

Deleting Unreported Innovation*

Ping-Sheng Koh: ESSEC Business School; kohp@essec.edu

David M. Reeb: National University of Singapore, Senior Fellow: *ABFER*; dmreeb@nus.edu.sg

Elvira Sojli: University of New South Wales; e.sojli@unsw.edu.au

Wing Wah Tham: University of New South Wales; w.tham@unsw.edu.au

Wendun Wang: Erasmus School of Economics, Erasmus University; wang@ese.eur.nl

October 26, 2020

Abstract

Innovation variables exhibit high rates of unobservability, often leading empirical studies to exclude firms that fail to report innovation. We assess the reliability of six methods for dealing with unobserved innovation using several different counterfactuals for firms without R&D or patents. These tests reveal that deleting firms without reported innovation or imputing them as zero innovators and including a dummy variable leads to biased parameter estimates for reported innovation and other *explanatory variables*. Deleting or ignoring firms without patents is especially problematic, leading to false-positive results in empirical tests. Our analysis suggests using both multiple imputation and instrumental variable estimates.

Keywords: Bias, Listwise Deletion, Innovation, Measuring Innovation, Multiple Imputation, Non-Patenting firms, Unreported R&D, Patents.

* An early version of this manuscript circulated under the title “Missing Innovation Around the World.” We are grateful to seminar participants at City University of Hong Kong, Hong Kong Polytechnic University, Louisiana State University, Maastricht University, the National University of Singapore, Rotterdam School of Management, Temple University, University of Queensland, University of South Carolina, University of Toronto, Virginia Tech, Wilfrid Laurier University, and the University of Technology Sydney. We appreciate discussions and advice from Renee Adams, Sumit Agarwal, Benito Arrunda, Richard Boylan, Lora Dimitrova, Bronwyn Hall, Gilles Hilary, Yael Hochberg, Qin Li, Gustavo Manso, Jiaming Mao, Naci Mocan, Randall Morck, Ivan Png, Wenlan Qian, Amit Seru, Vijay Singal, Xuan Tian and Rosemarie Ziedonis.

1. Introduction

Investors and academics exhibit substantial interest in understanding how corporate innovation influences firm growth and performance (Hochberg et al., 2018). Empirical studies typically use patents or R&D expenditures to measure firm innovation, often focusing on innovation as a variable of interest (e.g., Croce et al., 2018) or as a control variable (e.g., Huang, 2018). A well-known problem in this cross-disciplinary body of research is that most firms do not report their R&D spending nor obtain patents (Anton and Yao, 2004; Koh, Reeb, and Zhao, 2018). For instance, over 50% of US firms do not report their R&D spending and among firms with positive R&D for over a decade, roughly 60% of them do not obtain patents. Potential explanations for not reporting R&D or obtaining patents include negligible innovation inputs, unsuccessful innovation projects, or attempts to keep the innovation information secret (Png, 2017). Recent studies develop additional measures of corporate innovation, ranging from textual analysis in financial analysts reports (Bellstam et al., 2020) to firm disclosures of new products (Mukerjee et al., 2017). These new measures provide thoughtful approaches to capture different aspects of innovation, but typically suffer from the same concerns of truncation bias noted by Lerner and Seru (2017) regarding patents.

Recognizing that most firms fail to report R&D expenditures or seek patents, empirical researchers use a variety of methods to handle unreported innovation. The two most common approaches to dealing with this issue are excluding firms without R&D or patents (e.g., Hombert and Matray, 2018) or classifying these firms as zero innovators and including a dummy variable as suggested by Koh and Reeb (2015) (e.g., Masulis and Zhang, 2019). Are these the right approaches to handling unreported innovation?

Our analysis focuses on how these common methods for dealing with missing innovation data influence the economic conclusions in empirical finance research. To address this issue, we

investigate the reliability of different methods for handling unreported innovation in studies that focus on innovation as either an explanatory variable of interest or include it as a control variable. We compare six approaches to handling unreported innovation: Listwise deleting (discarding) firms without R&D or patents, deterministic imputation with either zero or industry average, inverse probability weighting, Heckman selection, and multiple imputation. Multiple imputation is arguably the least common method in corporate finance and relies on estimating the missing variable of interest using other observable covariates and explicitly adjusting for imputation uncertainty (see Internet Appendix I).

Our preliminary analysis reveals that firms with unreported innovation (R&D or patents) are predictable with known determinants of innovation and other corporate outcomes of interest, rejecting the hypothesis that unobservable innovation is *missing completely at random*.¹ These results raise the concern that the commonly used methods to handle unreported innovation could lead to biased parameter estimates due to the non-representativeness of the population under study (deletion) or distortions of the variance-covariance matrix (deterministic imputation).

Our empirical tests focus on two common measures of corporate innovation, namely R&D spending and patents (whether based on counts, citations, or market reactions). Firms without reported R&D could arise from a lack of R&D spending or a disclosure choice of the firm. Similarly, zero patents firms could stem from failed innovation projects, trade-secret based reporting choices, or from firms seeking regulatory protection in different jurisdictions. Each of these different types of missingness allows for differing tests of the different approaches to handling unreported innovation. Our empirical analysis focuses on both R&D and patents.

¹ Terminology in statistics differentiates between three types of missing data. *Missing Completely at Random* (MCAR) occurs when neither observables nor unobservables predict missing observations. *Missing at Random* (MAR) occurs when observables can predict missing observations and *Missing Not at Random* (MNAR) occurs when missing observations are related to observable and unobservable data (see Section 2 and Internet Appendix II for details).

We test the methods to handling firms without reported R&D in two ways. First, we use data on US firms that did not report R&D spending in a particular period but reported the amount in subsequent financial statements. As firms that initiate the reporting of R&D expenditures are required to report their R&D expenditures for prior years, we can directly compare several common treatments for missing R&D in the prior years. To the empiricist these firms do not appear to engage in R&D activity in the years R&D spending was not reported, creating a natural laboratory to evaluate different methods of handling unreported innovation. We denote this newly reported R&D spending in future financial statements as “*Recovered R&D*.” Using recovered R&D as a baseline, we compare the observed/counterfactual R&D with replacing these firms’ missing R&D with zero, the industry average, and multiple imputation. Further tests show that the R&D in firms with unreported R&D significantly differ from zero R&D firms, the average industry R&D, and positive R&D firms in aggregate. Notably, we find that on average multiple imputation gives estimates of R&D for the missing R&D data that is qualitatively similar (not statistically different) to their actual R&D reported in subsequent financial statements.

One potential issue with using recovered R&D to compare treatments for unreported observations is that these firms may differ from other firms that do not report innovation in the Compustat universe. To address this concern, we use an alternative counterfactual group based on the textual analysis measure of innovation. Specifically, we use the text-based ranking of corporate innovation in the S&P 500 firms from Bellstam et al. (2020). These analysts’ discussions about corporate innovation likely only arise in firms that engaged in R&D. Again, we find that multiple imputation is the best solution to handle missing R&D in this alternative counter-factual group.

Figure 1 shows the dramatically different innovation rankings of S&P 500 firms when using reported R&D and the textual analysis of Bellstam et al. (2020). The standard practice of classifying firms that fail to report R&D (or patents) as zero innovators leads to a large clumping of firms at

the bottom of the distribution. In contrast, the textual analysis approach, while only viable for a subset of firms, provides a more complete distribution of corporate innovation. Similarly, the multiple imputation approach to handling unreported R&D also mitigates the clumping problem inherent with deleting these firms or classifying them as zero innovators. Still, our results so far are based on recovered R&D and firms with analysts' coverage, which may differ from other firms with unreported innovation.

To further mitigate external validity concerns from the two different counter-factuals of R&D, we undertake two simulation studies. In the first simulation analysis, we use the empirical distribution of US Compustat data to evaluate the impact of differing levels of missingness of an innovation variable. This simulation approach approximates the analyses typically found in empirical studies of corporate innovation using panel data. The second simulation analysis uses clearly specified data generating processes, allowing us to gauge the impact of unreported innovation in a controlled, cross-sectional setting. This approach mitigates concerns about the comparability of the data on US firms in our first simulation with the data used in other studies. In both simulation exercises, we evaluate the six different approaches noted previously to handle unreported innovation. Our simulation analyses rely on two evaluation criteria: The bias (expressed as a proportion of the benchmark coefficient) and the root mean squared error (RSME) of the regression coefficient estimates. We evaluate the coefficient estimates on both the innovation variable (e.g., R&D or patents) and other control variables.

In both simulations, we find that deleting or excluding firms with unreported R&D leads to biased coefficient estimates for both R&D and any control variables correlated with R&D. Rather than providing a conservative approach, the deletion of firms without reported R&D is one of the worst methods for handling unreported innovation. For instance, if R&D is *missing at random*, then the average bias from excluding firms without reported R&D is almost ten times greater than

found using multiple imputation. Moreover, our analyses show that commonly used deterministic imputation models (e.g., replace missing R&D with zero or industry average and include a dummy variable) fare poorly in comparison to multiple imputation. In addition, we find that the RMSEs of the common methods for handling missing innovation are very large in comparison to multiple imputation in both simulation exercises (i.e. 70% larger). We also find that the bias and RMSE from deleting firms without observable innovation dramatically increase at higher rates of missingness. This simulation result suggests the high rate of unreported innovation in patents, relative to R&D spending, makes deletion especially challenging in studies focusing on patent-based metrics.

To illustrate the economic magnitude of inference problems with common approaches to handling unreported innovation, we replicate an influential finance study that uses R&D spending. Fama and French (2002) test the empirical predictions of the pecking order and trade-off models of capital structure and classify firms without reported R&D expenditure as zero R&D firms. We find that the coefficient estimates and standard errors for R&D and capital structure are significantly different when using multiple imputation to account for unreported innovation relative to the results when classifying these firms as zero innovators. Strikingly, under multiple imputation both R&D and market-to-book have positive coefficients, providing evidence consistent with pecking order theory predictions rather than the conflicting results reported in Fama and French (2002). We cannot categorically state which results are correct, but we do note that different approaches to handling missing R&D give opposing results, suggesting this is an important consideration in research design. The approach with the least bias and RMSE gives consistent results across different versions of this test. Studies that use corporate innovation variables should explicitly evaluate the appropriateness of the how they handle unreported innovation.

So far, our empirical analysis primarily concentrates on unreported R&D spending. Yet, patents from the United States Patent Office (USPTO) provide another common approach to measuring corporate innovation. The vast majority of US firms do not seek USPTO patents, which stems from limited innovation success, trade-secret choices, and firms that only seek patent protection outside the US. To evaluate different methods of handling firms without patents, we use new product announcements as the ground truth. We assume that firms have successful innovation when they make new product announcements, especially for major new product announcements. Yet, firms without patents tend to have fewer new product announcements than USPTO patenting firms. Consequently, classifying these firms as zero innovators, imputing them with the industry average number of patents, or deleting firms without patents is problematic. The multiple imputation estimates appear to place these non-patenting firms into the appropriate innovation categories.

Our final test focuses on firms without USPTO patents. The vast majority of patent-based empirical studies in the US rely on USPTO patents to measure innovation success. Among US firms, 69% of positive R&D firms never file for patents using USPTO data, while only 43% never file patent applications using the 30 global patent offices. This 26% wedge in unobserved patents for studies using USPTO patents provides another opportunity to examine different methods of handling unobserved innovation. In this particular case the nature of the missingness likely differs from trade-secret based reasons, limiting external validity. Yet, a benefit of focusing on this 26% wedge is that we can directly evaluate different methods of handling unobserved innovation in a large number of empirical studies. Strikingly, we again find that multiple imputation provides much closer estimates for the patents unobserved by the empirical researcher relying on USPTO patents than in other commonly used replacement methods.

