate returns and this result aligns with findings from previous studies (Polidano and Ryan, 2016, Dorsett, Lui and Weale, 2016, ?)

7.2 Sensitivity Analysis

We test the sensitivity of our results in two ways. First, we compare the ML models with traditional econometric models such as Fixed Effects, Difference-in-Difference and Ordinary Least Squares (OLS) with controls. Following from this, we replicate published work by Chesters (2015) and try to understand the sources of bias in their OLS model compared to the estimated ML models. Second, we recode the input features so they are measured within the 2 years before a person begins re-training (as opposed to 2001 for everyone). In this analysis, we also measure outcomes 4 years after retraining (as opposed to in 2019 for everyone). For both sensitivity analyses, we focus on the outcome of earnings. The Chesters (2015) paper does not look at entrepreneurship and the small samples of people who transition into entrepreneurship limit our ability to use it as an outcome in the latter analysis.

7.2.1 Comparing results between Machine Learning and Traditional Econometric Models

The ML models estimate smaller returns than the returns estimated in DD-FE or cross-sectional models (OLS with and without controls) where features have been selected based on theory or previous empirical learnings. For example, the ‘OLS Baseline model’ uses the features in models estimated in Chesters (2015). The DD-FE eliminates all selection effects that are fixed over time. Figure 13 displays the estimated returns from six different approaches.

A potential reason for the smaller results estimated in the ML models is that the additional features included, as well as the non-linear specifications of the features, more effectively account for selection into treatment. The smaller results suggest individuals positively select into further study i.e. the characteristics that lead one to complete further study are positively correlated to future earnings. Once we control for this upward selection bias, we thus estimate smaller returns to further education.

The smaller estimated results relative to the DD-FE model are likely to stem from the inclusion of key time-varying variables such as the ‘change in total gross income’ in the ML models, as well as other non-linear specifications. For example, the ML models allow the treatment effects to vary in a highly flexible fashion across different parts of the feature distributions rather than making linear extrapolations.

This points to a benefit of using ML models, compared to conventional models, because they can more effectively identify confounders. We show evidence of the types of con-

founders missed in conventional models in Table 4, as well as the direction of the bias stemming from their omission.

In addition, we show evidence that models which allow for more flexible functional-form specifications lead to differences in the ATE. Within our ML models, the GBR tree ensemble tended to perform better (in terms of the nested cv results) compared to the linear-based models. The former yielded a slightly smaller ATE compared to the LASSO and Ridge results, for example, and they were also consistent with results from the Bayesian Causal Forest.

7.2.2 Missing variables from the baseline model

As part of a replication exercise, we contrast the results from the ML model with published work using Ordinary Least Squares (OLS) and Fixed Effects models. We also contrast the features selected in the ML model with an approach that manually selects the variables as in the case of [Chesters \(2015\)](#). We call this the ‘baseline’ model.

As a descriptive exercise, Table 4 presents the features that were ‘missed’ by the baseline model. In the baseline model, we included features such as age, gender, state of residence, household weekly earnings, highest level of education attained, and current work schedule. This collection of variables have been informed by theory or previous empirical results.

The data-driven model identifies more salient variables compared to the baseline model. Additional variables include employment conditions such as work schedule, casual employment, firm size, tenure or years unemployed; financial measures such as weekly wage, investment income and mortgage debt; health measures such as limited vigorous activity and tobacco expenses; and work-life preferences related to working hours and child care.

We identify variables as missing from the baseline model if those variables explain the residual variation in the outcome. Specifically, we regress the residuals from the baseline models (without the treatment included) on the features included in the data-driven model and train a LASSO model to highlight the salient variables that were missed. The variables that are chosen are listed in Table 4. We also document how these variables are correlated to the outcome and to the treatment in order to give us a sense of the direction of the bias their omission may induce.

Most of the omitted variables bias in the OLS estimates is upwards.¹³ The upward bias is consistent with the ML-models estimating an economic return of further study that is significantly smaller than the return from an OLS model or a Difference-in-Difference -

¹³Exceptions include casual employment status, the presence of a past doctorate qualification, years unemployed, parental child care and dividend and business income.

Fixed Effects (DD-FE) model. In the DD-FE model, we use the same 5,441 individuals but they are followed over two waves: 2001 and 2019 (i.e. there are 10,882 person-wave observations). We control for individual and wave fixed-effects.

Figure 13 displays the estimated returns from six different models. The first three bars show significantly higher returns based on the OLS (no controls), OLS (with controls) and the DD-FE models compared to the last three bars, which are based on the ML models – Gradient Boosted Regression, Doubly Robust and Bayesian Causal Forest. We discuss these methods in more detail below.

It is important to highlight that our approach to identifying missing variables from the baseline model is a descriptive one. As previously mentioned, the ML algorithm randomly selects variables that are highly correlated thus we may have missed out on reporting the label of important variables omitted from the baseline model.

7.2.3 Feature Inputs and Outcomes Measured in a Narrower Time Frame Around Re-training

For sensitivity analysis, we repeated the T-learner estimations using feature inputs values taken from individuals two years before they began study. Thus, we examine if our main results are sensitive to changes in the mapping equations for the treatment and control outcome equations when features are measured closer to the event of study, compared to taking input values in 2001. We also measured outcomes four years after study began. This means that the timing between when the feature input values are measured, when a further degree commenced and was completed, as well as when the outcomes are measured, are all closer together. This necessarily leads us to estimate the short-term returns of obtaining a further degree.

Our results from the sensitivity analysis are similar to that of the main results. Specifically, the gains in gross earnings from a further degree in the sensitivity analysis are: \$74 per week (Ridge), \$117 per week (LASSO) and \$93 per week (GBR). The key takeaway from these results is that the average treatment effects in the main analysis are not sensitive to whether our features use 2001 as the input year or use the two years before study.

Furthermore, the main results are not sensitive to when outcomes are measured i.e. the returns measured four years after the start of a study spell are comparable to the returns averaged over 2 to 17 years after study completion. This may point to the fact that the returns to further study are accrued in the immediate years following the completion of the degree. It also suggests the returns may not atrophy over time, especially since

the majority of people who did complete a degree in the main analysis did so in the earlier years of the survey (Figure 5). Unfortunately, our sample sizes are not sufficient to explore heterogeneity in treatment effects by the year of completion.

The importance of employment-related features such as earnings (individual and household), wages, and hours worked are reiterated in the sensitivity analysis using the panel structure of the data. Namely, when we define our outcomes 4 years after the start of a study spell and where we define features two years before study started, we also see similar results to that of the main results. However, in Figure 18, it is clear that the ‘trend’ or ‘growth’ in the values of features such as individual earnings, hours worked and household income are also important. This finding of dynamic selection is echoed in the literature (Jacobson, LaLonde and Sullivan, 2005, Dynarski, Jacob and Kreisman, 2016, 2018).

In Figure 18, the feature mental health is also picked. This result may reflect the fact that the timing of the measurement of features, treatment and outcomes are all closer together compared to the main results. This means that mental health is an important factor in explaining the heterogeneity in relatively ‘short-term’ treatment effects.

8 Conclusions

Using a machine learning based methodology and data from the rich and representative Household Income and Labour Dynamics Australia survey we have shown that completing an additional degree later in life can increase the chance of becoming self-employed by 1-2 percentage points. It can also add \$55-110 (AUD, 2019) per week to an individual’s gross earnings, which represents roughly 5-8 percent of the weekly gross earning for the average worker in Australia. Our machine learning methodology has also uncovered sources of heterogeneity in this effect.

We make a significant methodological contribution to the economics literature by undertaking causal inference with machine learning techniques, and we show how these outputs can be interpreted and made into clear recommendations for practitioners.

Our methodology has allowed us to exploit the full set of background information on individuals from the HILDA survey, beginning with more than 3,400 variables, to control our analysis. We find that our automated feature selection method selects a set of controls/features that include those that have theoretical foundations and/or align with those chosen in past empirical studies. However, we also choose features that have been traditionally overlooked. These include variables such as household debt, wealth, hous-

ing, and geographic mobility variables. Other important predictors include the ages of both resident and non-resident children: non-resident children aged 15 or above matter and resident children aged 0-4 are important.

Qualification advancements do not benefit Australian workers in the same way. Those with lower weekly earnings appear to benefit more from later-life study than those with higher earnings. One possible reason is that ceiling effects limit the potential returns from additional education. However, for entrepreneurship, the reverse is true: those with lower personal and household earnings tend to benefit less from retraining in later life compared to those with higher earnings. This likely reflects the fact that entrepreneurship requires higher start-up costs, compared to being an employee, and thus higher starting incomes provide a boost to the returns to retraining. We also find that younger Australians (less than 45 years of age) benefit more than their older counterparts in terms of earnings gains from education. Similarly to our previous interpretation, a ceiling effect phenomenon may apply since age is highly correlated to weekly earnings. Again, by contrast, the reverse is true for entrepreneurship where it appears that older individuals are more likely to benefit from retraining.

Acquiring an additional qualification may increase earnings through a number of potential mechanisms. The first is that it lifts the chance an individual becomes an entrepreneur in later-life. The earnings gains here may arise because becoming self-employed extends one's career into later-retirement. We also find evidence that it increases the chance that individuals move from being unemployed or out of the labour force to being employed, which also aligns with the entrepreneurship channel for individuals who are made unemployed before retirement age. We also find evidence pointing to workers switching occupations or industries. This suggests that further education in later-life can support the economic goals of a larger workforce as well as a more mobile one.

9 Tables and Figures

Table 1: Summary Statistics

Variable label	Variable name	Mean	SD
Outcomes			
Annual Earnings individual in 2019	y_wscei	614.730	1044.717
Imputed wages			
Change in annual earnings between 2001 and 2019	y_dwscei	129.029	980.754
Entry into entrepreneurship	chsefin	0.055	0.228
Treatment Indicators			
Highest level of educ changed between 2001 and 2017	reduhl	0.097	0.296
Extra degree attained in 2002 to 2017	redufl	0.257	0.437
Extra degree Bachelor and/or above	bachab	0.072	0.259
Below bachelor	bbach	0.209	0.406
Technical degree	techdeg	0.151	0.358
Qualitative degree*	qualdeg	0.080	0.272
Covariates (features)			
Demographics			
Sex	hgsex	1.536	0.499
Section of State	hhsos	0.690	1.046
Age	hgage1	46.025	12.832
Age of youngest person in HH	hhyng	27.115	21.886
No. persons aged 0-4 years in HH	hh0_4	0.257	0.589
No. persons aged 10-14 years in HH	hh10_14	0.274	0.606
Age when first left home	fmagelh	21.502	11.230
Living circumstances	hgms	1.997	1.708
English fluency	hgeab	1.604	0.262
Unemployment rate in region	hhura	6.884	1.075
Education			
Highest year of school completed/attending	edhists	2.383	1.439
Bachelor degree (without honours) obtained	edqobd	0.211	0.330
Masters degree obtained	edqoms	0.041	0.160
Doctorate obtained	edqodc	0.011	0.085
No. qualifications unknown	edqunk	0.078	0.403
Employment			
Occupation	jbmo61	3.772	1.825
Years in paid work	ehtjbyr	21.963	11.907
Tenure with current employer	jbempt	8.505	7.369
Type of work schedule	jbmday	3.785	2.612
Current work schedule	jbmsch	2.255	1.819

Continued on next page

Continued from previous page

Variable label	Variable name	Mean	SD
Casual worker	jbcasab	1.797	0.291
Hours/week worked at home	jbmhrh	12.372	7.174
Hours/week travelling to and from work	lshrcom	3.052	3.716
Satisfaction with employment opportunities	losateo	6.693	2.557
Occupational status - current main job	jbmofs	50.177	19.199
No. persons employed at place of work	jbmwpss	3.746	1.961
Age intends to retire	rtiagel	345.709	230.208
Age retired/intends to retire	rtage	113.904	130.211
Prob. of losing job in next 12 months	jbmplaj	15.196	35.018
Prob. of accepting similar/better job	jbmppgj	59.585	26.196
Looked for work in last 4 weeks	jsl4wk	1.272	0.411
Years unemployed and looking for work	ehtujyr	0.464	1.647
Hours per week worked in last job	ujljhru	34.990	6.922
Industry of last job	ujljin1	9.373	1.822
Work preferences			
Total hours per week would choose to work	jbprhr	34.378	6.407
Importance of work situation to your life	loimpew	6.854	2.908
Childcare			
Child looks after self	chu_sf	0.128	0.144
Uses child care while at work	cpno	1.257	0.139
Parent provides child care	cpu_me	0.434	0.151
Work-family balance			
Do fair share of looking after children	pashare	2.411	0.671
Miss out on home/family activities	pawkmfh	3.904	1.069
Working makes me a better parent	pawkbp	4.038	0.979
Family			
No. dependent children aged 5-9	hhd5_9	0.261	0.584
No. dependent children aged 10-14	hhd1014	0.269	0.604
No. non-resident children	tcnr	0.993	1.373
Sex of non-resident child	ncsex1	1.509	0.320
Likely to have a child in the future	icprob	1.188	0.374
Finances			
Owned a home previously	hspown	1.368	0.424
Amount outstanding on home loans	hsmgowe	96803.720	43547.610
Time until home loan paid off	hsmgfin	2011.858	4.157
Food expenses outside the home	xposml	36.982	42.522
SEIFA (level of economic resources)	hhec10	5.463	2.897
Taxes on total income	txtottp	7476.727	14035.510
Change in total gross income since 1 year ago	wslya	2231.465	1950.065
Had an incorporated business	bifinc	1.715	0.199
Had a non-LLC or unincorporated business	bifuinc	1.259	0.193
Income			

Continued on next page

Continued from previous page

Variable label	Variable name	Mean	SD
HH current weekly gross wages - all jobs	hiwscei	992.666	918.261
Current weekly gross wages - main job	wscme	468.062	556.185
HH financial year gross wages	hiwsfei	52472.490	49458.180
Financial year gross wages	wsfe	25463.770	30265.630
Financial year regular market income	tifmktp	30734.790	33618.860
Financial year disposable total income	tifditp	27477.160	22701.270
Imputation flag: current weekly gross wages - all jobs	wscef	0.070	0.256
Imputation flag: current weekly gross wages - other jobs	wscoef	0.044	0.205
Imputation flag: financial year gross wages	wsfef	0.071	0.256
Other sources of income			
Receive superannuation/annuity payments	oifsup	0.059	0.232
Receive redundancy and severance payments	oifrsv	0.002	0.038
Receive other irregular payment	oifirr	0.001	0.027
Receive government pensions or allowances	bncyth	0.004	0.027
Receive Disability Support Pension	bnfdsp	0.151	0.181
Receive other regular public payments	oifpub	0.000	0.019
Financial year regular private income	tifprin	77.299	1409.625
Financial year investments	oifinvp	1951.052	10569.050
Financial year dividends	oidvry	744.263	4651.593
Financial year interest	oiint	666.116	3448.494
Financial year regular private pensions	oifpp	967.101	5055.004
Financial year business income (loss)	bifn	185.652	3274.511
Financial year business income (profit)	bifip	2597.792	13649.410
Financial year irregular transfers from non-resident parents	oifnpt	35.067	1305.812
Financial year public transfers	bnfapt	2865.540	4717.042
Financial year government non-income support payments	bnfnis	1025.031	2237.987
HH financial year public transfers	hifapti	5542.675	7937.136
HH financial year business income	hibifip	4880.589	18393.360
Health			
Imputation flag: current weekly public transfers	bncapuf	0.044	0.204
Imputation flag: financial year investments	oifinf	0.124	0.330
Imputation flag: financial year dividends	oidvryf	0.079	0.270
Imputation flag: financial year rental income	oirntf	0.071	0.257
Imputation flag: financial year business income	biff	0.071	0.258
Health limits vigorous activities	gh3a	2.108	0.718
How much pain interfered with normal work	gh8	1.704	0.971

Continued on next page

Continued from previous page

Variable label	Variable name	Mean	SD
Health condition/disability developed last 12 months	helthyr	1.870	0.151
Tobacco expense in average week	lstbca	37.771	10.690
Housing			
Years at current address	hsyrcad	9.541	10.226
External condition of dwelling	docond	1.970	0.870
No dwelling security	dosecno	0.552	0.497
No. homes lived in last 10 years	mhn10yr	3.456	1.107
Moved to be near place of work	mhreawp	0.084	0.111
Moved because I was travelling	mhrearo	0.009	0.038
Attitudes			
Importance of religion	loimprl	4.612	3.483
Working mothers care more about work success	atwkwms	3.729	1.807
Mothers who don't need money shouldn't work	atwkmsw	3.951	1.982
Identifiers			
Family number person 02	hhfam02	NA	NA
Relationship to person 03	rg03	NA	NA
ID of other responder for HH Questionnaire	hhp2	NA	NA

*Definition of technical and qualitative degree: Technical: STEM, Architecture, Agriculture and Environment, Medicine, Other Health-related Studies and Nursing, Management and Commerce and Law. Non-technical: Education, Society and Culture (includes economics!), Creative Arts, and Food, Hospitality and Personal Services.

Table 2: Average Treatment Effects: Entry into Entrepreneurship. Comparison across models.

Model	N	ATE	CI (ATE)
OLS (S-learner)	5441	0.016	[0.001, 0.030]
Probit (S-learner)	5441	0.013	[0.000, 0.026]
Logit (S-learner)	5441	0.012	[-0.001, 0.025]
T-learner (GBR)	5441	0.024	[0.010, 0.037]
T-learner (LASSO)	5441	0.024	[0.012, 0.037]
T-learner (Ridge)	5441	0.022	[0.008, 0.036]
Doubly Robust (GBR)	5441	0.019	[0.016, 0.021]
Doubly Robust (LASSO)	5441	0.019	[0.018, 0.021]
Doubly Robust (Ridge)	5441	0.018	[0.017, 0.020]
Hierarchical Logistic Model	5441	0.014	[0.000, 0.030]
Approximate Gaussian Process	5441	0.013	[-0.001, 0.027]

Notes: Sample of 25 or older respondents who had completed a degree at any point between 2002 and 2017. Total completions: 1,383.

Table 3: Average Treatment Effects: Level Earnings. Comparison across models.

Model	N	ATE	CI (ATE)
OLS (S-learner)	5441	64.41	[8.16, 120.66]
T-learner (GBR)	5441	88.38	[30.72, 137.15]
T-learner (LASSO)	5441	110.08	[4.01, 182.49]
T-learner (Ridge)	5441	108.95	[46.84, 183.05]
Doubly Robust (GBR)	5441	68.85	[50.91, 82.07]
Doubly Robust (LASSO)	5441	54.64	[27.97, 72.74]
Doubly Robust (Ridge)	5441	61.74	[45.7, 78.86]
Hierarchical Linear Model	5441	63.22	[0.63, 121.70]
Gaussian Process	5441	61.01	[12.63, 109.51]
Bayesian Causal Forests	5441	84.51	[26.28, 141.17]

Notes: Sample of 25 or older respondents who had completed a degree at any point between 2002 and 2017. Total completions: 1,383.

Figure 1: Selecting and Evaluating Model Class

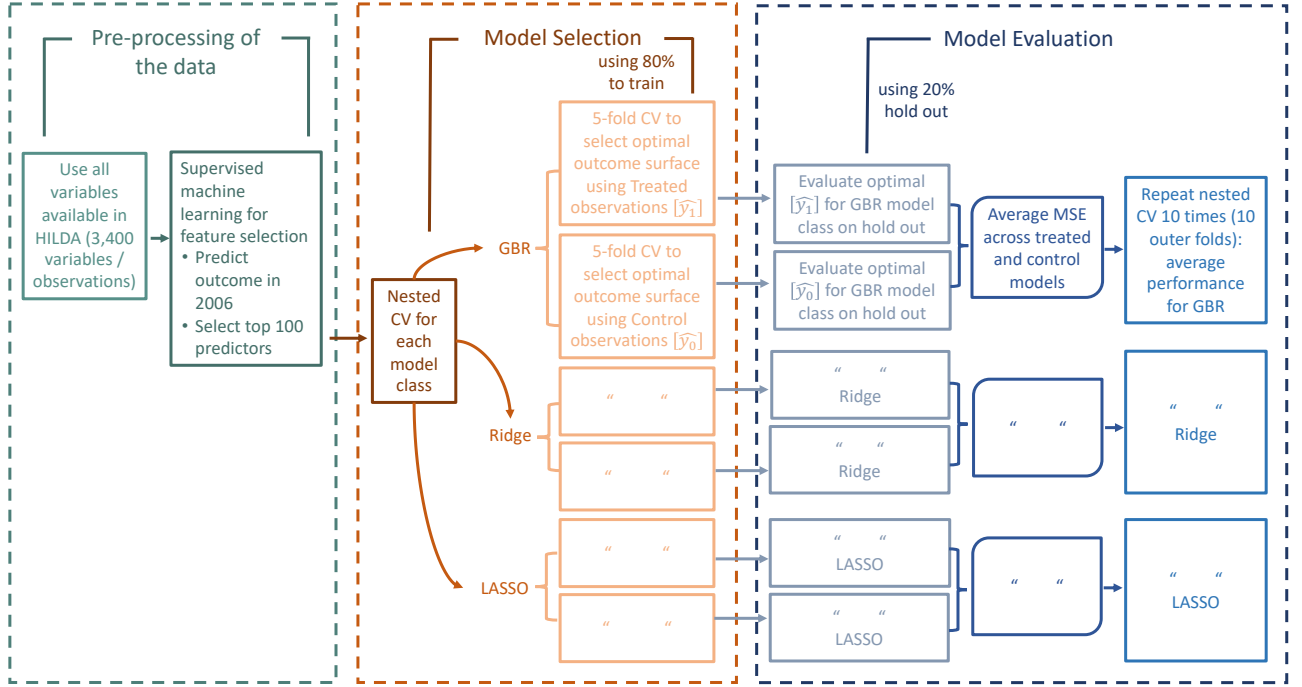


Figure 2: Generating Uncertainty Parameters

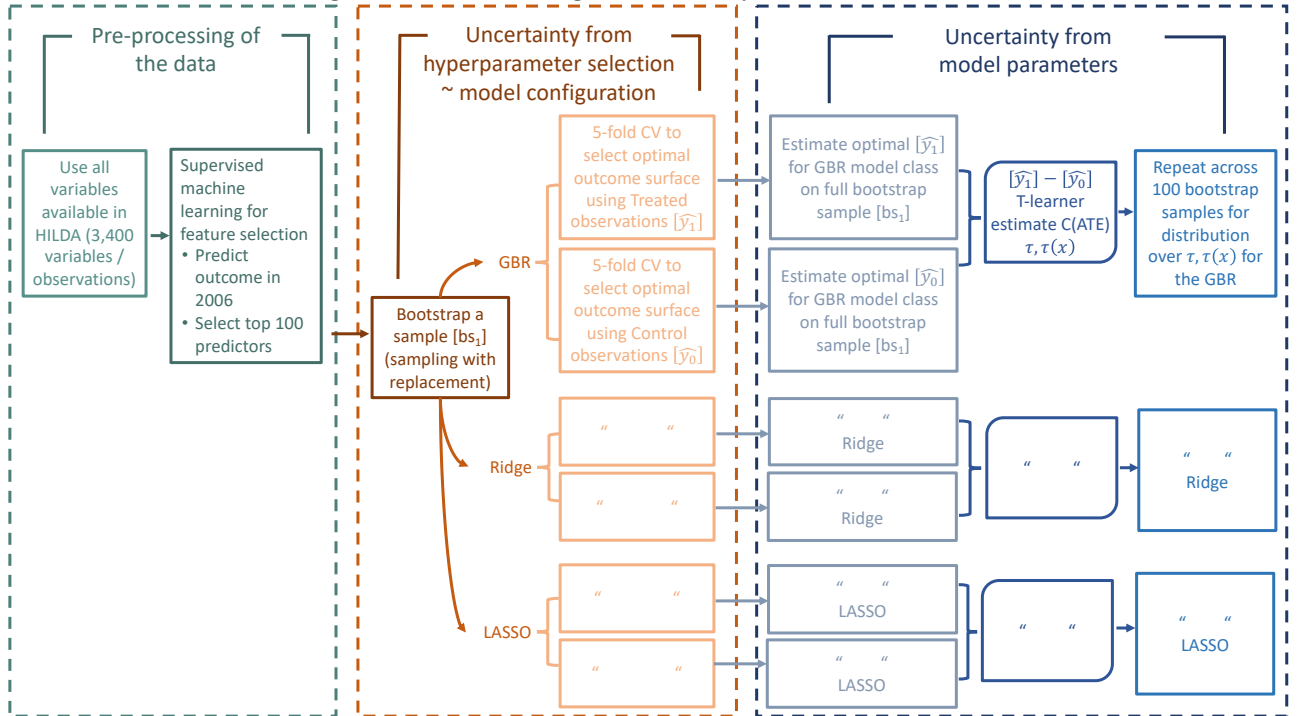
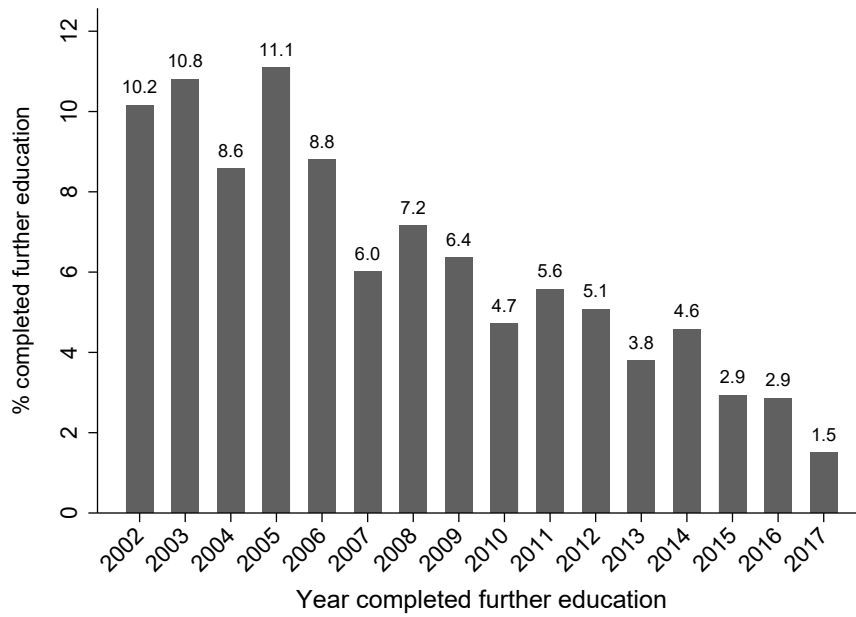


Figure 3: Timing of Completion



Notes: Sample of 25 or older respondents who had completed a degree at any point between 2002 and 2017. Total completions: 1,383.