The nature of unreported innovation is unknown to the researcher. How we handle this missing data problem ultimately comes down to our assumptions about the mechanisms of missingness. Implicitly, researchers deciding how to handle missing innovation data are making assumptions about whether missingness can or cannot be predicted by observables. For instance, the IV solution to missing innovation data relies on the ability to find truly exogenous shocks to overcome the selection bias (see discussion in Jiang, 2017). MI assumes that missingness can be predicted with observables, which implicitly facilitates the estimation of the average treatment effect under MAR. Given that the assumptions underlying IV and MI are both likely to be violated to some degree, the choice between assuming MAR or MNAR depends on the bias of the IV and MI estimates. Collins, Schafer, and Kam (2001) demonstrates that in many realistic cases, an erroneous assumption about MAR often has limited impact on estimates and standard errors because covariates included in the imputation models are often correlated with unobservable determinants of missingness. Consequently, we recommend using both MI and IV estimates when confronted with missing innovation data. Rather than simply deleting the firms with missing data (Branstetter, et al., 2019) or designating them as zero innovators (e.g Koch, et al., 2020), our analysis provides a roadmap for future researchers to adopt when using any common innovation measure as a treatment variable of interest or as a conditioning variable.

This study provides several insights and contributions to the innovation literature. First, studies on innovation should consider using several different approaches to handling unreported innovation (R&D and patents). In this context, we recommend that researchers provide some basic statistics for the degree or magnitude of the missing innovation data in their sample and how it relates to their key variable(s) of interest. Instead of simply deleting firms without patents, reported R&D, financial analysts' coverage, or new product announcements, we should attempt to adjust for the non-randomness in missing innovation data. We advise against the common approach of

performing the main analysis by replacing missing R&D with zero (or industry average) and then repeating the tests after excluding these non-reporting firms as sensitivity analysis. Both deterministic imputation and listwise deletion of firms with unreported corporate innovation can provide biased coefficient estimates if the missingness is non-random, making it difficult to evaluate how well one biased approach can provide a robustness test for another biased approach.

Second, this study contributes to the burgeoning work on the econometric challenges faced by researchers in finance. Bertrand, Duflo and Mullainathan (2004), Petersen (2009), and Thompson (2011) discuss methods to appropriately compute standard errors in the presence of cross-sectional and time-series dependence across residuals. Koh and Reeb (2015) compare two deterministic imputation methods and find that including a dummy variable for missing R&D firms improves imputing with zero or the industry average in their regressions. They are silent on the relative biasness of these two methods, and they do not evaluate excluding firms without reported R&D, multiple imputation, inverse probability weighting, or Heckman models. Our analysis shows that deterministic imputation solutions can lead to biased estimates and standard errors for both unreported R&D and patents, as well as other control variables. Importantly, our paper questions the foundations for deleting firms with unreported innovation (widely adopted in economics and finance) and the impact of using deterministic imputation models that classify these firms as zero innovation firms and including a dummy variable (frequently used by accounting and finance scholars). In addition, we show that the use of patent-based metrics, such as patent counts, citations, or market reactions, as alternative innovation measures does not resolve the unreported innovation problem. Instead, the problem of unobservable innovation is arguably more pronounced in studies that use patents to measure innovation than in ones using R&D expenditures, because of the higher rate of missingness in patents.

Third, studies that use R&D or patents as control variables also suffer from this missing data bias. Best practices for dealing with missing R&D and patents depend on the source or type of missing innovation data. If country, industry, or firm characteristics predict unreported R&D or missing patents (see Lerner and Seru, 2017), then our analysis suggests that using multiple imputation provides the most reasonable solution. Surprisingly, and across a wide variety of specifications and approaches, we find that both Heckman and inverse probability weighting rarely provide the best approaches to handling unreported innovation in our samples. For alternative data sets, researchers could consider undertaking simulations similar to ours to evaluate the various methods of handling unreported innovation. We provide our code for researchers interested in performing simulations on unreported innovation using their own unique data.

2. Handling Missing Innovation Data

There are numerous possible reasons for why we observe missing innovation data. Of course, unreported innovation could arise because the firm does not engage in innovation and has nothing to report. Unfortunately, the missingness mechanism cannot be positively identified from examining the observable data. Hence, as empiricists, we make either implicit or explicit assumptions about the missingness mechanism for firms without patents or R&D spending to draw inferences, which are separate from the statistical methods we use for parameter estimation. In general, missing data causes two problems: Bias in the parameter estimates and loss in efficiency (Rubin, 1976). Bias stems from the non-representativeness of the population under study. Loss of efficiency arises because information loss is a direct consequence of missing data, i.e. smaller samples.

We focus on the assumptions underlying different practices for handling unreported innovation and consider the econometric implications of these common approaches when the

assumptions are violated. Research in statistics has long recognized and studied the broad class of missing data problems (e.g., Robins and Wang, 2000), while research in machine learning also focuses on training and testing data that suffers from missing observations (e.g., Grangier and Melvin, 2010). Yet, many of these techniques and methods are relatively unused in research on corporate innovation.

To provide a framework for investigating unreported innovation, we consider the case where only one explanatory variable contains missing observations. Let y_i be the dependent variable and z_i be the innovation variable with missingness. We have the linear relation:

$$y_i = \alpha + \theta z_i + \varepsilon_i, \quad i = 1, \dots, N. \quad (1)$$

Let s_i be a selection indicator where $s_i = 1$, when z_i is not missing and firm i is included in the regression. Otherwise, when $s_i = 0$ firm i is deleted from the data. The validity of solutions to this problem depends on the missingness mechanism, thus we first present the three missing mechanisms. Rubin (1976) and Little and Rubin (2002) classify missing data mechanisms into *Missing Completely at Random*, *Missing at Random*, and *Missing Not at Random*.

1. **Missing completely at random (MCAR)**: The probability of missing can be formulated by:

$$P(s = 0|y, z) = P(s = 0).$$

This means that the missing probability does not depend on any random variables.

2. **Missing at random (MAR)**: The probability of missing can be formulated by:

$$P(s = 0|y, z, x) = P(s = 0|x, y).$$

The probability of missingness only depends on the set of *observed* variables x and y , but not on the missing variable itself or unobservable characteristics.

3. **Missing not at random (MNAR):** The missing mechanism depends on the value of z itself or on unobserved variables, e.g., high-income individuals tend to not participate in surveys related to income.

There are several potential mechanisms for missing innovation data, which likely differ between different measures of innovation. Missing R&D data could arise from firms seeking to avoid giving benchmark spending numbers to competitors, managerial decisions to create information asymmetry about their R&D intensity, a simple failure to report zero R&D, or difficulties in estimate R&D spending from the costs of goods sold. Unobservable innovation in patent data could arise from failed innovation projects, managerial decisions to facilitate private trading gains, firm attempts to keep detailed blueprints of their innovation out of the public domain, or because they focus on process rather than product innovation. Understanding the mechanism or the underlying reason for the missing innovation data is an especially important component in determining or assessing methods to handle the missing data. For instance, the Heckman selection approach requires the selection of an instrument based on the nature of the missingness.

2.1. Common Approaches to Unreported Innovation

One common approach to missing innovation data is to delete or exclude firms without R&D spending or patents. Listwise deletion only uses a subsample of observations, deleting firms or firm-years that contain missing values in the z -variable, in equation (1). This leads to estimating the following regression using a subsample of the data:

$$y_i = s_i\alpha + \theta s_i z_i + s_i \varepsilon_i, \tag{2}$$

where $s_i z_i$ is now the explanatory variable and $s_i \varepsilon_i$ is the error term. The OLS (ordinary least squares) estimator is unbiased if $E(s_i \varepsilon_i z_i) = 0$, which is by $E(\varepsilon_i | z_i, s_i) = 0$. If data is *Missing Completely at Random* and z_i is exogenous, then $E(\varepsilon_i | z_i, s_i) = E(\varepsilon_i | z_i) = 0$. Thus, deletion can lead to consistent estimates in the case of *Missing Completely at Random*. However, if the selection is driven by observed or even unobserved variables, then $E(\varepsilon_i | z_i, s_i) \neq 0$ in general because ε_i can be correlated with s_i even if one controls for z_i , leading to biased estimates produced by deletion. Thus, a preliminary test to consider the potential costs of deleting firms without observable innovation is to assess whether the missing data is correlated with key variables of interest.

Another common approach to dealing with missing innovation data is to impute the missing observations using various methods, and then treat the resulting data as given for further analysis. Frequently used deterministic imputation methods impute the missing values with zeros (i.e. firms without R&D are considered as having zero innovation), with the industry average level of innovation, or with fitted values based on some pre-specified model. The validity of this method depends on whether the specified imputation models are correct. If the imputation model perfectly coincides with the missing mechanism, then the resulting coefficient estimate using the imputed sample is consistent. The misspecification of a deterministic imputation model can lead to biased estimates because of the distortion of the variance-covariance matrices. This provides testable implications for analyzing missing R&D and patents, namely whether firms with unreported innovation have positive values of R&D or patentable innovations.

2.2. Alternative Methods for Handling Missing Innovation Data

Two other approaches to handling missing data are also viable candidates for unreported innovation. Inverse probability weighting (IPW) relies on assigning different weights to observed

points depending on their probability of being observed. As this probability is unknown for unreported innovation, we can estimate it using binary choice models, such as logit or probit, or with a nonparametric model.

The second approach is multiple imputation (MI). MI is essentially an iterative version of stochastic imputation, which aims at explicitly modeling the uncertainty/variability ignored by the deterministic imputation procedures. Instead of replacing with a single value (unrelated to other covariates/observed data), multiple imputation uses the (joint) distribution of the observed data to estimate the parameters of interest multiple times to capture the uncertainty/variability in the imputation procedure (see Internet Appendix I). MI methods and Heckman-type approaches to deal with unreported innovation arise from different assumptions about the nature of the missing data. In empirical studies, researchers face a tradeoff between the assumptions that underpin MI versus the assumptions about the exogeneity of the instruments used in Heckman models.

In our analysis, we investigate the relative performance of deleting firms without observable R&D or patents, common deterministic imputation methods to replace unobservable R&D or patents with zero or industry mean, Heckman, inverse probability weighting, and multiple imputation as different approaches to handle unreported innovation.

3. The Severity of Unreported Innovation

3.1 Data and Sample

The sample of patents is derived from the EPO-OECD-PATSTAT database. This database, also known as the EPO Worldwide Patent Statistical Database, contains a snapshot of the European Patent Office (EPO) master documentation database with worldwide coverage. It has more than 20 tables with bibliographic data, citations, and family links for about 70 million applications from more than 90 countries, including the EPO and the USPTO.

Our sample selection begins with the October 2013 version of the PATSTAT data. It contains 44,730,405 observations, including patentees who are individuals, governmental institution/universities, and companies for the sample period of 1999–2012.² Our analysis relies on the registered names on the original patent applications, rather than the ultimate patent owners, to better capture the entities that performed the innovation activities. We merge the patent data with all publicly-listed firms in the Compustat North America and Compustat Global database for 32 countries. Our matching algorithm consists of two main steps. First, we standardize patent assignee names and firm names, focusing on unifying suffixes and dampening the non-informative parts of firm names. Second, we apply multiple fuzzy string-matching techniques to identify the firm, if any, to which each patent belongs. We randomly selected firms to manually confirm the matching of patents to firms.

We focus on countries with at least 100 publicly-listed firms (excluding Hungary, Iceland, and Ireland).³ Thus, our primary sample contains 29 countries: Australia, Austria, Belgium, Brazil, Canada, China, Denmark, Finland, France, Germany, Greece, Hong Kong, India, Israel, Italy, Japan, Korea, Malaysia, the Netherlands, New Zealand, Norway, Singapore, South Africa, Spain, Sweden, Switzerland, Taiwan, the UK, and the US. There are 30 patent offices in the sample because the EPO is a separate entity from each European country's patent office; European firms sometimes patent in their home patent office and other times with the EPO. Our baseline sample includes 333,920 firm-year observations and 37,272 unique firms, of which 5,374 are cross-listed firms. All accounting variables are from Compustat (North America and Global) and are defined in Panel A of Table A1 in the Appendix.

² Our patent sample ends in 2012, because patents post 2012 may be affected by the truncation bias for citations. The truncation bias arises due to patents after 2012 not having enough time to receive citations and result in fewer citations in comparison to earlier patents (Hall et al., 2001).

³ Relaxing this 100-firm constraint or using a 1,000-firm constraint leads to similar inferences (see Table IA1 in Internet Appendix III).

Panel A in Table 1 reports the basic descriptive statistics of our sample firms. Only 35% of the observations in our sample report any information on R&D. Of those reporting R&D expenditures (118,264), 93% report positive R&D with an average R&D expenditure of 8% of their total assets. 7% of firms report zero R&D. The 75th percentile of R&D expenditures captures firms where R&D equates to roughly 6% of total assets. In addition, the sample firms invested an average of 6% of total assets in capital expenditure. Firms have an average of 9 patent applications, 4 patents granted, and 23 citations over the sample period.⁴ On average, firms are profitable with an average ROA (return on assets) of 1% (median of 5%) and are highly levered with median leverage of 52%. In our analysis, we focus on patent applications, as these capture the R&D activity happening around the firm, but find similar results using patents granted.

We apply the Adaptive Lasso procedure to identify any additional variables to those in Table 1 that may be relevant for the prediction of unreported R&D. Following the innovation-theory based work of Reeb and Zhao (2020), we use a ten-fold cross-validation and choose the two tuning parameters (lambda and gamma) to minimize the mean square error in the out-of-sample testing (Hui et al., 2015) to a set of 37 variables. Similarly to Reeb and Zhao (2020), the Lasso approach identifies total assets, stock liquidity, and industry patent intensity as relevant predictive variables for unreported R&D and total assets, stock liquidity, industry patent intensity, and R&D stock as relevant predictive variables for unreported patents.⁵

3.2. Univariate Comparison

⁴ The average time between filing a patent application and a patent being granted across different patent offices ranges between 2 and 4 years.

⁵ In untabulated analysis, we also use stepwise regressions to identify the relevant predictive variables for unreported innovation. The obtained relevant predictive variables are the same as those identified using the Lasso approach.

To better gauge the severity of the missing data problem and the potential impact of deleting firms without reported innovation, we compare samples with and without these firms. Specifically, we evaluate the effects of deleting innovation measures by comparing two approaches: deleting all observations without both R&D and patent applications and deleting all observations without either R&D or patent applications. The first group comprises only observations that have both reported R&D expenditures *and* patent information. Our benchmark group comprises of observations that have *either* reported R&D expenditures *or* patent applications with any of the 30 patent offices, R&D and patents. We conduct a univariate comparison under different samples.

Panel B in Table 1 reports the univariate characteristics of the full sample (Column 1), the sample that reports only R&D (Column 2), the sample that reports only patents (Column 3), and the sample with both R&D and patents (Column 4). Panel B shows that deleting missing innovation data substantially reduces the number of observations and paints a very different picture in comparison to the full sample. The samples with reported R&D or patents have less than a third of the observations of the full sample. These subsamples have higher total assets than the full sample, while the rest of the variables are significantly lower (Columns 5 and 6). The R&D and patent-only sample consists of 53,456 observations. Total assets, Tobin's Q , and sales growth are larger than those in the full sample, while the rest of the variables are smaller (Column 7). It is worth pointing out that ROA decreases by 400% from the full sample to the R&D and patenting sample. These results indicate that R&D and patenting are at least not *missing completely at random* and may depend on observables.

3.3. Tests of the Deletion Assumptions

Next, we evaluate the validity of the assumptions underlying the common practices of deleting missing innovation and replacement with zero. An example of an MCAR process (when

deletion of observations with innovation is valid) is one in which firms decide whether to report innovation based on coin flips. We test the underlying assumption behind deletion, where the estimates of interest are consistent, in two ways. First, we use the MCAR test of Little (1988) to investigate the missing-value pattern. Second, we study if unreported innovation is more prevalent across firms with certain firm characteristics, by examining the predictability of unreported innovation through regression analysis.

Whether missing data is MCAR can be tested by investigating if there are significant differences between the means of different missing-value patterns across variables of interest. This is formalized by Little (1988), who implements the Chi-square test of MCAR for multivariate quantitative data. The test statistic takes a form similar to the likelihood-ratio statistic for multivariate normal data and is asymptotically χ^2 distributed under the null hypothesis that there are no differences between the means of different missing-value patterns. Rejection of the null provides evidence that the missing data are not MCAR.

Table 2 reports Little's MCAR test statistics for unreported R&D and the number of patents with different covariates. All p -values for various specifications are smaller than 0.01 with the χ^2 statistic ranging between 297 and 22,889 for both the global and US sample, rejecting the null hypothesis that unreported R&D and non-patenting firms are unpredictable. The test provides strong evidence that unreported innovation is not MCAR.

3.3.1 Predicting Missing Innovation

Next, we investigate whether the observed variation in unreported R&D and patents at the firm-year level is systematically related to firm characteristics.⁶ We assess the existence of identifiable patterns in unreported innovation by conducting a regression analysis of unreported innovation on observable firm characteristics at the firm-year level for international and US firms. The chosen characteristics are based on innovation theory and the Lasso approach described above. Note that these tests do not seek to establish causality, but rather to emphasize association and predictability in the variation in unreported innovation to shed light on the nature of missingness in innovation. We estimate a panel regression model with year, industry, and country fixed effects, separately for unreported R&D and patents.⁷

In Table 3, the dependent variable is unreported R&D, which is equal to 1 when R&D is not reported and zero otherwise. For all firms, firm characteristics with country, industry, and year fixed effects explain up to 38% of the variation in unreported R&D (Column 3 of Table 3). Firm characteristics with firm and year fixed effects explain 81% of the variation in unreported R&D (Column 4). Unreported R&D increases at the firm level with property, plant and equipment (PPE) investment, ROA and sales growth, while it decreases with total assets, stock liquidity, and industry patent intensity. For US firms (Columns 5-7), industry and year fixed effects explain 53% of the variation in unreported R&D, while firm and year fixed effects explain 93% of this variation. Unreported R&D increases with PPE, ROA, and leverage, while it decreases with industry patent intensity for US firms.

⁶ Cross-country regressions show that percentage of firms with unreported R&D and patents is predictable with macroeconomic variables, including economic openness, manufacturing intensity, government subsidies, labor regulations, intellectual property rights, university ties, skilled labor, honesty, regulatory efficacy, and Commonwealth countries.

⁷ We report the results using least square estimation because it allows us to easily incorporate multi-level fixed effects. We also estimate the determinants of unreported innovation using binary choice models, logit and probit, with various specifications of fixed effects. The results remain qualitatively the same and are available upon request.

Table 4 presents the prediction results for unreported patents, which is set equal to 1 when a firm does not file for a USPTO patent in a given year, and zero otherwise. USPTO patents are the benchmark in this analysis, as they are widely used to measure both US and non-US firm innovation. Firm, industry, and country characteristics explain up to 26% of the variation in unreported patents (Columns 1-3). Unreported patents increase at the firm level with ROA and sales growth and decrease with total assets. Focusing on just the subset of US firms, firm and industry characteristics explain a substantial amount of the variation in unreported patents (32% to 77%; Columns 6-7). Unreported patents increase with ROA and sales growth, while decreasing with total assets and stock liquidity for US firms.

Collectively, the evidence in this section indicates a significant correlation between unreported innovation (R&D and patents) and firm-specific factors. Thus, the result appears to be inconsistent with innovation *missing completely at random*.⁸

4. Unique Setting to Investigate Imputing Unreported R&D

4.1. The Setting

The main challenge to imputation approaches relates to how close are the imputed estimates to the true yet unobservable values. In this section, we adopt an innovative approach that partially overcomes the unobservable true value problem to examine the efficacy of the various common methods used to handle missing R&D in studies of corporate innovation.

Except for the first year of operation, firms are required to disclose their prior-year financial numbers on their financial statements to enable across time comparisons by users of general-purpose financial statements (Statement of Financial Accounting Concepts No. 8). This enables us

⁸ The prediction relation remains statistically and economically strong when using either contemporaneous or lagged explanatory variables. These results are available from the authors upon request. Results in Table IA2 in Internet Appendix III show that focusing only on the Lasso inferred variables does not qualitatively affect the inference.

to identify a unique (albeit narrower) setting where we can “recover” the previously unreported R&D expenditure information that serves as the true (yet previously unobservable) value. Specifically, when firms switch from not reporting to reporting R&D expenditures, they are required to report both the current year and prior year R&D expenditure amounts. In this instance, we can identify the previously unreported R&D expenditures. Our unique setting is thus especially appropriate to investigate how close the imputed estimates from various imputation methods are to the “recovered” true values.

Using the sample of US firms for the period 1992 to 2016, we identify firms that switch between reporting and not reporting R&D expenditure.⁹ We find 738 unique firms that switch between reporting and not reporting R&D. We then manually collect data from the annual reports (10Ks) of these firms on their prior years’ R&D expenditure, collecting information on the reported R&D in the year of the switch and up to two years prior to the switch in reporting. We restrict our analysis to firms without any major corporate events (e.g., merger and acquisitions) over the past two years that would have altered the underlying business operations of the firm (e.g., Bena and Li, 2014).¹⁰ We denote these as “Recovered R&D” firms. This provides us with 763 observations for the switch year (some firms switch between reporting and not reporting R&D more than once during our sample period) and 1,032 recovered observations (some firms report amounts for one year, while others for two years before the switch).

4.2 Comparing Recovered R&D Firms to Zero

We begin our analysis by comparing the characteristics of Recovered R&D firms with zero

⁹ Our initial analysis uses a window which covers 1999-2012 due to data limitations for pre-1999 international data. Our PATSTAT sample ends in 2012 and determines the end of the main sample. In tests focusing strictly on US firms, we use a longer sample period (1992-2016).

¹⁰ This is to ensure that the prior year figures disclosed in the switch year reflect only business operations that existed in the prior year 10K filings where R&D spending was not reported.