Figure 4: Degree completions by age

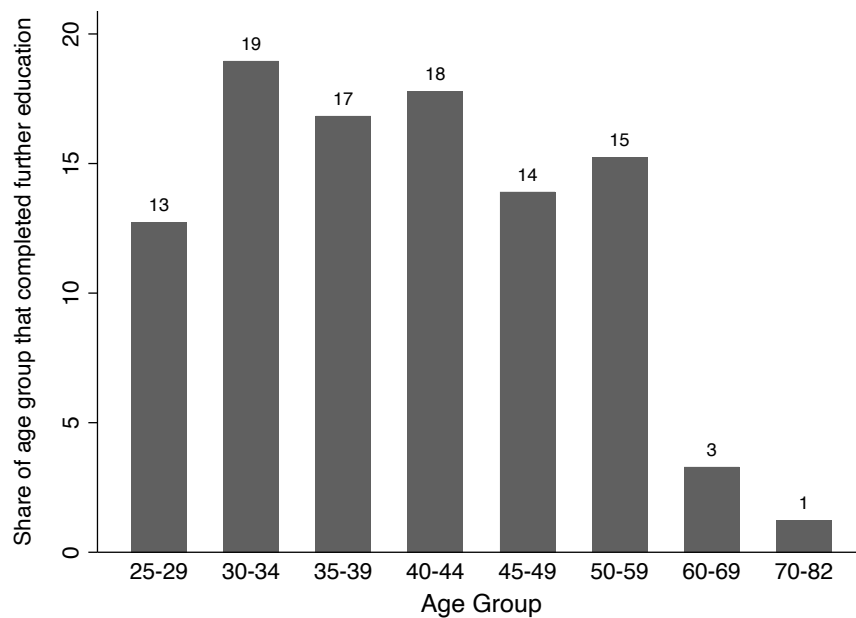
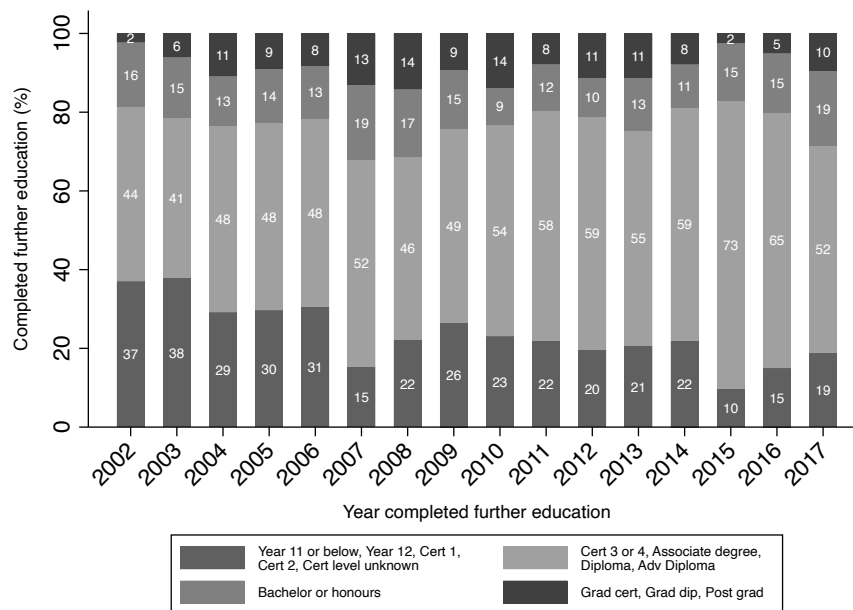


Figure 5: Timing of Completion by Type of Degree



Notes: Sample of 25 or older respondents who had completed a degree at any point between 2002 and 2017. Total completions: 1,383.

Figure 6: Degree completions by sex

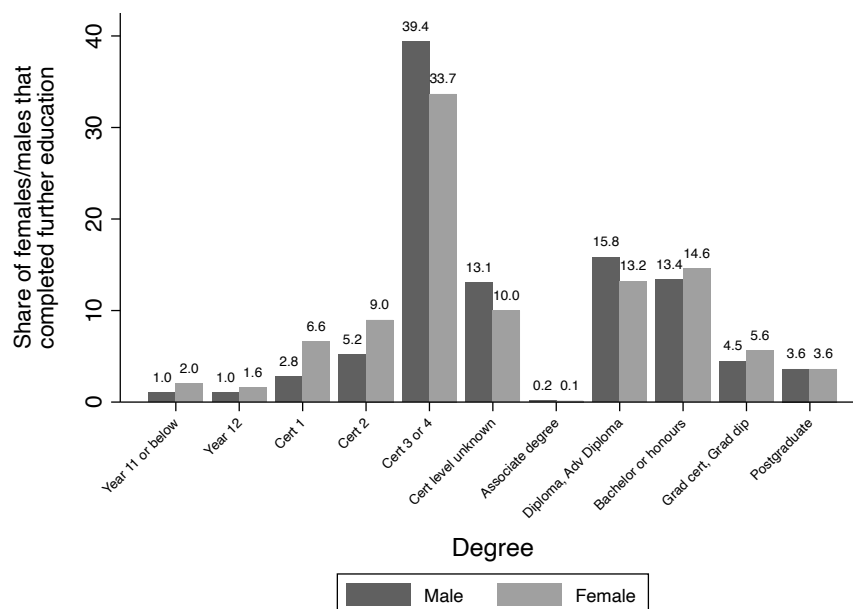


Figure 7: Earnings and Entrepreneurship by year

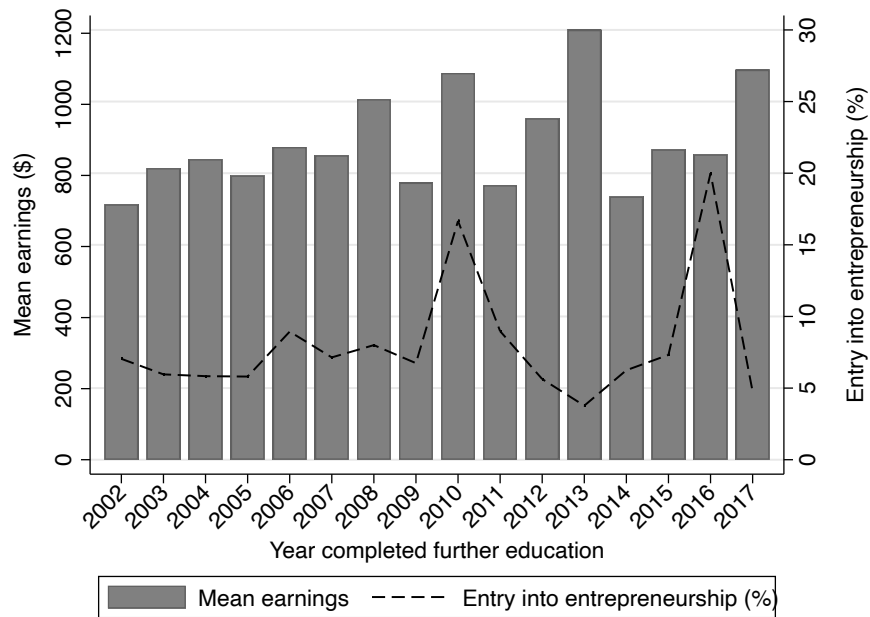
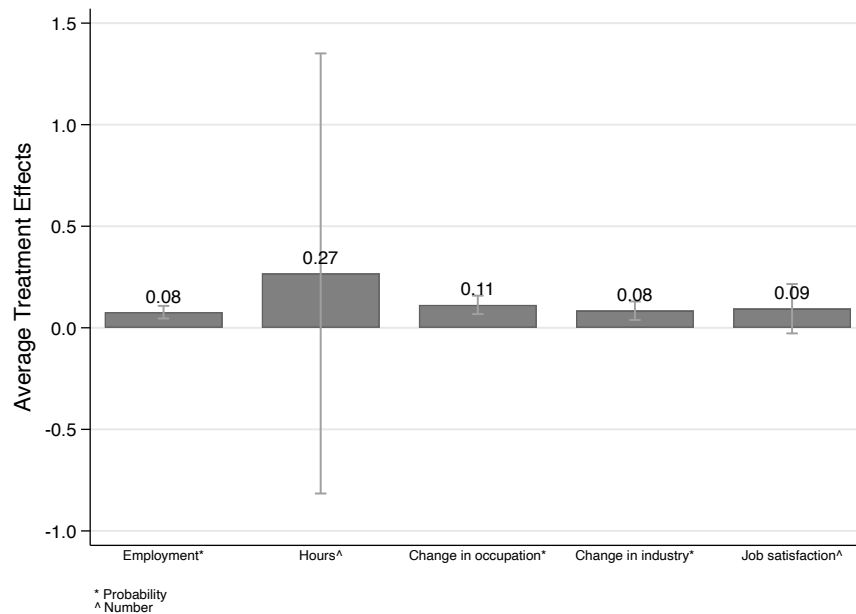
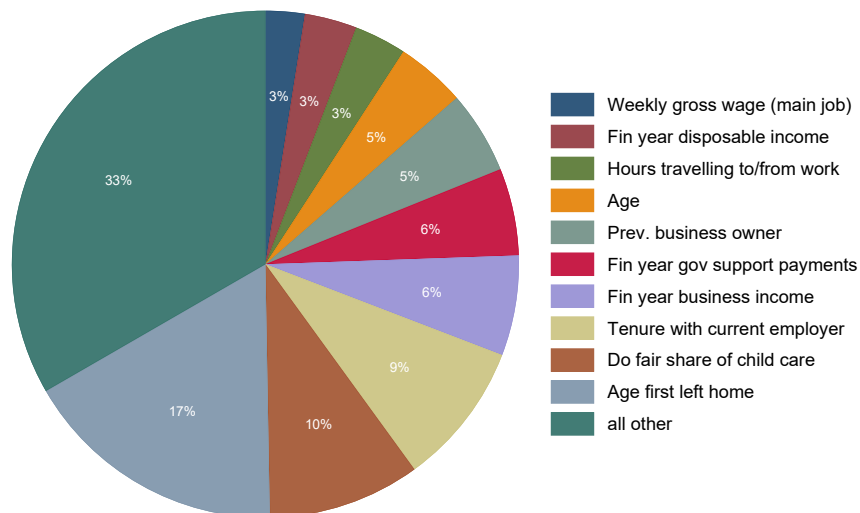


Figure 8: Other Employment Outcomes



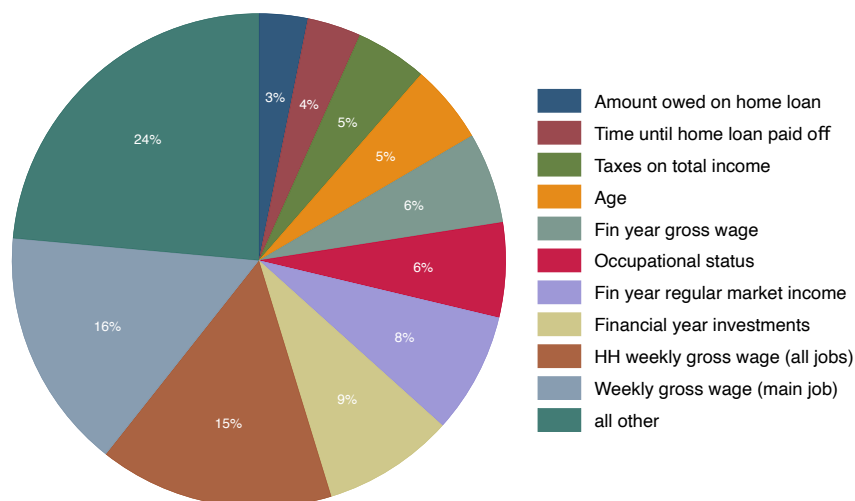
Notes: The impact of a new qualification. Sample of people who are 25 or older in 2001. Observation sizes vary depending on the outcome variable. All results are estimated using the LASSO algorithm.