R&D firms and positive R&D firms. Panel A of Table 5 presents the results. Panel A shows that the average R&D investment for the switching firms with Recovered R&D is \$6.69 million a year and compares the “Recovered R&D” firms to firms that report zero R&D and positive R&D for the comparative years (t-1 and t-2). The R&D expenditure and R&D value of recovered firms are statistically different from the zero R&D firms. In addition, the recovered firms differ from zero R&D firms across several different dimensions like total assets, PPE, and leverage. The R&D absolute investment for recovered firms is significantly lower than positive R&D firms, but the R&D expenditure of the two groups are not distinguishable from each other. Recovered R&D firms also differ from positive R&D firms in total assets and PPE. Untabulated multivariate tests provide similar inferences, showing that recovered R&D is predictable by many common firm characteristics. Overall, results in Panel A of Table 5 show that unreported R&D expenditure firms differ from both zero R&D firms and positive R&D firms, suggesting that deleting them or classifying them as zero innovators is problematic. More specifically, if an innovation covariate is correlated with any of the variables predicting recovered R&D firms, then excluding or classifying these firms as zero innovators can lead to biased inferences.

4.3. Comparing Different Imputation Methods

Potential methods of handling missing data are listwise deletion, imputation with zero or industry mean, and multiple imputation.¹¹ We test the different imputation techniques using the “Recovered R&D” sample as a counterfactual for the true R&D in Panel B of Table 5. In the Compustat data, this recovered R&D appears as missing, and we impute this R&D with zero, with

¹¹ See Internet Appendix I for a detailed exposition of all the methods.

average industry R&D (two-digits), and with multiple imputed R&D. We compare the recovered R&D with the imputed R&D and calculate the difference and related t-statistics.

We use two samples for the R&D multiple imputation. First, we base our multiple imputation estimation on the whole US sample for the period 1992 to 2016 (not just the recovered R&D sample), *MI Full Sample*. Second, we base the multiple imputation only on the sample of recovered R&D firms and firms matched within industry and size quartile, *MI Sub Sample*. Our imputation is based on the three Adaptive Lasso determined variables: total assets, stock liquidity, and industry patent intensity and estimated by industry (two-digit). We use 200 iterations for the imputation in the analysis. If one were to be sceptical about only using Lasso-based variables and missing important variables related to corporate finance regressions of interest, Table IA3 in Internet Appendix III presents the results for MI using ROA, PPE, sales growth, leverage, lagged R&D expenditure at the firm level, as conditioning information on their own and in addition to the Lasso variables. The results remain quantitatively similar.

Panel B of Table 5 shows that recovered R&D is statistically different from zero, i.e. replacing with zero underestimates the recovered R&D values. In terms of the dollar amount of R&D, the average recovered R&D is \$6.91 million, while imputing with the industry average, gives an estimate of \$77.86 million. The industry average imputed value is over 10 times their actual R&D spending and significantly different from the recovered value. On the other hand, the two multiple imputation methods generate an average of \$6.36 million and \$8.66 million, which are not statistically different from the recovered R&D values. The relatively large variance in the MI values points to the difficulty in using these as exact point estimates for innovation in firms with missing R&D, even though it could provide less biased coefficient estimates in OLS models.

Panel C of Table 5 compares multiple imputation with an alternative innovation benchmark, text-based innovation (Bellstam, Cookson, and Bhagat, 2020). Text-based innovation

measures firm innovation using Latent Dirichlet Allocation (LDA) analysis for analyst reports of the S&P500 firms. This measure is standardized with mean 0 and variance of one, and therefore we can effectively only compare the rank correlation between multiple imputation and the text-based innovation measure. The results in Panel C of Table 5 show that there is a strong and statistically significant rank correlation, 30%, between MI and text-based innovation. Multiple imputed R&D is actually more highly correlated with the text based innovation measure than USPTO patents, global patents, or reported R&D.

Overall, results in Table 5 show that firms with recovered R&D differ from firms that explicitly report zero R&D, they are not similar to the average firm in the industry, and multiple imputation provides the closest imputation to the true value of their R&D investment. Although this design provides a sharp setting to generate the above-mentioned insights, it may not be fully representative of the broader set of firms with missing innovation data. For a more representative set of firms, we focus on the relation between MI and text-based innovation for S&P500 firms, which is large and positive. However, to provide a broader and more comprehensive analysis, in the next section, we turn to two simulation-based analyses to alleviate this concern.

5. Simulation Analysis

We consider two simulation studies, one based on the empirical distribution of Compustat (US) data and one on simulated data, to compare different methods of dealing with missing data in various data generating processes (DGPs). The first approach mimics current empirical exercises involving R&D. The second approach allows us to determine the distribution of all variables and their correlations and to examine the performance of methods in a well-controlled environment.

In both cases, we compare six methods to handle missing values. First, we consider listwise deletion. Second, we impute the missing R&D expenditure by zeros (ImpZero). Third, we impute

the missing R&D by the industry average (ImpMean). Specifically, if an observation of firm i at time t is missing, we impute the missing observation by the industry average (two-digit SIC code) for the firm in the same year. For both ImpZero and ImpMean, we also include a dummy variable indicating missingness as an explanatory variable. Fourth, we use Heckman's two-stage procedure with the selection variables containing all the observed covariates W . Heckman's procedure first predicts firms' selection probabilities by W , then corrects the selection bias by including a transformation of these predicted probabilities as an additional explanatory variable. Next, we consider the inverse probability weighting method (IPW) that weights each i -th observed point by the inverse of its conditional selection probability in least square estimation. We use the standard package of Heckman's procedure and IPW in STATA. Finally, we consider multiple imputation (MI). Since the variables generating the missingness are not known a priori, we use all observables including the outcome variable as selection variables in the imputation model.¹² To implement MI, we use 200 imputations based on a Markov chain Monte Carlo (MCMC) procedure and employ a multivariate normal regression for each imputation.

We evaluate the performance of the six methods with two criteria: The bias (B) and root mean squared error (RMSE) of coefficient estimates of the main regression. In particular, let θ be the coefficient vector of the main regression of interest. We calculate the bias and RMSE of the estimate $\hat{\theta}$ respectively, by:

$$B(\hat{\theta}) = \frac{1}{R} \sum_{r=1}^R |\hat{\theta}^r - \theta^0| / \theta^0 \quad (3a)$$

$$RMSE(\hat{\theta}) = \left[\left[\frac{1}{R} \sum_{r=1}^R (\hat{\theta}^r - \theta^0) \right]^2 + var(\hat{\theta}^r) \right]^{1/2}, \quad (3b)$$

¹² Multiple imputation draws the missing variable from a joint (predictive) distribution of observables for multiple times, and thus the set of observables should include all variables that are potentially correlated with the missing variable. Since sales growth is correlated with R&D, we also include it in the imputation model. Ignoring sales growth in imputing unreported R&D leads to incomplete conditioning of observables and biased estimates in the regression of interest (Moons et al., 2006; Sterne et al., 2009; Bartlett et al., 2011).

where $\hat{\theta}^r$ is the estimate in the r -th replication, $var(\hat{\theta}^r)$ is the estimated robust variance of $\hat{\theta}^r$, θ^0 is the true value of the parameter, R is the number of simulations. Note that we present the bias as a proportion of the benchmark θ^0 to compare across coefficients, thus one cannot use the reported bias (B) to calculate the RMSE. Based on the observed missingness of R&D and patents in the US, we consider two levels of missingness relevant for innovation variables: 50% and 70%. We perform 500 simulations.

5.1. Empirical Distribution-Based Simulation

For the empirical distribution-based simulation, we begin with a panel sample of 783 firms in Compustat over the period 1992-2012, where we have non-missing information on all financial variables of interest, except for R&D. The data include: natural log of total assets (A), leverage (L), intangible assets (I), Tobin's Q (Q), return on assets (R), R&D expenditure (RD), liquidity (V), industry patent intensity (PI), and sales growth (S). To investigate the effects of R&D missingness on the coefficient estimates of our evaluation model and how different methods of handling missing R&D perform, we generate R&D expenditure with missing observations that incorporate the three types of missingness. The resulting estimated coefficients ($\hat{\theta}^r$) under each condition are used to calculate the bias and RMSE per Eqns. (3a) and (3b). Our baseline regression uses simulated sales growth as the dependent variable and R&D expenditures, the natural log of total assets, Tobin's Q, leverage, and return on assets as explanatory variables. This approach enables us to obtain a clean set of benchmark coefficients that are free from researcher intervention except for the balanced, non-missing data criteria. Next, we describe the data generating process (DGP) for i) R&D expenditure with missing observations, and ii) the outcome variable of interest, sales growth.

5.1.1. Generate Missing R&D Expenditure

To simulate the missing R&D, we employ a subsample of complete balanced panel data, without missingness in R&D, that contains 311 firms over 21 years from 1992 to 2012. A clear advantage of this approach is that we do not need to make assumptions or estimate the conditional distribution of the R&D given that it is not missing, which is typically difficult to obtain. More importantly, it allows us to introduce the three types of missingness more precisely into the data as described below while providing us with “true” values as benchmark cases. We generate a missing indicator for R&D, denoted by M that equals 1 if R&D is missing and 0 otherwise. Once we model and assign missing R&D observations, we can obtain the simulated R&D since the data are complete, and the non-missing observations are given by their original values.

To create a missing indicator for R&D, we consider the three missing mechanisms: *Missing completely at random*, *missing at random*, and *missing not at random*. Let α_i be the individual firm effects and denote $\tilde{\eta}_{it}$ as an idiosyncratic error. The three missing patterns can be summarized by:

- Missing completely at random: $M_{it} = \tilde{\eta}_{it}$, (4)

- Missing at random: $M_{it} = \alpha_i + \beta'_O X_{it}^O + \tilde{\eta}_{it}$, (5)

- Missing not at random: $M_{it} = \alpha_i + \beta'_O X_{it}^O + \beta'_U X_{it}^U + \tilde{\eta}_{it}$, (6)

where X_{it}^O contains observed variables by researchers, while X_{it}^U is unobserved and only appears in the DGP but is omitted in imputation models. For MAR, we consider $X_{it}^O = (A_{it}, V_{it}, PI_{it})'$, which are the Lasso derived variables, and for MNAR, we add $X_{it}^U = I_{it}$ to X_{it}^O (where I_{it} represents intangible assets).¹³ To generate the missing indicator M_{it} , we need to know the true values of the parameters α_i , β_O , and β_U . However, it is well recognized that modeling binary variables is

¹³ The choice of conditioning observed variables for the missingness specification does not affect the simulation results. Table A2 in the Appendix presents robustness analysis with other specifications for X_{it}^O .