Figure 9: Important Features in Heterogeneous Treatment Effects Estimation using DR: Entry into Entrepreneurship



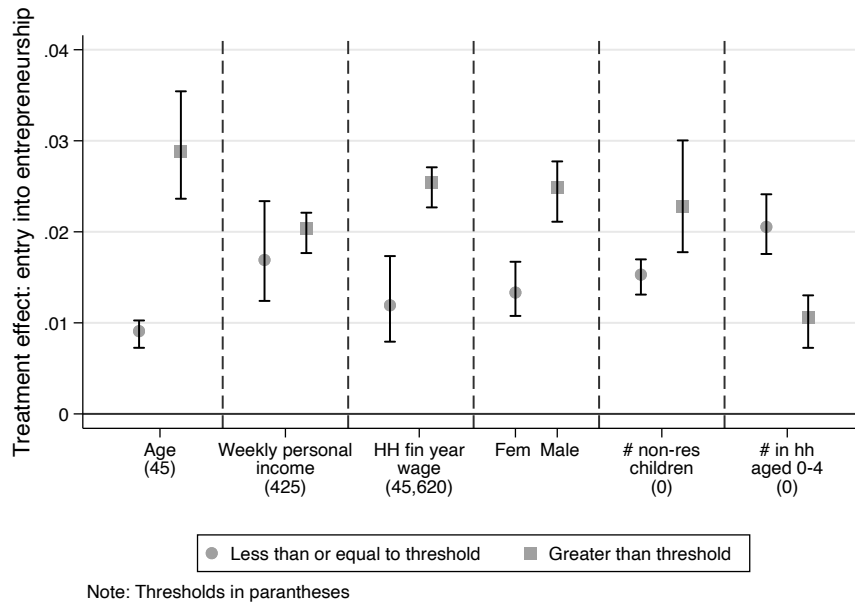
Notes: Sample of 25 or older who had completed a degree at any point between 2002 and 2017. Total number of observations 5,441.

Figure 10: Important Features in Heterogeneous Treatment Effects Estimation using DR: Level Earnings



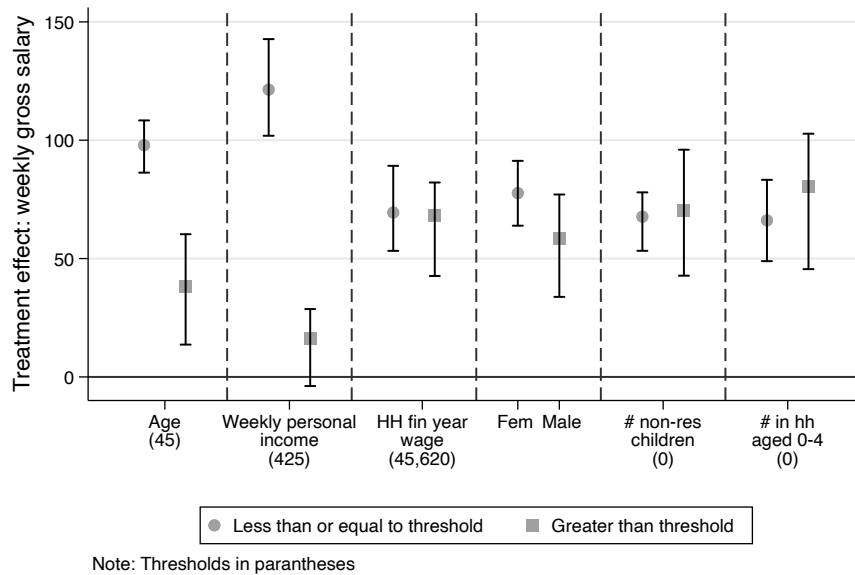
Notes: Sample of 25 or older who had completed a degree at any point between 2002 and 2017. Total number of observations 5,441.

Figure 11: Entrepreneurship HTEs: DR



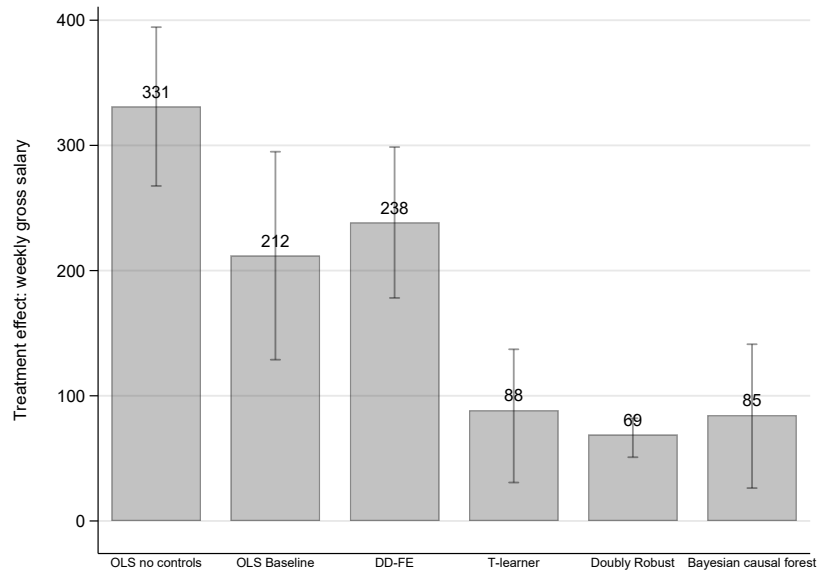
Notes: Sample of 25 or older who had completed a degree at any point between 2002 and 2017. Total number of observations 5,441.

Figure 12: Earnings HTEs: DR



Notes: Sample of 25 or older who had completed a degree at any point between 2002 and 2017. Total number of observations 5,441.

Figure 13: Comparison of Treatment Effects across Different Methods



Notes: Unless stated otherwise, the method uses a sample of 25 or older respondents who had completed a degree at any point between 2002 and 2017. Total number of observations 5,441. The OLS Baseline model uses the features manually selected in models by Chesters (2015). The Difference-in-Difference Fixed Effects (DD-FE) model uses the same individuals as the other methods but follows them over two waves: 2001 and 2019 (i.e. there are 10,882 person-wave observations); person and wave fixed effects included. The T-learner and Doubly Robust results are based on the Gradient Boosted Regression. The last bar is based on the Bayesian Causal Forest.

References

- Abel, Jaison R., and Richard Deitz.** 2014. “Do the Benefits of College Still Outweigh the Costs?” *Current Issues in Economics and Finance*, 20(3).
- Acemoglu, Daron, and David Autor.** 2011. “Skills, tasks and technologies: Implications for employment and earnings.” In *Handbook of Labor Economics*. Vol. 4, 1043–1171. Elsevier.
- Astebro, Thomas, and Jing Chen.** 2014. “The Entrepreneurial Earnings Puzzle: Mismeasurement or Real?” *Journal of Business Venturing*, 29(1): 88–105.
- Athey, Susan, and Guido W Imbens.** 2017. “The state of applied econometrics: Causality and policy evaluation.” *Journal of Economic Perspectives*, 31(2): 3–32.
- Athey, Susan, and Stefan Wager.** 2019. “Estimating treatment effects with causal forests: An application.” *Observational Studies*, 5(2): 37–51.
- Atkinson, Georgina, and John Stanwick.** 2016. “Trends in VET: Policy and participation.” *Aelaide: NCVER*.
- Autor, David H, Lawrence F Katz, and Melissa S Kearney.** 2008. “Trends in US wage inequality: Revising the revisionists.” *Review of Economics and Statistics*, 90(2): 300–323.
- Becker, Gary S.** 1964. *Human Capital: A Theoretical and Empirical Analysis, with Special Reference to Education*. New York:National Bureau of Economic Research.
- Belfield, Clive, and Thomas Bailey.** 2017. “The labor market returns to sub-baccalaureate college: A review. A CAPSEE working paper.” *Center for Analysis of Postsecondary Education and Employment*.
- Bingham, Eli, Jonathan P Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall, and Noah D Goodman.** 2019. “Pyro: Deep universal probabilistic programming.” *Journal of Machine Learning Research*, 20(1): 973–978.
- Bishop, Christopher M.** 2006. *Pattern recognition and machine learning*. Vol. 4, Springer.
- Breiman, Leo.** 2001. “Random forests.” *Machine Learning*, 45(1): 5–32.
- Caruso, Stephanie.** 2018. “The changing face of a student: Returning to education at a mature age in Australia.” <https://www.shortcourses.com.au/ed/studying-as-a-mature-age-student/>.
- Cassar, Gavin.** 2006. “Entrepreneur Opportunity Costs and Intended Venture Growth.” *Journal of Business Venturing*, 21(5): 610–632.
- Cherastidtham, Ittima, and Andrew Norton.** 2018. “University Attrition: What Helps and What Hinders University Completion?” Grattan Institute Report.
- Chernozhukov, Victor, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins.** 2018. “Double/debiased machine learning for treatment and structural parameters.” *Econometrics Journal*, 21(1).
- Chesters, Jenny.** 2015. “Within-generation social mobility in Australia: The effect of returning to education on occupational status and earnings.” *Journal of Sociology*, 51(2): 385–400.

- Chipman, Hugh A, Edward I George, and Robert E McCulloch.** 2010. “BART: Bayesian additive regression trees.”
- Coelli, Michael, Domenico Tabasso, and Rezida Zakirova.** 2012. *Studying beyond Age 25: Who does it and what do they gain? Research report.* ERIC.
- Deci, Edward L., and Richard M. Ryan.** 1985. “The General Causality Orientations Scale: Self-determination in Personality.” *Journal of Research in Personality*, 19(2): 109–134.
- Dorsett, Richard, Silvia Lui, and Martin Weale.** 2016. “The effect of lifelong learning on men’s wages.” *Empirical Economics*, 51(2): 737–762.
- Dynarski, Susan, Brian Jacob, and Daniel Kreisman.** 2016. “The fixed-effects model in returns to schooling and its application to community colleges: A methodological note.” *Center for Analysis of Postsecondary Education and Employment*.
- Dynarski, Susan, Brian Jacob, and Daniel Kreisman.** 2018. “How important are fixed effects and time trends in estimating returns to schooling? Evidence from a replication of Jacobson, Lalonde, and Sullivan, 2005.” *Journal of Applied Econometrics*, 33(7): 1098–1108.
- Efron, Bradley, and Robert Tibshirani.** 1986. “Bootstrap methods for standard errors, confidence intervals, and other measures of statistical accuracy.” *Statistical science*, 54–75.
- Fossen, Frank M, and Tobias JM Büttner.** 2013. “The returns to education for opportunity entrepreneurs, necessity entrepreneurs, and paid employees.” *Economics of Education Review*, 37: 66–84.
- Friedman, Jerome H.** 2001. “Greedy function approximation: a gradient boosting machine.” *Annals of Statistics*, 1189–1232.
- Guerra, Lenin C.** 2022. “Non-completion in Postsecondary Education: Why Are So Many Students Not Finishing Their Courses?”
- Haber, Sigal, and Arie Reichel.** 2007. “The Cumulative Nature of the Entrepreneurial Process: The Contribution of Human Capital, Planning and Environment Resources to Small Venture Performance.” *Journal of Business Venturing*, 22(1): 119–145.
- Hahn, P Richard, Carlos M Carvalho, David Puelz, and Jingyu He.** 2018. “Regularization and confounding in linear regression for treatment effect estimation.” *Bayesian Analysis*, 13(1): 163–182.
- Hahn, P Richard, Jared S Murray, and Carlos M Carvalho.** 2020. “Bayesian regression tree models for causal inference: Regularization, confounding, and heterogeneous effects (with discussion).” *Bayesian Analysis*, 15(3): 965–1056.
- Hanson, Melanie.** 2022. “College Dropout Rates.”
- Heckman, James J.** 1993. “What Has Been Learned About Labor Supply in the Past Twenty Years?” *The American Economic Review*, 83(2): 116–121.
- Herneas, Erik, Simen Markussen, John Piggott, and Knut Roed.** 2016. “Pension Reform and Labor Supply.” *Journal of Public Economics*, 142: 39–55.
- Hoffman, Matthew D, and Andrew Gelman.** 2014. “The No-U-Turn sampler: Adaptively setting path lengths in Hamiltonian Monte Carlo.” *Journal of Machine Learning Research*, 15: 1593–1623.