difficult in econometrics, and this is even more complicated in panel data models due to the difficulty of estimating individual fixed effects; see Lahiri and Yang (2013) for a review. To incorporate the firm fixed effects, we adopt the commonly used assumption that the firm fixed effects are correlated with the time-average of covariates in a linear manner (see Chamberlain, 1984), i.e. $\alpha_i = c + \gamma'_O \bar{X}_i^O + u_i$ in MAR and $\alpha_i = c + \gamma'_O \bar{X}_i^O + \gamma'_U \bar{X}_i^U + u_i$ in MNAR, where $\bar{X}_i^O = 1/T \sum_{t=1}^T X_{it}^O$, $\bar{X}_i^U = 1/T \sum_{t=1}^T X_{it}^U$ and u_i is the idiosyncratic noise. This assumption implies that we can incorporate the firm fixed-effects by augmenting regressions (5) and (6) by the time-series averages of covariates, respectively, as:

$$M_{it} = c + \beta'_O X_{it}^O + \gamma'_O \bar{X}_i^O + \eta_{it}, \quad (7)$$

$$M_{it} = c + \beta'_O X_{it}^O + \beta'_U X_{it}^U + \gamma'_O \bar{X}_i^O + \gamma'_U \bar{X}_i^U + \eta_{it}, \quad (8)$$

where $\eta_{it} = \tilde{\eta}_{it} + u_i$. Since there are no fixed effects in (7) and (8), we can estimate all parameters in these two models and predict M_{it} based on these estimates. Specifically, we first estimate (7) and (8), respectively, by a probit regression of the missing data indicator for R&D using the panel data sample (783 firms). We set the estimates \hat{c} , $\hat{\beta}_O$, $\hat{\beta}_U$, $\hat{\gamma}_O$, and $\hat{\gamma}_U$, as the true parameters to generate the missing probability M_{it}^* in the *complete subsample* of the data:

$$M_{it}^* = \Phi(\mathbf{p}m_{it}). \quad (9)$$

Φ is the normal CDF function and $\mathbf{p}m_{it}$ is obtained for the three scenarios by:

1. Missing completely at random: $\mathbf{p}m_{it} = \eta_{it}$, (10)

2. Missing at random: $\mathbf{p}m_{it} = \hat{c} + \hat{\beta}_O X_{it}^O + \hat{\gamma}_O \bar{X}_i^O + \eta_{it}$, (11)

3. Missing not at random: $\mathbf{p}m_{it} = \hat{c} + \hat{\beta}_O X_{it}^O + \hat{\beta}_U X_{it}^U + \hat{\gamma}_O \bar{X}_i^O + \hat{\gamma}_U \bar{X}_i^U + \eta_{it}$, (12)

where $\eta_{it} \sim IID N(0, \sigma_\eta^2)$ and $\sigma_\eta^2 = 0.15$ based on the empirical distribution of the error term.

Once we obtain M_{it}^* , we set the (i,t)-th observation of R&D as missing ($M_{it}^* = 1$) depending on

$M_{it}^* > Q_\tau(M_{it}^*)$, where $Q_\tau(M_{it}^*)$ is the τ -th quantile of M_{it}^* , and τ controls the percentage of missing.

5.1.2 Generating Sales Growth

We simulate the outcome variable of interest, i.e. sales growth S , because observable growth is potentially influenced by variables omitted from our empirical specification. We want to isolate the impact of missing innovation data from the errors from omitted variables in our regression of sales growth on innovation.¹⁴ We generate S in the complete subsample without any missingness (311 firms over 21 years). The DGP of S is based on the following model:

$$S_{it} = \mu_i + \delta'RD_{it} + \theta'Z_{it} + \varepsilon_{it}, \quad (13)$$

where μ_i is firm fixed effects, Z_{it} contains the determinants of sales growth, $Z_{it} = \{A_{it}, Q_{it}, R_{it}, L_{it}\}'$, and ε_{it} is the error term. Note that intangible assets are not observed and thus also not included in the DGP of S . The firm fixed effects μ_i are generated by $\mu_i = 0.1\iota'\bar{Z}_i$, where ι is a 4×1 vector of ones and $\bar{Z}_i = 1/T \sum_{t=1}^T Z_{it}$, and thus μ_i is correlated with sales growth determinants. To obtain the parameters for δ' and θ' , we estimate (13) using the same complete subsample without missingness and fix the estimated values in the simulation. To allow the idiosyncratic error to be correlated with selection instruments, we generate $\varepsilon_{it} = \tilde{\varepsilon}_{it} + \bar{Q}_i$ in MAR and $\varepsilon_{it} = \tilde{\varepsilon}_{it} + 0.5(\bar{Q}_i + \bar{I}_i)$ in MNAR. Here \bar{Q}_i and \bar{I}_i are the time average of Tobin's Q and intangible assets for firm i , respectively, which drive the missingness of R&D as discussed in Section 5.1.1.

¹⁴ We use simulated sales growth rather than observed sales growth in the benchmarking exercise because it allows us to explicitly compare the estimated coefficients to the true values. In contrast, using observed sales growth in our tests allows bias from two sources: imputation bias and misspecification bias (e.g., omitted variables), rendering the comparison between various imputation methods less clear.

$\tilde{\varepsilon}_{it} \sim IID N(0, \sigma_{\varepsilon}^2)$ and $\sigma_{\varepsilon}^2 = 0.18$ based on the empirical distribution of the residual from estimating equation (13).

5.1.3 Simulation Results

Table 6 reports the simulation results under three missing mechanisms and two levels of missingness (50% and 70%). When R&D is missing completely at random, we find that both bias and RMSE increase with increasing missingness in R&D (Panel A). All methods show a relatively small bias under MCAR, except IPW and Heckman. IPW and Heckman, typically do not include fixed effects due to the difficulty in estimating fixed effects in binary model settings, which potentially explains part of their relatively poor performance (we use the standard packages in STATA for these two methods). In *missing completely at random*, multiple imputation has the lowest bias. Multiple imputation exhibits relatively smaller RMSE than other methods too. Deterministic imputation methods (ImpZero and ImpMean) generate double the bias in multiple imputations and RMSEs that are similar to MI. Still, MI has both the lowest average bias and RMSE under MCAR.

Panel B shows the results for MAR, where the bias of all methods increases from MCAR. Under MAR, all methods lead to biased estimates, not only for R&D (which has missing observations), but also for the other explanatory variables that do not have any missingness. MI on average produces the lowest bias across all of six methods followed by ImpMean and ImpZero. The average absolute bias in listwise deletion is over ten times greater than the bias in multiple imputation, while bias in IPW and Heckman are over 170 times and 80 times greater than MI. The common imputation methods, on average, exhibit similar RMSEs, where ImpZero, ImpMean, and MI have the lowest RMSEs. Panel C shows the results when missingness is driven by unobservables (MNAR). Under MNAR, MI continues to produce the lowest bias among all six

methods followed by ImpZero and ImpMean. The bias in LD is six times larger than the bias in MI. Focusing on RMSE, once again ImpZero, ImpMean, and MI all exhibit similar magnitudes.

For robustness, we also investigate other settings for both the determinants of R&D missingness in equations (11) and (12), as well as the sales growth DGP in (13). We introduce a larger set of conditioning variables in equations (11) and (12) in Panel A of Table A2, a larger set of conditioning variables in equation (13) in Panel B of Table A2, and Lasso as a variable selection procedure to determine which covariates should be included in the MI models in Panel C of Table A2. In the last setting, we consider a double Lasso procedure that applies Lasso to both R&D expenditure and sales growth regressions.¹⁵ The double Lasso procedure is theoretically justified by Belloni et al. (2014), and it allows us to select variables that are correlated with both R&D expenditure and sales growth for accurate imputation. The results in Table A2 are quantitatively similar to those in Table 6. MI results in the lowest average bias consistently across methods.

Our simulations focus on two separate levels of R&D missingness, namely 50% and 70%. However, our cross-country sample, which underlies Tables 1 to 4 reveals that the level of missingness varies by country. Specifically, the rate of missing R&D data ranges from 5% missing in Japan to 85% missing in Italy. Consequently, we repeat the simulation analysis across a wide selection of missingness levels in 5% increments. Figure 2 shows the relative bias in the R&D coefficient estimate in using multiple imputation and listwise deletion as the rate of missing R&D increases from 5% to 85% for MAR. Across the entire range of missing R&D, multiple imputation exhibits substantially lower bias in the R&D coefficient estimate relative to listwise deletion.

¹⁵ In particular, we first employ Lasso to select covariates in the model regressing R&D expenditure on all available covariates: Tobin's q, total assets, leverage, ROA, liquidity, industry patent intensity, and their time-series averages for each firm. We denote the selected covariates as X_{RD} . Next, we use Lasso to select the covariates in the model regressing sales growth on all available covariates specified above and obtain the selected covariates denoted as X_{SG} . We use the union of X_{RD} and X_{SG} as variables in multiple imputation.

It is worth noting that these results constitute a lower bound on bias generated by LD and deterministic imputation methods for two reasons. First, we include all the missingness determinants (A, V, and PI) as control variables, which implies that even if one knows the missingness mechanism and correctly controls for it, the estimated coefficients are still biased. Second, we have assumed that the errors in the sales growth and the selection regressions are not correlated, which is most likely not the case in reality. In untabulated results, we show that if the errors of the two regressions are correlated, then the bias of deletion and deterministic imputation increases.

The analysis of simulations based on the empirical distribution, albeit realistic and informative, does not allow us to clearly infer how the correlation between variables, which might differ across data samples, influences the performance of methods. Therefore, in the next subsection, we conduct simulation analysis using generated data, where we can precisely specify the correlation among errors and compare the magnitude of the effects of various methods in a well-controlled environment.

5.2 Simulation with Generated Data

5.2.1 Data Generating Process

We generate the dependent variable of interest as follows:

$$Y_i = z_{1i}\theta_1 + z_{2i}\theta_2 + \varepsilon_i, \quad i = 1, \dots, N, \quad (14)$$

where $\theta_1 = \theta_2 = 1$, $\varepsilon_i \sim IID N(0,1)$, and $N=1,000$. The two covariates z_{1i} and z_{2i} are generated by a multivariate normal distribution with unit means and variance-covariance matrix specified later. z_{1i} contains missing observations, while z_{2i} is completely observed. Let M_i be the missing indicator of x_{1i} that equals 1 if x_{1i} is missing and 0 otherwise, which is determined by $M_i =$

$1[M_i^* > Q_\tau(M^*)]$, where $1[\cdot]$ is an indicator function, M_i^* is a latent variable, $Q_\tau(M^*)$ is the τ -th quantile of M^* . We consider two values of $Q_\tau(M^*)$, 0.7 and 0.5, which correspond to 70% and 50% of missing observations in x_{1i} , respectively. We consider three missing mechanisms for x_{1i} :

- Missing completely at random: $M_i^* = \eta_i$,
- Missing at random: $M_i^* = x_{1i}\gamma_1 + x_{2i}\gamma_2 + \eta_i$,
- Missing not at random: $M_i^* = x_{1i}\gamma_1 + x_{2i}\gamma_2 + x_{3i}\gamma_3 + x_{4i}\gamma_4 + \eta_i$.

We set $\{\gamma_1, \gamma_2, \gamma_3, \gamma_4\} = \{2, 1, 1, 1\}$. x_{1i} and x_{2i} are observed covariates that drive the missing pattern, while x_{3i} and x_{4i} are unobserved. η_i is the error term, independently generated from $N(0,1)$ in MCAR, but correlated with ε_i in MAR and MNAR. We consider various patterns of correlations between the generated variables. In the benchmark case, we set the covariance matrix for the multivariate normally distributed $\{z_1, z_2, x_1, x_2, x_3, x_4, \varepsilon, \eta\}$ as:

$$\begin{pmatrix} 1 & & & & & & & & \\ 0.4 & 1 & & & & & & & \\ 0.5 & 0.4 & 1 & & & & & & \\ 0.4 & 0.4 & -0.2 & 1 & & & & & \\ 0.2 & 0.1 & 0.2 & 0.3 & 1 & & & & \\ 0.1 & 0.2 & 0.1 & 0.1 & 0.1 & 1 & & & \\ 0 & 0 & 0 & 0 & 0 & 0.4 & 1 & & \\ 0 & 0 & 0 & 0 & 0 & 0 & 0.4 & 1 & \end{pmatrix}.$$

Note that all covariates $\{z_1, z_2, x_1, x_2, x_3, x_4, \varepsilon, \eta\}$ are correlated with each other. The two error terms are correlated with each other, but they are independent of the observed covariates. In MNAR, the missingness is also driven by two unobservables, which may be correlated with the errors. Hence, the unobserved selection variable x_3 is uncorrelated with both errors, and x_4 correlated with ε .¹⁶

¹⁶ We tried different parameters and generated different densities, reaching quantitatively similar conclusions. We also considered alternative specifications of the covariance matrix to investigate how the correlation between variables affects the performance of different methods. The results are available upon request from the authors.

5.2.2 Results

Table A3 in the Appendix presents the results for the simulation with generated data. As in the previous simulation, both bias and RMSE increase with missingness in z_1 and from MCAR to MNAR. We first focus on the results with 70% missingness. Panel A shows the results for MCAR, where all methods produce negligible bias and small RMSE for both explanatory variables. There are only marginal differences across LD, Heckman, IPW and MI. However, the two deterministic imputation methods (using zero and industry average) produce the largest biases and RMSE.

Panel B shows the results for MAR, where listwise deletion exhibits a substantial sample selection bias, and both coefficient estimates are downward biased at 15% and 12% respectively for θ_1 and θ_2 . Imputation using zeros or industry means increase the bias in θ_1 from -19% to -23% but decreases the bias in θ_2 from 28% to 12%. The Heckman procedure exhibits among the smallest bias that is comparable to MI, but at the cost of variance. This is reflected in the large RMSE of the Heckman estimates, suggesting that the two-step procedure is rather inefficient. On the contrary, MI performs well, despite the increase in biasness, it continues to have among the lowest bias. The bias of MI is around half as large as that of listwise deletion, and almost two times smaller than that of imputation using zeros or means, and MI has the smallest RMSE.

For MNAR, all methods produce biased estimates due to the non-random missing pattern, but the degree of bias differs substantially across methods (Panel C). Imputation using zeros or mean leads to the largest bias and RMSE for θ_1 among the six methods (θ_1 bias is around 28%); while the bias and RMSE for θ_2 are generally in the middle of the six methods (θ_2 bias is 14%). The biasness in LD and IPW methods also deteriorates in comparison to the MAR setting, leading to more than 17% and 15% downward bias in θ_1 and θ_2 respectively. Both the bias and RMSE

for the Heckman procedure deteriorated by 62.5% in comparison to MAR (the largest deterioration among the six methods). Despite the observed deterioration compared to MAR, MI continues to produce the lowest bias and RMSE among the six methods.

In general, all six methods show lower bias and RMSE at lower levels of missingness (50%) than at higher levels of missingness (70%). As the level of missingness varies across different data sets, the relative efficacy of the methods we investigated could differ. Replicating our analysis across different levels of missingness reveals that multiple imputation consistently exhibits the smallest bias. In contrast, the bias from listwise deletion increases substantially with the level of missingness.

6. Impact of Bias on Inferences through Replication

Our analysis so far has conceptually demonstrated the problems with the various methods of handling unreported R&D. Yet, we do not know the economic significance of the effect of the treatment of unreported R&D for empirical studies. This is a cross disciplinary problem and we use the analysis of Fama and French (2002, FF02 hereafter) to assess the economic significance of the effect of different methods of handling unreported R&D on economic inference. One of the most important issues in corporate finance is understanding how firms chose their capital structure. The two prevailing models are the trade-off and pecking order models. Fama and French (2002) test the implications of these models for firm dividends and leverage. They report a positive relation between leverage and profitability for dividend and non-dividend paying firms. FF02 find ambiguous results on the relation between investments and leverage, as the two proxies for investment have opposite signs: market-to-book is positively correlated with leverage and R&D expenditures are negatively correlated with leverage. For expositional simplicity, we focus this analysis on comparing multiple imputation to the two most commonly used solutions to missing R&D, namely deleting firms without R&D and classifying them as zero innovators.

We replicate their sample and note that 60% of the firms in their sample do not report R&D expenditures. FF02 classifies all firms with unreported R&D as having zero R&D, and they include a dummy variable equal to one to differentiate firms with unreported R&D from firms that report zero R&D. We estimate the leverage regression (Equation 15) below to evaluate if leverage differs across firms in the manner predicted by the trade-off or pecking order model using three approaches—listwise deletion, zero imputation with a missing dummy, and multiple imputation—and compare the resulting estimates:

$$\frac{L_t}{A_t} = \beta_0 + \beta_1 \frac{V_t}{A_t} + \beta_2 \frac{ET_t}{A_t} + \beta_3 \frac{Dp_t}{A_t} + \beta_4 RDD_t + \beta_5 \frac{RD_t}{A_t} + \beta_6 \ln(A_t) + e_t. \quad (15)$$

We follow FF02 in the choice of the sample period, variables of interest, and notation. $\frac{ET_t}{A_t}$ the ratio of annual pre-interest pre-tax earnings to end-of-year total assets, is a proxy for the expected profitability of assets in place.¹⁷ $\frac{V_t}{A_t}$, the ratio of a firm's total market value to its book value, is a proxy for expected investment opportunities. $\frac{RD_t}{A_t}$, the ratio of R&D expenditures to assets, is an additional proxy for expected investment. Unreported R&D is imputed with zero. RDD_t is a dummy variable equal to 1 for unreported R&D, and zero otherwise. $\frac{Dp_t}{A_t}$, the ratio of depreciation expense to assets serves as a proxy for non-debt tax shields. $\ln(A_t)$, the natural logarithm of total assets is a proxy for volatility. The sample period is 1965-1999 as in FF02.

Table 7 replicates FF02 using a contemporaneous regression with two-way fixed effects, double clustered standard errors, and five additional treatments for unreported R&D. Panel A presents results for dividend-payer firms and Panel B for non-dividend payer firms. We use listwise deletion (*LD*), multiple imputation with only the variables in the regression (*MI*), multiple

¹⁷ *ET*, earnings before taxes, preferred dividends, and interest payments is the income that could be sheltered from corporate taxes by interest deductions. Thus $\frac{ET_t}{A_t}$ is a measure of profitability when we look for tax effects in the trade-off model.

imputation with Lasso variables volume and patent intensity (*MI Lasso*), pseudo R&D, and text-based innovation. In order to implement pseudo R&D, we still need to impute R&D for firms that do not report pseudo R&D, which is 95% of the sample. We impute R&D with zero for those cases. There are only 2,016 observations where text-based innovation is not missing in this sample, because text-based innovation is available for the period 1990 to 2012.

The various estimation techniques lead to very different estimates for β , confirming the importance of how violations of the MCAR assumption and the distortion of the variance-covariance matrix with zero imputation that cause listwise deletion and zero imputation to yield inconsistent estimates. The estimates based on zero imputation and listwise deletion reported in Columns (1) and (2) are negative and differ considerably from estimates using multiple imputation (Columns 3 and 4), which are positive. This suggests that the inferences made by the researcher in innovation could be driven by how they chose to deal with missing innovation data.

The results using zero imputation and a dummy variable show no relation between market-to-book ratio and leverage, and a marginal effect of profitability on leverage for dividend-paying firms (Column 1 Panel A). Listwise deletion leads to an insignificant relation between market-to-book ratio and leverage and between depreciation and leverage for dividend-paying firms (Column 2 Panel A). Columns (3) and (4) of Table 7 presents the results for unreported R&D imputed using multiple imputation. In this case, there is a substantial change in the magnitude of the coefficients of all the explanatory variables. Most importantly, the relation between both investment variables and leverage is positive and internally consistent. Now R&D expenditure has a positive impact on leverage and not negative, as in Columns (1) and (2), which is congruent with pecking order theory. Using alternative measures of innovation like pseudo R&D and text-based innovation leads to similar results to zero-imputation and listwise deletion. The pseudo R&D result is heavily driven by the zero imputation of the rest of the observations. The estimates in Table 7 illustrate that the

method used to handle missing R&D can lead to substantially different inferences. Using multiple imputation for missing R&D in this setting potentially explains the puzzling findings in the original FF02 study.

7. Patents

So far, we have presented analyses using unreported R&D. Yet, many studies of corporate innovation rely on patent data from the USPTO, with studies of international firms also tending to rely on this patent database. Unlike disclosure requirements for R&D expenditures, firms do not face an affirmative duty to seek patents, which potentially explains why most US firms do not submit patent applications. Consequently, we conduct a similar analysis for patents as for unreported R&D. We investigate the performance of different imputation methods using two sets of counterfactuals: first we use non-USPTO patents of US firms as counterfactuals, and second, we use the product innovation measure Mukherjee et al. (2016). We conclude with an empirical data-based simulation for patents.

7.2 Relevance of Unreported Innovation via Patents

Innovation-related studies, across accounting, economics, and finance, focus on patenting as the most important outcome of the research and development process. These studies mainly use data from the USPTO-NBER dataset. This dataset includes all firms that have applied for USPTO patents and the NBER has conducted extensive disambiguation of firm data. While this dataset has been instrumental in conducting the first pieces of research, it provides a partial view of innovation and patenting activity. For researchers interested in understanding and/or capturing a fuller extent of firm innovation activities, investigating patenting only through the USPTO will

underestimate the innovative activity of many firms. Next, we investigate the properties of patents filed outside USPTO jurisdictions to understand the importance of non-USPTO patents.

Over 14,000 US-firms applied for USPTO patents in the sample period, 9,518 US-firms received patents abroad, while 1,676 non-US firms received USPTO patents and 1,758 non-US firms received non-USPTO patents. The total number of patents granted in non-USPTO jurisdictions per year is substantial for both US and foreign firms. For instance, foreign firms are granted 20% more patents in non-USPTO jurisdictions than in the US. US firms also are granted, on average, 22 patents a year outside the US and 28.2 patents in the US. Firms without USPTO patents are typically deleted or counted as non-innovative firms when using USPTO data and generate a bias in the coverage of patenting.

We use the sample of US firms that patent abroad to investigate the different methods of handling unreported patents, similar to Panel B of Table 6. We impute observations without USPTO patents with zero, industry mean (two-digit SIC code), and multiple imputation. Multiple imputation is carried out with all the variables in Table 4, by industry (two-digit), and the Lasso variables of stock volume, patent intensity, and R&D stock. Results in Figure 3 show that US patents abroad are not equal to zero, they are different from the USPTO industry mean, but they are not statistically different from MI.

As an alternative counterfactual, we use MI to predict the number patents of firms with new major products, as in Mukherjee et al. (2016), without USPTO patent applications. Table 8 presents the comparison between two different MI methods. MI M1 includes the Lasso variables only: stock liquidity, R&D stock, and industry patent intensity using PATSTAT patents, MI M2 is the multiply imputed non-USPTO patents using the same model as M1 and $\ln(\text{total assets})$, ROA, PPE, capital expenditure, sales growth, and leverage as conditioning information.

Panel A presents the innovation characteristics of firms with various coverage of patents and new products. The majority of the sample has no patents and no new products (Column 2). These firms have the lowest R&D expenditure and lowest percentage of R&D reporting. Just over 3% of the sample has both USPTO patents and new products (Column 3), but these firms have the highest reporting rates for R&D (91%) and the largest R&D expenditures. About 6% of the sample has no USPTO patents but announce new products (Column 1). Finally 14% of the sample has patents but does not report any new products. The difference between firms with no patents and new products and firms with patents and no new products is large and statistically significant.