- Hundley, Greg.** 2001. “Why and When Are the Self-Employed More Satisfied with Their Work?” *Industrial Relations: A Journal of Economy and Society*, 40(2): 293–316.
- Imbens, Guido W, and Donald B Rubin.** 2015. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press.
- Jacobson, Louis, Robert LaLonde, and Daniel G Sullivan.** 2005. “Estimating the returns to community college schooling for displaced workers.” *Journal of Econometrics*, 125(1-2): 271–304.
- Kautonen, Teemu, Ewald Kibler, and Maria Minniti.** 2017. “Late-Career Entrepreneurship, Income and Quality of Life.” *Journal of Business Venturing*, 32(3): 318–333.
- Kennedy, Edward H.** 2020. “Optimal doubly robust estimation of heterogeneous causal effects.” *arXiv preprint arXiv:2004.14497*.
- Kibler, Ewald, Thomas Wainwright, Teemu Kautonen, and Robert Blackburn.** 2012. “(Work) Life after Work?: Older Entrepreneurship in London - Motivations and Barriers.” Small Business Research Centre Kingston University.
- Knaus, Michael C, Michael Lechner, and Anthony Strittmatter.** 2021. “Machine learning estimation of heterogeneous causal effects: Empirical monte carlo evidence.” *Econometrics Journal*, 24(1): 134–161.
- Krantsevich, Nikolay, Jingyu He, and P Richard Hahn.** 2022. “Stochastic tree ensembles for estimating heterogeneous effects.” *arXiv preprint arXiv:2209.06998*.
- Krishnan, Karthik, and Pinshuo Wang.** 2019. “The Cost of Financing Education: Can Student Debt Hinder Entrepreneurship?” *Management Science*, 65(10): 4522–4554.
- Künzel, Sören R, Jasjeet S Sekhon, Peter J Bickel, and Bin Yu.** 2019. “Metalearners for estimating heterogeneous treatment effects using machine learning.” *Proceedings of the National Academy of Sciences*, 116(10): 4156–4165.
- Levine, Ross, and Yona Rubinstein.** 2017. “Smart and Illicit: Who Becomes an Entrepreneur and Do They Earn More?” *The Quarterly Journal of Economics*, 132(2): 963–1018.
- Lombard, Karen V.** 2001. “Female Self-employment and Demand for Flexible, Non-standard Work Schedules.” *Economic Inquiry*, 39(2): 214–237.
- Molnar, Christoph.** 2020. *Interpretable machine learning*. <https://christophm.github.io/interpretable-ml-book/>.
- Montgomery, Mark, and Irene Powell.** 2006. “The Effect of Tuition and Opportunity Cost on the Pursuit and Completion of a Graduate Management Degree.” *Journal of Education for Business*, 81(4): 190–200.
- Mountjoy, Jack.** 2022. “Community colleges and upward mobility.” *American Economic Review*, 112(8): 2580–2630.
- NCVER DataBuilder.** 2021. “Total VET students and courses 2020: program enrolments.” *Department of Education, Skills and Employment*. <https://www.ncver.edu.au/research-and-statistics/data/databuilder>.

- OECD.** 2016. “Indicator C1: Who participates in education?” *Education at a Glance*.
- Oosterbeek, Hessel, Mirjam Van Praag, and Auke Ijsselstein.** 2010. “The Impact of Entrepreneurship Education on Entrepreneurship Skills and Motivation.” *European Economic Review*, 54(3): 442–454.
- O’Shea, S, J May, and C Stone.** 2015. “Breaking the barriers: Supporting and engaging mature age first-in-family university learners and their families (Final Report).”
- Parker, Simon C., and Jonathan C. Rougier.** 2007. “The Retirement Behaviour of the Self-employed in Britain.” *Applied Economics*, 39(6): 697–713.
- Pearl, Judea.** 2012. “On a class of bias-amplifying variables that endanger effect estimates.” *arXiv preprint arXiv:1203.3503*.
- Perales, Francisco, and Jenny Chesters.** 2017. “The returns to mature-age education in Australia.” *International Journal of Educational Research*, 85: 87–98.
- Phan, Du, Neeraj Pradhan, and Martin Jankowiak.** 2019. “Composable Effects for Flexible and Accelerated Probabilistic Programming in NumPyro.” *arXiv preprint arXiv:1912.11554*.
- Polidano, Cain, and Chris Ryan.** 2016. “Long-term outcomes from Australian vocational education.” *Melbourne Institute of Applied Economic and Social Research, University of Melbourne*.
- Stephan, Ute, and Ulrike Roesler.** 2010. “Health of Entrepreneurs versus Employees in a National Representative Sample.” *Journal of Occupational and Organizational Psychology*, 83(3): 717–738.
- Studies in Australia.** 2018. “Study Costs.” <https://www.studiesinaustralia.com/studying-in-australia/how-to-study-in-australia/study-costs>.
- Unger, Jens M., Andreas Rauch, Michael Frese, and Nina Rosenbusch.** 2011. “Human Capital and Entrepreneurial Success: A Meta-analytical Review.” *Journal of Business Venturing*, 26(3): 341–358.
- Universities Australia.** 2019. “2019 Higher Education Facts and Figures.” <https://www.universitiesaustralia.edu.au/wp-content/uploads/2019/08/190716-Facts-and-Figures-2019-Final-v2.pdf>.
- Universities Australia.** 2020. “2020 Higher Education Facts and Figures.” <https://www.universitiesaustralia.edu.au/wp-content/uploads/2020/11/200917-HE-Facts-and-Figures-2020.pdf>.
- Weber, Paull, and Michael Schaper.** 2004. “Understanding the Grey Entrepreneur.” *Journal of Enterprising Culture*, 12(02): 147–164.
- Williams, Christopher, and Matthias Seeger.** 2000. “Using the Nyström method to speed up kernel machines.” *Advances in neural information processing systems*, 13.
- Williams, Christopher KI, and Carl Edward Rasmussen.** 2006. *Gaussian processes for machine learning*. Vol. 2, MIT press Cambridge, MA.

- Xu, Di, and Madeline Trimble.** 2016. "What about certificates? Evidence on the labor market returns to nondegree community college awards in two states." *Educational Evaluation and Policy Analysis*, 38(2): 272–292.
- Zacharakis, Andrew L., and G. Dale Meyer.** 2000. "The Potential of Actuarial Decision Models: Can They Improve the Venture Capital Investment Decision?" *Journal of Business Venturing*, 15(4): 323–346.
- Zeidenberg, Matthew, Marc Scott, and Clive Belfield.** 2015. "What about the non-completers? The labor market returns to progress in community college." *Economics of Education Review*, 49: 142–156.
- Zhao, H., G. O'Connor, J. Wu, and G. T. Lumpkin.** 2021. "Age and Entrepreneurial Career Success: A Review and a Meta-analysis." *Journal of Business Venturing*, 36(1): 106007.

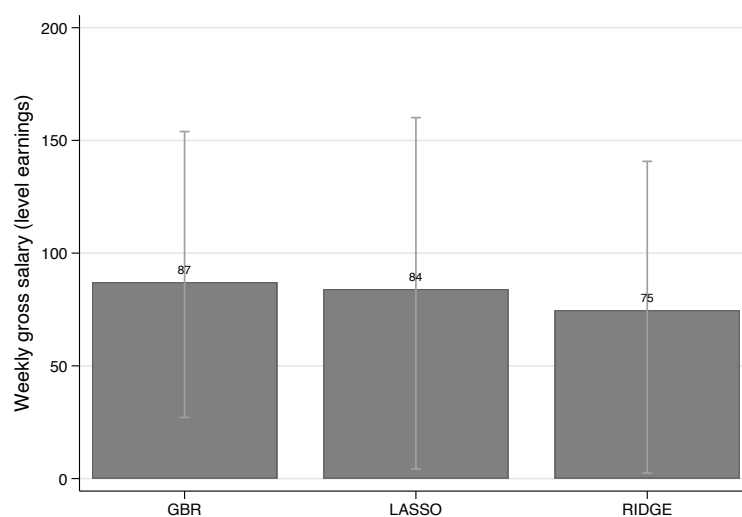
Appendix A: Tables

Table 4: ML variables omitted by OLS Baseline model

Variable label	Variable name	Relationship with retraining (redufl)	Relationship with outcome (y_wscei)	Bias direction in OLS models
<i>Education</i>				
Doctorate obtained	edqode	-	+	-
<i>Employment</i>				
Tenure with current employer	jbempt	-	-	+
Current work schedule	jbmsch	-	-	+
Casual worker	jbcasab	-	+	-
Occupational status - current main job	jbmo6s	+	+	+
No. persons employed at place of work	jbmwpsz	+	+	+
Prob. of accepting similar/better job	jbmpgj	+	+	+
Years unemployed and looking for work	ehtujyr	+	-	-
<i>Work-life balance</i>				
Hours per week would choose to work	jbprhr	+	+	+
Parent provides child care	cpu_me			-
Do fair share of looking after children	pashare	-	+	-
Miss out on home/family activities	pawkmfh	+	+	+
<i>Income</i>				
Current weekly gross wages - main job	wscme	+	+	+
Imputation flag: current weekly gross wages - all jobs	wscef	+	+	+
Change in total gross income since 1 year ago	wslya	+	+	+
Fin year investments	oifinvp	-	-	+
Fin year business income (profit)	bifip	-	-	+
Amount outstanding on home loans	hsmgowe	+	+	+
Imputation flag: fin year dividends	oidvryf	+	-	-
Imputation flag: fin year rental income	oirntf	+	+	+
Imputation flag: fin year business income	biff	+	-	-
<i>Health</i>				
Health limits vigorous activities	gh3a	+	+	+
Tobacco expense in average week	lstbca	-	-	+
<i>Identifiers</i>				
ID of other responder in HH qnn	hhp2	-	-	+

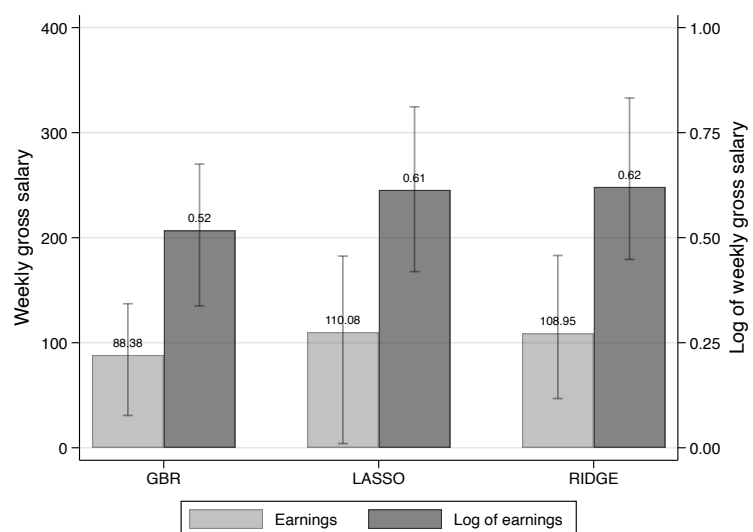
Appendix B: Figures

Figure 14: Value-add in earnings: 25-45 year-old sample



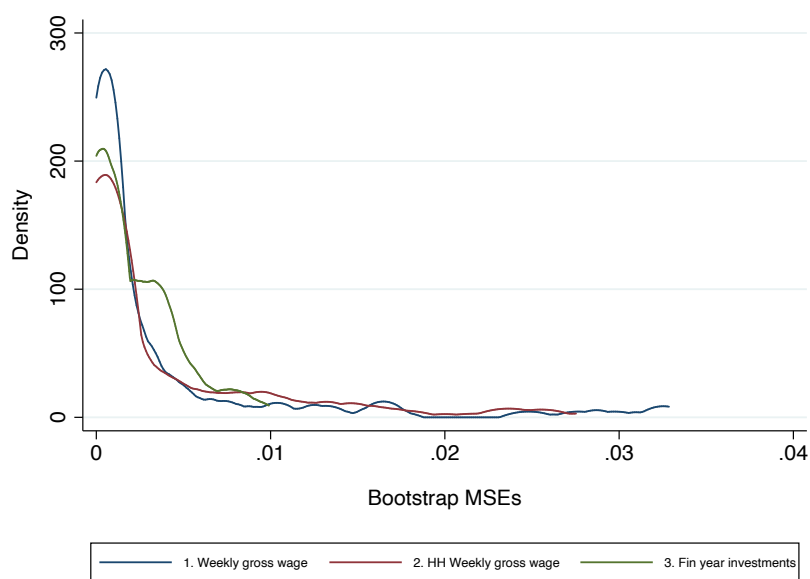
Notes: Sample of 25-45 who had completed a degree at any point between 2002 and 2017. Total number of observations 3,684.

Figure 15: Value-add in log earnings



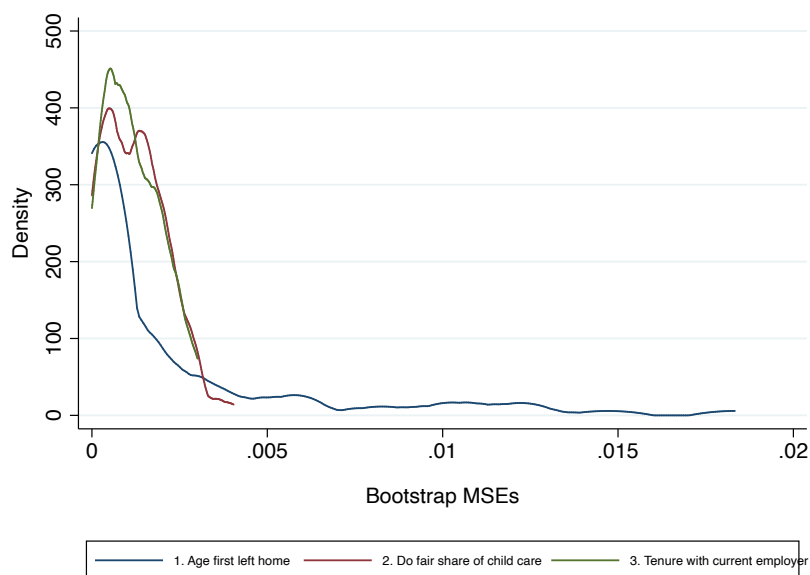
Notes: Sample of 25 or older who had completed a degree at any point between 2002 and 2017. Total number of observations 5,441.

Figure 16: Top 3 Features Distribution of Importance using DR (GBR): Level Earnings



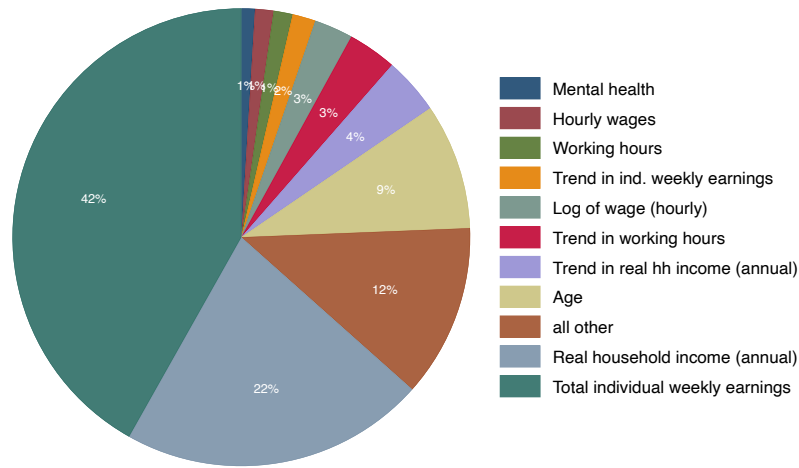
Notes: Sample of 25 or older who had completed a degree at any point between 2002 and 2017. Total number of observations 5,441.

Figure 17: Top 3 Features Distribution of Importance using DR (GBR): Entry into Entrepreneurship



Notes: Sample of 25 or older who had completed a degree at any point between 2002 and 2017. Total number of observations 5,441.

Figure 18: Important Features in Heterogeneous Treatment Effects Estimation using panel sample (GBR): Level Earnings



Notes: Sample of 21 or older individuals who had completed a degree at any point between 2003 and 2015, inclusive. Outcomes are defined 4 years after a study spell began and features are defined in both the two years preceding the start of a study spell. There were 1,814 individuals who started and completed a further educational degree, and 60,945 non-unique control observations who never completed a further degree.

Appendix C: Bootstrapping Procedure

To estimate the parameters and their associated uncertainty we use a bootstrapped validation procedure.

As a first step we obtain the 91 top predictors from the initial pre-processing of the full dataset, shown in Figure 2. That is, we train a supervised machine learning LASSO model to extract the features that best predict earnings in 2006.

The second step involves training our models using the 91 top predictors on a bootstrapped sample, s , to select the best models for $\mu_1^{(s)}(x)$ and $\mu_0^{(s)}(x)$.

Third, and once we have these predicted outcome surfaces, $\mu_1^{(s)}(x)$ and $\mu_0^{(s)}(x)$, we are able to calculate the individual treatment effect, $\tau(x_i)$, for each person, i , in the original sample (not the individuals from the bootstrap sample) by substituting the values of their features into the LASSO, Ridge or tree estimators for the outcome surfaces. We can obtain a sample mean, $\bar{\tau}^{(s)}$, by averaging $\frac{1}{n} \sum_{i=1}^n \tau^{(s)}(x_i)$ using the bootstrapped effect model. We repeat this procedure over $S = 100$ bootstrap samples. This provides an empirical distribution of $\bar{\tau}$ and $\tau(x_i)$. The grand mean over the bootstrap sample means, $\bar{\tau}_G = \frac{1}{S} \sum_{s=1}^S \bar{\tau}^{(s)}$, will converge to the sample treatment effect mean. We use $\bar{\tau}_G$ as an estimate of the ATE, and $\frac{1}{S} \sum_{s=1}^S \tau^{(s)}(x_i)$ as an estimate of the individual CATE. The bootstrap resample is the same size as the original sample because the variation of the ATE depends on the size of the sample. Thus, to approximate this variation we need to use resamples of the same size.

To obtain confidence intervals for the ATE and CATE estimates we use standard empirical bootstrap confidence interval estimators (Efron and Tibshirani, 1986).

For the DR-learners, similar to the T-learner, we train $\mu_1(x)$ and $\mu_0(x)$ models across 100 bootstrap samples and weight these outcome surfaces by the propensity score model, $\rho(x)$, which is estimated using logistic regression (as described previously).

Appendix D: Permutation Importance Procedure

The permutation importance procedure involves testing the performance of a model after permuting the order of samples of each individual feature, thereby keeping the underlying distribution of that feature intact but breaking the predictive relationship learned by the model with that feature. The model performance we are interested in is the one that maps the features to the individual treatment effects.

Following this approach, we compute the individual treatment effects. Note that we train the model on the bootstrapped sample but estimate the individual treatment effects using the feature values for individuals from the original sample. Thus, for every individual we have a distribution of values of their individual treatment effects.

After obtaining the individual treatment effects, we train another model that maps the features to the individual treatment effects. We use cross-validation to select our hyperparameters and obtain the optimal model.

Using the original data, we take a single column among the features and permute the order of the data and calculate a new set of individual treatment effects. We compare the new and original individual treatment effects (based on the permuted data and those from the non-permuted data) and calculate the Mean Squared Errors (MSE).

We repeat this for all the features, permuting them individually and evaluating how they change the prediction of the individual treatment effect target. Features that yield the largest MSEs are likely to be more important than those features with lower MSEs since permuting those features breaks the most informative predictive relationships.

We then repeat the above steps across all the bootstrap samples. Note that a different bootstrap sample will change the value of the individual treatment effects since we train different outcome surfaces for $\mu_0^{(s)}(x)$ and $\mu_1^{(s)}(x)$ for each bootstrap sample.

We embed the permutation importance selection method in a bootstrapping procedure in order to capture hyperparameter uncertainty. For example, a different ‘tree depth’ could be chosen between different bootstrap samples. This would affect the type of non-linear/interaction relationships that would be captured by the models, which in turn would affect which features turn out to be important.

Finally, we obtain an average MSE for each feature, averaged across all bootstrap samples. This average value allows us to rank the features by their importance. Again, those with the largest average MSE values are the most important. We can also evaluate the uncertainty of this estimate since we obtain a distribution of MSE values across the different bootstrap samples.

Appendix E: Bayesian Model Details

Online Appendix F: Main Sample

F.1 Sample Selection

Our analysis sample includes everyone who was 25 or above and not currently studying in 2001, who are observed in both 2001 and 2019 in terms of the outcome and treatment variables.

We delete any individuals who were currently studying in 2001 if:

- They reported currently studying full part or part time for the main survey
- According to the calendar, they have undertaken any full time or part time studies
- They are currently receiving Abstudy/Austudy payment or had received these last financial year
- They have cited study as the reason for not looking for work

1078 individuals were deleted after applying this sample exclusion.

F.2 Variable description

F.2.1 Outcome Variables

Weekly earnings from main job in 2019 (*w19_wscmei*) records the weekly earnings from the main job for the individual in 2019.

Employed in 2019 (*w19_employed*) records whether the individual is employed in 2019 or not.

Weekly earnings from all jobs in 2019 (*w19_earning*) records the weekly earnings from all jobs for the individual in 2019.

Working hours in 2019 (*w19_wkhr*) records the total number of hours the individual works in all jobs in a week on average. Working hours are set to 0 for those not working.

F.2.2 Treatment Variables: Retraining

Retraining based on highest attainment¹⁴ (*reduhl*) records whether the individual has had retraining between 2002 and 2017, based on whether there was a change in the highest education level attained stated in the two years.

Retraining completion based on detailed qualifications (*redudl*) records whether the individual has completed any one of the following qualifications since last interviewed between 2002 and 2017¹⁵:

- Trade certificate or apprenticeship

¹⁴The HILDA variable on highest attainment was constructed using three components: the age the individual left school, the highest education attainment in the previous wave and the current level of secondary school attained or currently studying for.

¹⁵Refer to https://en.wikipedia.org/wiki/Australian_Qualifications_Framework for how most of these degrees are situated relative to each other in a hierarchy and the duration of these qualifications.

- Technicians cert/Advanced certificate
- Teaching qualification
- Nursing qualification
- Associate Degree
- Advance Diploma (3 years full time or equivalent)
- Bachelor degree but not honours
- Certificate I
- Certificate II
- Certificate III
- Certificate IV
- Certificate of unknown level
- Doctorate
- Diploma NFI
- Diploma (2 years full time or equivalent)
- Graduate Certificate
- Graduate Diploma
- Honours
- Masters
- Other

Retraining completion based on both highest attainment and detailed qualifications (*reduft*) records whether the individual has completed retraining based on both the variables *reduhl* and *redudl*. When either of these variables has a value of 1, this variable will take on the value of 1.

F.2.3 Input Variables

For each variable, missing values (if any) have been set to zero and a new binary variable has been generated to indicate the observations that are missing.

Demographics

Female (*p_fem*) records whether the individual is female.

Age group in 2001 records whether in 2001 the individual was:

- Aged 25-34 (*p_age1*)
- Aged 35-44 (*p_age2*)
- Aged 45-54 (*p_age3*)
- Aged 55-64 (*p_age4*)

- Aged 65 and above (*p_age5*)

Country of birth records whether or not an individual was born in:

- Australia and not indigenous (*p_cob1*)
- English speaking countries (*p_cob2*)
- Non-English speaking countries (*p_cob3*)
- Indigenous (*p_cob4*)

Poor English speaking abilities (*p_poeng*) records whether the individual has poor English speaking abilities.

Remoteness records whether the individual lives in:

- A major city (*p_urdg1*)
- An inner region (*p_urdg2*)
- Outer and remote areas or migratory in nature (*p_urdg3*)

Marital status in 2001 records whether in 2001 the individual was:

- Married (*p_mar1*)
- De facto (*p_mar2*)
- Separated (*p_mar3*)
- Divorced (*p_mar4*)
- Widowed (*p_mar5*)
- Single and never been married (*p_mar6*)

Parental Status

Number of dependents in 2001 (*p_noch*) records the number of dependent children the individual had in 2001.

Physical Health

Severity of health conditions in 2001 records whether the individual had:

- No health conditions (*p_ddeg1*)
- A mild condition (*p_ddeg2*)
- A moderate condition (*p_ddeg3*)
- A severe condition (*p_ddeg4*)

Labour Force Variables

Labour market status in 2001 records whether the individual was:

- Employed (*p-lfs1*)
- Unemployed (*p-lfs2*)
- Not in the labour market (*p-lfs3*)

Extent of working hour match with preferences in 2001 records whether the match between the individual's total weekly working hours across all jobs and their preferred number of working hours made them:

- Not working (*p-whp1*)
- Underemployed by at least 4 hours a week (*p-whp2*)
- Roughly Matched: Preferred and Actual Hours Worked differ by less than 4 hours a week (*p-whp3*)
- Overemployed by at least 4 hours a week (*p-whp4*)

Employee type in 2001 records whether the individual was:

- Not working (*p-emp1*)
- An employee (*p-emp2*)
- An employee of own business (*p-emp3*)
- Self Employed (*p-emp4*)
- Unpaid family worker (*p-emp5*)

Contract type in 2001 records whether the individual was:

- Not working (*p-con1*)
- On a fixed term contract (*p-con2*)
- On a casual contract (*p-con3*)
- On a permanent contract (*p-con4*)
- On other types of contracts (*p-con5*)

Occupation in 2001 records whether the individual was working as:

- Not working (*p-occ1*)
- Armed forces (*p-occ2*)
- Legislators, Senior Officials and Managers (*p-occ3*)
- Professionals (*p-occ4*)
- Technicians and Associate Professionals (*p-occ5*)
- Clerks (*p-occ6*)
- Service Workers and Shop and Market Sales Workers (*p-occ7*)

- Skilled Agriculture and Fishery Workers (*p_occ8*)
- Craft and Related Trades Workers (*p_occ9*)
- Plant and Machine Operators and Assemblers (*p_occ10*)
- Elementary Occupations (*p_occ11*)

Household income in 2001 (*p_rehdi*) records the real value of the individual's total household income indexed at 2012 price levels and adjusted for household size.