Panel B of Table 8 presents the single and multiple imputations for the four above categories and compares firms without patents and new products to firms with patents and no new products. Imputing the number of patents with zero or industry average is statistically different from the counterfactual of firms with patents and no new products. In contrast, imputing with multiple imputation results in patent numbers that are not statistically different from firms with patents and no new products.

7.3 Empirical Data-Based Simulation

Patents and R&D expenditures may have different determinants and missingness levels. To further understand the properties of the different methods for handling missing data in the patent setting, we replicate the empirical distribution-based simulation, with the USPTO patent data distribution. Table A4 in the Appendix presents the results of the simulation based on the patent empirical distribution. Under MAR, IPW and Heckman generate the highest biasness in coefficient estimates relative to both imputation and deletion. Focusing on MNAR, deterministic imputation and multiple imputation both perform better than listwise deletion, IPW and Heckman approaches.

8. Conclusions and Recommendations

Most public firms do not report R&D expenditures, do not obtain patents, nor receive patent citations. Studies across accounting, economics, and finance typically exclude firms without reported R&D or patent activity or classify them as zero innovators (e.g., Autor et al., 2020; Corrado et al., 2020; DeSimone et al., 2020; Koch et al., 2020). Given the wide range of approaches used in recent empirical studies, there is a substantial need for research that evaluates the relative role of various solutions to the case of missing R&D or patents, giving guidelines for future research for studies that use these variables as conditioning variables. We study how various methods of handling unreported innovation affect our inferences about corporate research and development. More specifically, we explore the assumptions underlying different methods of handling unreported innovation, assess the biases that each of these methods introduces, and provide guidance for future research.

Instead of arising randomly, we document that unreported innovation is systematically correlated with several firm, industry, and country characteristics. Accordingly, eliminating firms without R&D or patents provides biased results, if a proposed innovation covariate is correlated with any of these predictor variables (e.g., firm size, leverage, profits, etc.). Because patent prevalence is even lower than the frequency of reported R&D, concerns about biases from deleting firms without patents is especially pronounced.

Using recovered R&D, which allows us to accurately measure unreported R&D expenditures in prior unreported years, we compare different methods of handling firms without reported R&D. These recovered R&D firms do not look like zero R&D firms nor do they appear similar to positive reporting R&D firms. The recovered R&D is also statistically different from the industry average. This finding is problematic for the common methods for handling unreported innovation, namely deleting the firms, classifying them as zero innovators or setting their R&D to

the industry average and including a dummy variable. Simulation results allow us to rank the biases created from different methods of coping with firms without patents or reported R&D. For instance, replacing missing innovation with zeros underestimates true innovation and leads to biased R&D coefficient estimates (e.g., Table 6). To demonstrate the economic impact of these findings, we replicate an influential finance study (Fama and French, 2002) and explicitly show how different approaches to unreported innovation affect empirical inferences.

Innovation variables exhibit very high rates of unobservability. The most common methods to handle firms without observable innovation (R&D or patents) are excluding them (listwise deletion) and deterministic imputation (with zero or the industry mean). Our results show that unreported R&D and firms without patents are predictable and that the variables used to predict this missingness are known determinants of both innovation and other corporate outcomes of interest. Consequently, in studies that rely on the traditional methods of handling unobservable innovation, the residual in the regressions will likely be correlated with other explanatory variables. The deletion of firms with unobservable innovation and their classification as non-innovators, even after including a dummy variable, can lead to biased coefficients of not only innovation, but also other explanatory variables. These traditional methods of handling unreported innovation do not work well in addressing unreported innovation when the selection is correlated with outcome variables of interest. One of the most important takeaways from these findings is that commonly used solutions to handle unreported innovation can lead to biased parameter estimates that make prior inferences about corporate innovation difficult to assess.

It is difficult to give definitive solutions to dealing with missing innovation across different datasets, countries, and research settings. Our results using US data reveal that the two common methods to handling missing innovation data provide biased coefficient estimates and standard errors. Strikingly, across a wide range of specifications, multiple imputation exhibits the least bias

and RMSE among the six methods we investigate. Importantly, MI is the solution to unreported innovation data that is used the least in finance and economics studies. Of course, in a single industry analysis with limited numbers of positive R&D firms (e.g., Real estate renting and leasing, SIC 53) or where upward of 90% of the missing R&D firms arise from zero R&D expenditures, replacing missing with zero will provide a reasonable solution. Yet, multiple imputation still performs well in this scenario as well.

The results allow us to provide some general guidelines and recommendations for economics and finance scholars confronted with unreported innovation.

1. In studies of innovation, missing R&D and patents can arise from: i) random collection error from data providers, ii) managers not reporting R&D expenses due to zero (near zero) innovation, iii) strategic disclosure choices in reporting R&D expenses and patenting, iv) unsuccessful R&D, or v) firms filing for patents in alternative patent offices. Consequently, researchers should report both full and partial sample characteristics of the variables of interest. The level or degree of missingness of the innovation variable being used should be noted.
2. Researchers with missing innovation data should test if the missing data is predictable or MCAR. Little (1988) provides a test to determine if the data is *missing completely at random*. For Stata users, the *mcartest* command implements this test.
3. If the missing data is unpredictable or MCAR (maybe because the missing data stems from random collection errors by the data provider), then researchers could potentially delete or exclude the observations with missing data.
4. If the missing data is predictable, then researchers should attempt to predict missing innovation data using economically motivated observable variables. The predictive variables should be included as covariates in the regression and selection model. The

researcher could use multiple imputation (for Stata users the **MI** command) to handle the missing observations.

5. If the missing data is predictable and there are both observable and unobservable characteristics that lead to missing innovation data, the problem is more challenging. Schafer and Graham (2002) show that multiple imputation can often be unbiased for MNAR and MAR data even though the researcher assumes the data to be MAR. Conceptually, both Heckman correction and multiple imputation remain appealing, with both approaches involving assumptions and tradeoffs. Surprisingly, Heckman correction and Inverse Probability Weighting are the worst performers under MNAR in our simulations, with MI typically performing the best.

Overall, when missingness is beyond the researcher's control and its distribution is unknown, handling this missing data ultimately boils down to the assumptions and mechanisms of missingness. In finance studies, a researcher must often decide on the assumptions of MAR versus MNAR, which involves the use of MI versus IV respectively. Yet, the assumptions underlying both IV and MI are likely to be violated. In practice, violating the assumptions of MNAR often has only a minor impact on estimates and standard errors because the covariates included in imputation models are often correlated with the determinants of missingness (Collins et al., 2001). Consequently, one might tilt towards the use of MI for missing innovation data. Yet, we recommend reporting both MI and IV estimates, coupled with a discussion of the plausibility of the underlying assumptions in the spirit of partial identification.

In summary, the proportion and distribution of missing innovation data in the sample should be reported. Researchers should conduct an analysis of the randomness and predictability of the missing innovation data in their sample. Researchers should consider performing simulations

similar to ours, based on their own data, to choose between the various methods of handling unreported innovation.

References

- Anton, J. and D. Yao, 2004. Little patents and big secrets: Managing intellectual property, *Rand Journal of Economics* 35(1), 1-22.
- Autor, D., D. Dorn, G. Hanson, G. Pisano, and P. Shu, 2020. Foreign competition and domestic innovation: Evidence from US patents. *American Economic Review: Insights*, 2(3), 357-74.
- Bartlett, J.W., C. Frost, and J.R. Carpenter, 2011. Multiple imputation models should incorporate the outcome in the model of interest, *Brain* 134(11), 189.
- Belloni, A., V. Chernozhukov, and C. Hansen, 2014. Inference on treatment effects after selection amongst high-dimensional controls, *Review of Economic Studies* 81(2), 608-650.
- Bellstam, G., A. Cookson, and S. Bhagat, 2020. A text-based analysis of corporate innovation, *Management Science*, forthcoming.
- Bena, J., and K. Li, 2014. Corporate innovations and mergers and acquisitions, *Journal of Finance* 69(5), 1923-1960.
- Bertrand, M., E. Duflo, and S. Mullainathan, 2004. How much should we trust differences-in-differences estimates? *The Quarterly Journal of Economics* 119(1), 249–275.
- Branstetter, L., M. Drev, and N. Kwon, 2019. Get with the program: Software-driven innovation in traditional manufacturing, *Management Science* 65(2), 541-558.
- Chamberlain, G., 1984. Panel data, in *Handbook of Econometrics*, eds. Z. Griliches and M. D. Intrilligator, North-Holland Amsterdam, 1248-1318.
- Collins, L.M., J. L. Schafer, and C. Kam, 2001. A comparison of inclusive and restrictive strategies in modern missing data procedures, *Psychological Methods* 6(4), 330-351.
- Corrado, C., D. Martin, and Q. Wu, 2020. Innovation α : What do IP-intensive stock price indexes tell us about innovation? *American Economic Review* 110, 31-35.
- Croce, M., T. Nguyen, S. Raymond, and L. Schmid, 2018. Government debt and the returns to innovation, *Journal of Financial Economics* 132(2), 205-225.
- De Simone, L., J. Huang, and L. Krull, 2020. R&D and the rising foreign profitability of U.S. multinational corporations, *The Accounting Review*, forthcoming.
- Fama, E. and K. French, 2002. Testing the trade-off and pecking order predictions about dividends and debt, *Review of Financial Studies* 15(1), 1-33.
- Grangier, D. and I. Melvin, 2010, Feature set embedding for incomplete data. *Advances in Neural Information Processing Systems*.

- Hall, B. H., Jaffe, A. B., and Trajtenberg, M., 2001, The NBER patent citation data file: lessons, insights and methodological tools, Working Paper 8498, National Bureau of Economic Research.
- Hochberg, Y., C. Serrano, and R. Ziedonis, 2018. Patent collateral, investor commitment, and the market for venture lending, *Journal of Financial Economics* 130(1), 74-94.
- Hombert, J. and A. Matray, 2018. Can innovation help U.S. manufacturing firms escape import competition from China? *Journal of Finance* 73(5), 2003-2039.
- Huang, J., 2018. The customer knows best: The investment value of consumer opinions, *Journal of Financial Economics* 128(1), 164-187.
- Hui, F. K. C., D. I. Warton, and S. D. Foster, 2015. Tuning parameter selection for the adaptive Lasso using ERIC, *Journal of the American Statistical Association* 110(509), 262-269.
- Jiang, W., 2017. Have instrumental variables brought us closer to the truth, *The Review of Corporate Finance Studies* 6(2), 127-140.
- Koch, A., M. Panayides, and S. Thomas, 2020. Common ownership and competition in product markets, *Journal of Financial Economics*, forthcoming.
- Koh, P. S. and D. Reeb, 2015. Missing R&D, *Journal of Accounting and Economics* 60(1), 73-94.
- Koh, P. S., D. Reeb, and W. Zhao, 2018. CEO confidence and unreported R&D, *Management Science* 64(12), 5461-5959.
- Lahiri, K., and L. Yang, 2013. Forecasting binary outcomes, in *Handbook of Economic Forecasting* (Vol. 2B), eds. A. Timmermann and G. Elliott, Amsterdam: North-Holland, 1025-1106.
- Lerner, J. and A. Seru, 2017. The use and misuse of patent data: Issues for corporate finance and beyond, Working Paper, Harvard University.
- Little, R.J.A., 1988. A test of missing completely at random for multivariate data with missing values, *Journal of the American Statistical Association* 83(404), 1198-1202.
- Little, R.J.A. and D.B. Rubin, 2002. Statistical analysis with missing data, 2nd Edition. New York, NY: John Wiley & Sons, Inc.
- Masulis, R. and E. Zhang, 2019. How valuable are independent directors: Evidence from external distractions, *Journal of Financial Economics* 132(3), 226-256.
- Moons, K., R. Donders, T. Stijnen, and Jr F. Harrel, 2006. Using the outcome for imputation of missing predictor values was preferred, *Journal of Clinical Epidemiology* 59(10), 1092-1101.
- Mukherjee, A., M. Singh, and A. Zaldokas, 2017. Do corporate taxes hinder innovation? *Journal of Financial Economics* 124 (1), 195-221.
- Petersen, M.A., 2009. Estimating standard errors in finance panel data sets: Comparing approaches, *Review of Financial Studies* 22(1), 435-480.
- Png, I., 2017. Law and innovation: Evidence from state trade secrets laws, *Review of Economics and Statistics* 99(1), 167-179.
- Reeb, D.M. and W. Zhao, 2020. Disregarding the shoulders of the giants: Evidence from innovation research, Working Paper.