Partner labour force status in 2001 records whether the individual had:

- No partner or no resident partner (*p_plfs1*)
- A partner who was employed (*p_plfs2*)
- A partner who was unemployed (*p_plfs3*)
- A partner who was not in the labour force (*p_plfs4*)

Parental information

Father's country of birth records whether or not the individual's father was born in:

- Australia (*p_fcob1*)
- English speaking countries (*p_fcob2*)
- Non-English speaking countries or indigenous (*p_fcob3*)

Mother's country of birth records whether or not the individual's mother was born in:

- Australia (*p_mcob1*)
- English speaking countries (*p_mcob2*)
- Non-English speaking countries or indigenous (*p_mcob3*)

Father's education records whether the individual's father's highest education, as reported in 2005, was:

- None (*p_fedu1*)
- Primary (*p_fedu2*)
- Below secondary (*p_fedu3*)
- Secondary (*p_fedu4*)
- Post-secondary, non-university (*p_fedu5*)
- Post-secondary, university (*p_fedu6*)

Mother's education records whether the individual's mother's highest education, as reported in 2005, was:

- None (*p_medu1*)

- Primary (*p_medu2*)
- Below secondary (*p_medu3*)
- Secondary (*p_medu4*)
- Post-secondary, non-university (*p_medu5*)
- Post-secondary, university (*p_medu6*)

Father undertaken post-school qualification through employer or non-tertiary means (*p_fpsm*) records whether the individual's father had undertaken his highest qualification through employers or other channels other than tertiary education, as reported in 2005.

Mother undertaken post-school qualification through employer or non-tertiary means (*p_mpsm*) records whether the individual's mother had undertaken his highest qualification through employers or other channels other than tertiary education, as reported in 2005.

Father's Employment at age 14 records whether the individual's father was working when they were aged 14, in the following categories:

- Father deceased or not living with respondent (*p_femp1*)
- Father not employed (*p_femp2*)
- Father employed (*p_femp3*)

Mother's Employment at age 14 (*p_memp*) records whether the individual's mother was working when they were aged 14, in the following categories:

- Mother deceased or not living with respondent (*p_memp1*)
- Mother not employed (*p_memp2*)
- Mother employed (*p_memp3*)

Father substantially unemployed growing up records whether the individual's father had been unemployed for 6 months or more when they were growing up, in the following categories:

- Father not living with respondent (*p_fsue1*)
- Father not substantially unemployed (*p_fsue2*)
- Father substantially unemployed (*p_fsue3*)

Father's Occupation records whether at age 14 the individual's father was last known working as:

- Father not in household (*p_focc1*)
- Armed forces (*p_focc2*)
- Legislators, Senior Officials and Managers (*p_focc3*)
- Professionals (*p_focc4*)

- Technicians and Associate Professionals (*p_focc5*)
- Clerks (*p_focc6*)
- Service Workers and Shop and Market Sales Workers (*p_focc7*)
- Skilled Agriculture and Fishery Workers (*p_focc8*)
- Craft and Related Trades Workers (*p_focc9*)
- Plant and Machine Operators and Assemblers (*p_focc10*)
- Elementary Occupations (*p_focc11*)

Mother's Occupation records whether at age 14 the individual's mother last known working as:

- Mother not in household (*p_mocc1*)
- Armed forces (*p_mocc2*)
- Legislators, Senior Officials and Managers (*p_mocc3*)
- Professionals (*p_mocc4*)
- Technicians and Associate Professionals (*p_mocc5*)
- Clerks (*p_mocc6*)
- Service Workers and Shop and Market Sales Workers (*p_mocc7*)
- Skilled Agriculture and Fishery Workers (*p_mocc8*)
- Craft and Related Trades Workers (*p_mocc9*)
- Plant and Machine Operators and Assemblers (*p_mocc10*)
- Elementary Occupations (*p_mocc11*)

Non-cognitive variables

Well-being in 2001 (*p_losat*) records the life satisfaction score, which ranges from 0 to 10, of the individual reported in 2001. A higher score means the individual is more satisfied with his/her life.

Attitude towards having job in 2001 (*p_jbwk*) records the average score of attitude towards having a job reported by the individual in 2001 across two items (*p_jadnm* and *p_jahpj*), in a scale ranging from 1 to 7, with a higher score indicating a more favourable attitude towards having a job.

Enjoy job without needing money in 2001 (*p_jadnm*) records the extent the individual agreed with the statement that the person would enjoy having a job even if they did not need the money in 2001, in a scale ranging from 1 to 7, with a higher score indicating more agreement.

Important to have paying job in 2001 (*p_jahpj*) records the extent the individual agreed with the statement that in order to be happy in life it is important to have a paying job in 2001, in a scale ranging from 1 to 7, with a higher score indicating more agreement.

Prior Year Outcome variables

Mental health in 2001 (*p_mh01*). This is the transformed mental health scores from the aggregation of mental health items of the SF-36 Health Survey, as reported by the individual in 2001. It ranges from 0 to 100, with higher scores indicating better mental health.

Mental health in 2001 below norm (*p_mb01*) records whether the individual's mental health scores for 2001 was below the average of mental health scores across our analytical sample for that year.

Working hours in 2001 (*p_wh01*) records the number of hours the individual works across all jobs in a week on average. Working hours are set to 0 for those not working.

Hourly Wages in 2001 (*p_hrw01*) records the average hourly wage of the individual's main job in 2001. Hourly wages are set to 0 for those not working and set to missing for those reporting working more than 100 hours a week.

F.2.4 Variables that are not included in the model

The unique person identifier (*xwaveid*).

Completed retraining after 2017 based on highest education (*rehllt*) records whether the individual had only completed their retraining after 2017, comparing their education level in 2017 and 2019.

Completed retraining after 2017 based on detailed qualifications (*redllt*) records whether the individual has completed any one of the following qualifications since last interviewed between 2018 and 2019:

- Trade certificate or apprenticeship
- Technicians cert/Advanced certificate
- Teaching qualification
- Nursing qualification
- Associate Degree
- Advance Diploma (3 years full time or equivalent)
- Bachelor degree but not honours
- Certificate I
- Certificate II
- Certificate III
- Certificate IV
- Certificate of unknown level
- Doctorate
- Diploma NFI
- Diploma (2 years full time or equivalent)
- Graduate Certificate
- Graduate Diploma

- Honours
- Masters
- Other

Completed retraining after 2017 based on both highest attainment and detailed qualifications (*refllt*) records whether the individual has completed retraining after 2017 based on both the variables *rehllt* and *redllt*. When either of these variables has a value of 1, this variable will take on the value of 1.

Timing of Education Completion

Year of first retraining completion records the year of the first reported instance of retraining completion as provided by the detailed qualification variables and include the following categories:

- 2002 (*p_rcom1*)
- 2003 (*p_rcom2*)
- 2004 (*p_rcom3*)
- 2005 (*p_rcom4*)
- 2006 (*p_rcom5*)
- 2007 (*p_rcom6*)
- 2008 (*p_rcom7*)
- 2009 (*p_rcom8*)
- 2010 (*p_rcom9*)
- 2011 (*p_rcom10*)
- 2012 (*p_rcom11*)
- 2013 (*p_rcom12*)
- 2014 (*p_rcom13*)
- 2015 (*p_rcom14*)
- 2016 (*p_rcom15*)
- 2017 (*p_rcom16*)
- 2018 (*p_rcom17*)
- 2019 (*p_rcom18*)

Locus of control in 2003 (*p_cotrl*) records the transformed composite score¹⁶ for locus of control items reported by the individual in 2003, the first year in HILDA for which this information becomes available. The transformation results in a variable that is ranged between 7 and 49. Locus of control measures the degree to which individuals attribute outcomes to internal versus external factors or the extent their welfare are in their own control compared to external circumstances. A higher score indicates having a more external locus of control, which is considered as a favourable personality trait.

¹⁶See Buddlemeyer and Powdthavee (2015) for details of the transformation.

Frequency of reading books in 2012 (*p_rdf*) records the frequency the individual reads books in 2012, the first year in HILDA for which this information becomes available. This is a proxy for love of learning¹⁷. This is a categorical variable encompassing the following frequencies:

- Every day or most days (*p_rdf1*)
- Several times a week (*p_rdf2*)
- About once a week (*p_rdf3*)
- 2 or 3 times a month (*p_rdf4*)
- About once a month (*p_rdf5*)
- Less than once a month (*p_rdf6*)
- Never (*p_rdf7*)

¹⁷HILDA contains a question on reading newspapers and magazines but we feel that reflects a care for or understanding of current issues more than a love of learning.

Online Appendix G: Panel Sample: Sensitivity Analysis

G.1 Sample Selection

Treated sample: For any person in HILDA who ever reported *starting* a degree (determined by taking a person who switches from reporting “not currently studying” in one wave to “currently studying” in the next wave) and/or *completing* a degree, we select their first study event as a treatment observation if it satisfies three other conditions.

They are: (1) at least 21 years old in the starting year of study¹⁸, (2) they were present in the two years before the start of study (in order to have information on their feature values), (3) there were not currently studying in any of the two years before the starting year of further study (to avoid reverse-causation issues), (4) they completed their further degree and (5) they were present in the survey and had a non-missing outcome 4 years after the start of study.

If a study event does not satisfy these conditions, we look to the next study event that satisfies these conditions or (if unavailable) delete the person from our sample completely. Conditions (3) and (5) together mean that we analyse a sample of individuals who started their degrees anytime between 2003 and 2015.

In our treated group, 1,814 individuals **started and completed a further educational degree**.

Control sample: These are those who had never started retraining throughout HILDA. From these control observations, we assign a time stamp to them for the year the control person theoretically started to study. We do this for every year from 2003 to 2019. This implies that never re-educated individuals can be duplicated and used multiple times. For example, if a control individual is observed throughout the years 2001 to 2016, then they will be a control for the separate treated individuals that started retraining in 2003, in 2004, 2005 and up to 2017 i.e. the control individual will be duplicated 15 times.

There are 60,945 control observations i.e. individuals who never completed a further degree. However, as described above, these are non-unique observations in the sense that a control individual can be duplicated up to 15 times.

G.2 Variable description

G.2.1 Outcome Variables

Weekly earnings from main job in fourth year after the individual started their retraining (*f4_wscmei*) records the weekly earnings from the main job for the individual in the fourth year after the individual started their retraining.

¹⁸Note that we expanded the age range in this sensitivity analysis to ensure sufficient treatment observations for the estimation of the treatment outcome surfaces.

G.2.2 Treatment Variables: Retraining

Retraining completion based on both highest attainment and detailed qualifications (*reduft*) records whether the individual has completed retraining based on a comparison of the highest education attainment and the number of qualifications gained across waves 1 and 17. If either of these have gone up, *reduft* takes a value of 1 and 0 otherwise.

G.2.3 Input Variables

Characteristics in the Year Prior to Retraining Start

Demographics

Gender (*p1_hgsex*) records the gender of the individual. The value of 1 denotes males whereas the value 2 denotes females.

Age (*p1_hgage*) records the age of the individual in the year prior to retraining start.

Country of birth (*p1_anbcob*) records whether or not an individual was born in:

- Australia (value=1)
- English speaking countries (value=2)
- Non-English speaking countries (value=3)

Indigenous Status (*p1_anatsi*) records whether or not an individual is:

- Not indigenous (value=1)
- Aboriginal (value=2)
- Torres Islander (value=3)
- Both Aboriginal and Torres Islander (value=4)

Poor English speaking abilities (*p1_poeng*) records whether the individual has poor English speaking abilities in the year prior to retraining start.

State of residence (*p1_hhstate*) records the state of residence of the individual in the year prior to retraining start:

- NSW (value=1)
- VIC (value=2)
- QLD (value=3)
- SA (value=4)
- WA (value=5)
- TAS (value=6)
- NT (value=7)

- ACT (value=8)

Remoteness (*p1_hhsos*) records whether, in the year prior to retraining start, the individual lives in:

- A major city (value=0)
- An inner region (value=1)
- Outer and remote areas (value=2)
- migratory in nature (value=3)

Marital status (*p1_mrcurr*) records whether, in the year prior to retraining start, the individual was:

- Married (value=1)
- De facto (value=2)
- Separated (value=3)
- Divorced (value=4)
- Widowed (value=5)
- Single and never been married (value=6)

Household size (*p1_hhsize*) records the total number of individuals living in the same household as the individual (including the individual) in the year prior to retraining start.

Sexual orientation (*p1_lgtb*) records that the individual's sexual orientation is not heterosexual. The variable is constructed from the Sexual Identity question that is only asked in waves 12 and 16. We combine answers from both waves to create a binary indicator for the individual ever reporting a sexual identity that is not heterosexual, treating sexual orientation as a fixed trait for a given individual.

Parental Status

Number of dependents (*p1_totalkids*) records the number of children under 15 the individual had in the household in the year prior to retraining start.

Having children (*p1_anykid*) records the individual had any dependents in the household in the year prior to retraining start.

Children under 5 (*p1_kidu5*) records the individual had children under 5 in the household in the year prior to retraining start.

Age of youngest (*p1_rcyng*) records the age of the youngest children living with the respondent in the year prior to retraining start (including adult children).

Physical Health

Severity of health conditions (*p1_disdeg*) records whether, in the year prior to retraining start, the individual had:

- No health conditions (value=0)

- A mild condition (value=1)
- A moderate condition (value=2)
- A severe condition (value=3)

Labour Force Variables

Labour market status (*p1_lfs*) records whether the individual was:

- Employed (value=1)
- Unemployed (value=2)
- Not in the labour market (value=3)

Extent of working hour match with preferences (*p1_whpref*) records whether, in the year prior to retraining start, the match between the individual's total weekly working hours across all jobs and their preferred number of working hours made them:

- Underemployed by at least 4 hours a week (value=1)
- Roughly Matched: Preferred and Actual Hours Worked differ by less than 4 hours a week (value=2)
- Overemployed by at least 4 hours a week (value=3)

Employee type (*p1_emptytype*) records whether, in the year prior to retraining start, the individual was:

- An employee (value=1)
- An employee of own business (value=2)
- Self Employed (value=3)
- Unpaid family worker (value=4)

Contract type (*p1_conttype*) records whether, in the year prior to retraining start, the individual was:

- On a fixed term contract (value=1)
- On a casual contract (value=2)
- On a permanent contract (value=3)
- On other types of contracts (value=4)

Occupation (*p1_occ*) records whether, in the year prior to retraining start, the individual was working as:

- Armed forces (value=0)
- Legislators, Senior Officials and Managers (value=1)
- Professionals (value=2)
- Technicians and Associate Professionals (value=3)

- Clerks (value=4)
- Service Workers and Shop and Market Sales Workers (value=5)
- Skilled Agriculture and Fishery Workers (value=6)
- Craft and Related Trades Workers (value=7)
- Plant and Machine Operators and Assemblers (value=8)
- Elementary Occupations (value=9)

Union membership (*p1_union*) records whether the individual was a union member in the year prior to retraining start.

Real household income (*p1_rhdi*) records the real value of the individual's total household income indexed at 2012 price levels and adjusted for household size in the year prior to retraining start.

Partner labour force status (*p1_plfs*) records whether, in the year prior to retraining start, the individual:

- Had no partner or no resident partner (value=0)
- Had a partner who was employed (value=1)
- Had a partner who was unemployed (value=2)
- Had a partner who was not in the labour force (value=3)

Years in paid work (*p1_ehtjb*) records the total number of years in paid work the individual has spent in the year prior to retraining start.

Percent finding as least as good a job (*p1_jbmggj*) records, for employees, the percentage that they will find as least as good a job as they currently have in their own estimation in the year prior to retraining start

Occupational scale (*p1_jbmo6s*) records the Australian Socioeconomic Index 2006 ranking of the individual's occupation in the year prior to retraining start. It ranges from 0 to 100, with higher scores indicating higher occupational status.

Tenure with employer (*p1_jbempt*) records the total years spent with the current employer for the individual in the year prior to starting retraining.

Parental information

Father's country of birth (*p1_fcob*) records whether or not the individual's father was born in:

- Australia (value=1)
- English speaking countries (value=2)
- Non-English speaking countries or indigenous (value=3)

Mother's country of birth (*p1_mcob*) records whether or not the individual's mother was born in:

- Australia (value=1)

- English speaking countries (value=2)
- Non-English speaking countries or indigenous (value=3)

Father's education records whether the individual's father's highest education, as reported in 2005, was:

- None (value=1)
- Primary (value=2)
- Below secondary (value=3)
- Secondary (value=4)
- Post-secondary, non-university (value=5)
- Post-secondary, university (value=6)

Mother's education records whether the individual's mother's highest education, as reported in 2005, was:

- None (value=1)
- Primary (value=2)
- Below secondary (value=3)
- Secondary (value=4)
- Post-secondary, non-university (value=5)
- Post-secondary, university (value=6)

Father undertaken post-school qualification through employer or non-tertiary means (*p_fpsm*) records whether the individual's father had undertaken his highest qualification through employers or other channels other than tertiary education, as reported in 2005.

Mother undertaken post-school qualification through employer or non-tertiary means (*p_mpsm*) records whether the individual's mother had undertaken his highest qualification through employers or other channels other than tertiary education, as reported in 2005.

Father's Employment at age 14 (*p1_femp*) records whether the individual's father was working or not when they were aged 14.

Mother's Employment at age 14 (*p1_memp*) records whether the individual's mother was working or not when they were aged 14.

Father substantially unemployed growing up (*p1_fsue*) records whether the individual's father had been unemployed or 6 months or more when they were aged 14.

Father's Occupation (*p1_focc*) records whether at age 14 the individual's father was last known working as:

- Armed forces (value=0)
- Legislators, Senior Officials and Managers (value=1)

- Professionals (value=2)
- Technicians and Associate Professionals (value=3)
- Clerks (value=4)
- Service Workers and Shop and Market Sales Workers (value=5)
- Skilled Agriculture and Fishery Workers (value=6)
- Craft and Related Trades Workers (value=7)
- Plant and Machine Operators and Assemblers (value=8)
- Elementary Occupations (value=9)

Mother's Occupation (*p1_mocc*) records whether at age 14 the individual's mother last known working as:

- Armed forces (value=0)
- Legislators, Senior Officials and Managers (value=1)
- Professionals (value=2)
- Technicians and Associate Professionals (value=3)
- Clerks (value=4)
- Service Workers and Shop and Market Sales Workers (value=5)
- Skilled Agriculture and Fishery Workers (value=6)
- Craft and Related Trades Workers (value=7)
- Plant and Machine Operators and Assemblers (value=8)
- Elementary Occupations (value=9)

Income Support

On income support (*p1_onis*) records the individual was on income support in the year prior to starting retraining

On Newstart (*p1_onnsa*) records the individual was on Newstart Allowance in the year prior to starting retraining

On Age Pension (*p1_onap*) records the individual was on Age Pension in the year prior to starting retraining

On DSP (*p1_ondsp*) records the individual was on Disability Support Pension in the year prior to starting retraining

On Carer Payment (*p1_oncp*) records the individual was on Carer Payment in the year prior to starting retraining

On Widow Allowance/Wife Pension (*p1_onww*) records the individual was on Widow Allowance/Wife Pension in the year prior to starting retraining

On Youth Allowance (*p1_ony*) records the individual was on Youth Allowance in the year prior to starting retraining

On Mature Age Allowance (*p1_onma*) records the individual was on Mature Age Allowance in the year prior to starting retraining

On Mature Age Partner Allowance (*p1_onmap*) records the individual was on Mature Age Partner Allowance in the year prior to starting retraining

On Ab/Austudy (*p1_onsdy*) records the individual was on Ab/Austudy in the year prior to starting retraining

On Bereavement Allowance (*p1_onba*) records the individual was on Bereavement Allowance in the year prior to starting retraining

On Sickness Allowance/Special Benefits (*p1_onsab*) records the individual was on Sickness Allowance/Special Benefits in the year prior to starting retraining

On Partner Allowance (*p1_onpa*) records the individual was on Partner Allowance in the year prior to starting retraining

On Parenting Payments (*p1_onpp*) records the individual was on Parenting Payments in the year prior to starting retraining

Housing situation

Mortgage balance (*p1_hsmgowe*) records the amount still owing on the mortgage that the individual had in the year prior to retraining start. For those without a mortgage or not home owner, the mortgage balance is set to 0.

Non home owners (*p1_renter*) records whether the individual was renting or not living in their own homes in the year prior to retraining start.

Prior Year Outcomes

Weekly income from all jobs (*p1_earning*) records the weekly earnings from all jobs for the individual in the year prior to the individual starting their retraining.

Weekly income from main job (*p1_wscmei*) records the weekly earnings from the main job for the individual in the year prior to the individual starting their retraining.

Weekly working hours (*p1_wkhr*) records the total number of hours the individual works in all jobs in a week on average in the year prior to the individual started their retraining. Working hours are set to 0 for those not working.

Real hourly wage (*p1_rlwage*) records the real hourly wage of the individual in the year prior to the individual starting their retraining, indexed at 2012 price levels. Hourly wages are set to 0 for those not working and set to missing for those reporting working more than 100 hours a week. All wages have then been adjusted up by \$1 to preserve sample size for the logarithm transformation.

Log hourly wage (*p1_lnwage*) records the log of *p1_rlwage*.

Mental health (*p1_ghmh*) records the transformed mental health scores from the aggregation of mental health items of the SF-36 Health Survey, as reported by the individual in the year prior to the individual started their retraining. It ranges from 0 to 100, with higher scores indicating better mental health.

Life satisfaction (*p1_losat*) records the life satisfaction score reported by the individual in the year prior to the individual started their retraining. It ranges from 0 to 10, with higher scores indicating higher life satisfaction.

Delta variables

For all the variables described in the preceding section titled Characteristics in the Year Prior to Retraining Start, we create a further set of change or delta variables. Specifically, each delta variable is the subtracting of the value of a given characteristic in the two years prior to starting retraining from the value of this characteristic in the year prior to retraining start.

All delta variables are denoted by the *d_* prefix.

Education-related variables

Level of retraining completed: Bachelor and above (*bachab*) records whether the individual had completed retraining at bachelor and above. The variable is set to 0 for the control group and missing for those who had completed certificates.

Level of retraining completed: Below Bachelor (*bbach*) records whether the individual had completed retraining that is below bachelor level. The variable is set to 0 for the control group and missing for those who had completed a bachelor or higher qualification.

Main field of study: technical degree (*techdeg*) records whether the individual's main field of study was a technical degree. The variable is set to 0 for the control group and missing for those whose main field of study was a qualitative degree. Technical degrees include:

- Natural and physical sciences
- Information technology
- Engineering and related technologies
- Architecture and building
- Agriculture, environment and related studies
- Medicine
- Nursing
- Other health-related (e.g. Pharmacy, Dental studies, Rehabilitation therapies, Optical science, Veterinary studies)
- Management and commerce (e.g. Accounting, Business, Sales and marketing, Banking and finance, Office studies)
- Law

Main field of study: qualitative degree (*qualdeg*) records whether the individual's main field of study was a qualitative degree. The variable is set to 0 for the control group and missing for those whose main field of study was a technical degree. Qualitative degrees include:

- Education
- Society and culture (e.g. Economics, Political science, Social work, History, Psychology, Languages, Religion, Sport)
- Creative arts
- Food, hospitality and personal services
- Other

Study duration (*fsddur*) records the total number of waves an individual had spent studying from the start of their first study event counted in our sample.

Starting Study intensity (*csftsd*) records whether the individual was studying full time or not when they started their retraining.

Finishing Study intensity (*fsftsd*) records whether the individual was studying full time or not when they completed their retraining.

Other variables

Number of waves in HILDA (*numwave*) records the number of waves in which the respondent has submitted a valid response for the HILDA survey.

G.2.4 Variables that are not included in the model

The unique person identifier (*xwaveid*)

Wave started retraining (*icswave*)

Wave completed retraining (*ifswave*)

Control group indicator (*control*)

Started but did not complete retraining between 2003-2017 (*ncomp*)

Starting year of retraining imputed (*impute*) is a binary indicator for individuals for which we observe their retraining completion but they never reported ever starting retraining and so we had to impute a starting wave for these individuals.

Started re-education in wave 2018/19 (*lateststart*) is an indicator for those individuals who had started their retraining in 2018 or 2019.

Online Appendix H: Nested CV Holdout Sample

Table 5: Nested CV Holdout Sample: Level Earnings

Model	Outcome surface	Negative MSE	NMSE Std	R-squared	R-squared Std	ATE	ATE Std
GBR	Treated	-886515	452077	0.22	0.06	68.2	28.4
	Control	-659056	107251	0.36	0.07		
LASSO	Treated	-955958	361911	0.15	0.09	94.1	14.5
	Control	-710521	178030	0.32	0.05		
Ridge	Treated	-966849	434518	0.16	0.08	97.8	14.5
	Control	-712374	174033	0.32	0.04		

Notes: 5 fold CV performed on 80% train sample. All statistics presented in this table are based on the 20% holdout sample. Ten outer folds are used. See Figure 1 for more details.

Table 6: Nested CV Holdout Sample: Entry into Entrepreneurship

Model	Outcome surface	Negative MSE	NMSE Std	R-squared	R-squared Std	ATE	ATE Std
GBR	Treated	-0.071	0.021	-0.017	0.016	0.026	0.004
	Control	-0.044	0.008	0.022	0.012		
LASSO	Treated	-0.070	0.015	-0.002	0.017	0.025	0.003
	Control	-0.045	0.009	0.009	0.007		
Ridge	Treated	-0.070	0.021	-0.015	0.022	0.022	0.002
	Control	-0.044	0.006	0.024	0.018		

Notes: 5 fold CV performed on 80% train sample. All statistics presented in this table are based on the 20% holdout sample. Ten outer folds are used. See Figure 1 for more details.