- Robins, J. and N. Wang, 2000. Inference for imputation estimators, *Biometrika* 87(1), 113-124.
- Rubin, D.B. 1976. Inference and missing data, *Biometrika* 63(3), 581–592.
- Rubin, D.B. 1987. Multiple imputation for nonresponse in surveys. New York, NY: John Wiley.
- Schafer, J. and J. Graham, 2002. Missing data: Our view of the state of the art, *Psychological Methods* 7(2), 147-177.
- Sterne, J., I. White, J. Carlin, M. Spratt, P. Royston, M. Kenward, A. Wood, and J. Carpenter, 2009. Multiple imputation for missing data in epidemiological and clinical research: Potential and pitfalls, *British Medical Journal* 338, b2393.
- Thompson, S.B., 2011. Simple formulas for standard errors that cluster by both firm and time, *Journal of Financial Economics* 99(1), 1-10.

Table 1
Sample Characteristics and Univariate Comparisons

This table shows the sample characteristics and univariate comparisons. Panel A presents the sample characteristics. The sample period is 1999–2012. Panel B shows the difference in characteristics across different deletion methods. “Full Sample” uses all available observations without deletion based on either reported R&D or patent application information. “Report R&D” includes only observations that report R&D. “Report Patent” includes only observations that patent applications in any patent office, “R&D and Patent” includes only observations that have positive R&D and patent filings in the PATSTAT. Firm-years represent the maximum number of observations available for each subsample. Variable definitions are presented in Table A1. *, **, and *** represent significance at the 10%, 5%, and 1% levels, respectively.

Panel A. Sample Characteristics

Variables	N	Mean	Median	Std. Dev.	25th	75th
	(1)	(2)	(3)	(4)	(5)	(6)
R&D Expenditure	118,264	0.08	0.02	0.60	0.00	0.06
Report R&D	333,920	0.35	0.00	0.48	0.00	1.00
Ln(Total Assets)	330,790	6.74	6.64	2.96	4.75	8.61
PPE	328,021	0.28	0.23	0.23	0.01	0.43
Tobin’s Q	225,349	1.67	0.64	19.97	0.31	1.30
Leverage	330,580	0.95	0.52	63.21	0.32	0.69
Capital Expenditure	311,017	0.06	0.03	0.78	0.01	0.07
ROA	328,801	0.01	0.05	0.22	0.01	0.10
Sales Growth	302,442	0.26	0.07	1.05	-0.04	0.25
No. of Patent Applications	333,920	9.36	0.00	140.78	0.00	0.00
No. of Patents Granted	333,920	4.50	0.00	69.54	0.00	0.00
Citations	333,920	23.43	0.00	442.67	0.00	0.00

Panel B. Univariate Comparison of Samples

	Full Sample	Report R&D	Report Patent	R&D and Patent	Differences		
	(1)	(2)	(3)	(4)	(5) = ((1)-(2))/(1)	(6) = ((1)-(3))/(1)	(7) = ((1)-(4))/(1)
Ln(Total Assets)	6.74	7.25	7.47	7.40	-8%***	-11%***	-10%***
PPE	0.28	0.24	0.23	0.20	14%***	18%***	29%***
Tobin’s Q	1.67	1.55	1.74	1.86	7%**	-4%	-11%***
Leverage	0.95	0.53	0.57	0.48	44%***	40%***	49%***
Capital Expenditure	0.06	0.05	0.05	0.05	17%***	17%***	17%***
ROA	0.01	0.00	-0.01	-0.03	100%***	200%***	400%***
Sales Growth	0.26	0.23	0.25	0.31	12%***	4%**	-19%***
N (Firm-years)	330,790	122,546	118,264	53,456			

Table 2
Testing MCAR

The table presents the missing completely at random test for the predictability of unreported innovation. The test is based on Little (1988) test for MCAR. Columns (1)-(4) *World* present the results for all countries in the sample. Columns (5)-(7) present the results for the *US* only. D.o.F. is the number of degrees of freedom, Prob> χ^2 is the probability of the null hypothesis that the data is MCAR. Variable definitions are presented in Table A1.

	World				US		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
R&D(\$value)	X	X	X	X	X	X	X
Num. Patent Appl.	X	X	X	X	X	X	X
Ln(Total Assets)		X	X	X		X	X
PPE			X	X		X	X
Leverage			X	X		X	X
CapEx			X	X			X
ROA				X			X
Sales Growth				X			X
χ^2 dist.	297	7,431	25,062	42,971	6,359	9,857	22,889
D.o.F.	2	9	87	326	5	23	180
Prob> χ^2	0.00	0.00	0.00	0.00	0.00	0.00	0.00

Table 3
Predictability of Unreported R&D

The table presents the OLS regression results for predictability of unreported R&D. Columns (1)-(4) *World* present the results for all countries in the sample. Columns (5)-(7) present the results for the *US* only. Standard errors are double clustered at firm and year level. T-statistics are presented in brackets. Variable definitions are presented in Table A1. *, **, and *** represent significance at the 10%, 5%, and 1% levels, respectively. Adj. R² is the adjusted R².

	World				US		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
Ln(Total Assets)	-0.028*** (-12.17)	-0.024*** (-10.29)	-0.018*** (-11.24)	-0.009** (-2.67)	0.033*** (8.91)	-0.001 (-0.25)	-0.009** (-2.71)
PPE	0.228*** (13.83)	0.206*** (14.01)	0.179*** (17.31)	0.016* (1.95)	0.237*** (6.44)	0.309*** (9.20)	0.049** (2.24)
Leverage	-0.000 (-1.04)	-0.000 (-1.06)	0.000** (2.31)	-0.000 (-0.11)	0.002** (2.18)	0.001** (2.32)	0.001*** (6.73)
CapEx	0.004 (1.35)	0.003 (1.22)	-0.000 (-0.62)	0.000 (0.22)	0.019 (0.31)	-0.159*** (-3.02)	-0.015 (-1.18)
ROA	0.182*** (15.34)	0.172*** (14.66)	0.128*** (10.79)	0.022*** (3.79)	0.178*** (8.13)	0.166*** (10.19)	0.039*** (3.92)
Sales Growth	0.011*** (3.29)	0.006* (1.79)	-0.002 (-1.11)	0.002** (2.32)	-0.006* (-1.77)	-0.007** (-3.17)	-0.001 (-0.88)
Stock Liquidity	-0.006*** (-7.47)	-0.006*** (-8.54)	-0.005*** (-12.00)	-0.000 (-1.38)	-0.008*** (-12.37)	-0.003*** (-5.43)	-0.000 (-1.53)
Patent Intensity	-556.167*** (-15.05)	-6.925 (-0.50)	13.555 (1.08)	37.054** (2.22)	-573.734*** (-12.74)	-21.605*** (-4.42)	-19.704** (-2.36)
Country FE			Yes				
Industry FE		Yes	Yes			Yes	
Year FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Firm FE				Yes			Yes
N	283,987	283,987	283,987	281,243	64,386	64,386	63,086
Adj. R ²	0.14	0.23	0.38	0.81	0.23	0.53	0.93

Table 4
Predicting Non-patent Seeking Firms

The table presents OLS regressions of unreported USPTO patents and explanatory variables. The dependent variable is an indicator variable equal to 1 when a firm does not have USPTO patents, and zero otherwise. Columns (1)-(4) *World* present the regression results for all countries in the sample. Columns (5)-(7) present the regression for *US listed firms* only. Standard errors are double clustered at firm and year level. T-statistics are presented in brackets. Variable definitions are presented in Table A1. *, **, and *** represent significance at the 10%, 5%, and 1% levels, respectively. Adj. R² is the adjusted R².

	World				US		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
Ln(Total Assets)	-0.018*** (-15.49)	-0.014*** (-13.42)	-0.031*** (-23.36)	-0.010*** (-8.03)	-0.018*** (-4.97)	-0.036*** (-12.75)	-0.021*** (-5.04)
PPE	0.109*** (12.66)	0.149*** (15.24)	0.117*** (13.60)	-0.007 (-1.72)	0.123*** (5.68)	0.225*** (8.12)	-0.034 (-1.71)
Leverage	0.000 (0.62)	0.000 (0.54)	0.000 (1.06)	-0.000 (-1.08)	0.001 (1.79)	0.000 (0.71)	-0.000 (-1.23)
CapEX	0.000 (0.17)	-0.000 (-0.15)	-0.001 (-1.36)	-0.000 (-0.85)	0.071 (1.55)	-0.122** (-2.40)	0.025 (1.27)
ROA	0.141*** (12.08)	0.121*** (13.15)	0.145*** (16.50)	0.015*** (3.17)	0.286*** (8.77)	0.221*** (8.25)	0.019 (1.88)
Sales Growth	0.003** (2.18)	0.001 (0.73)	-0.005*** (-4.62)	0.002*** (3.78)	-0.004 (-0.94)	-0.004** (-1.96)	0.005*** (3.07)
Stock Liquidity	-0.007*** (-19.71)	-0.007*** (-21.10)	-0.004*** (-14.43)	-0.000** (-2.55)	-0.007*** (-14.81)	-0.005*** (-11.10)	-0.001** (-2.04)
Patent Intensity	-355.432*** (-16.00)	-40.196** (-3.26)	-30.734** (-3.38)	-10.067 (-1.64)	-518.481*** (-14.05)	-44.799** (-3.11)	-11.023 (-0.74)
R&D Stock	-0.000** (-2.11)	-0.000** (-2.13)	-0.000** (-2.25)	-0.000 (-1.12)	-0.000*** (-6.37)	-0.000*** (-4.71)	0.000 (1.78)
Country FE			Yes				
Industry FE		Yes	Yes			Yes	
Year FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Firm FE				Yes			Yes
N	281,763	281,763	281,763	278,999	64,383	64,383	63,086
R ²	0.11	0.16	0.26	0.76	0.18	0.32	0.77

Table 5
Recovered R&D and Imputed Unreported R&D

The table presents the Recovered R&D (in the years t-1 and t-2 from switch year) statistics and its comparison with different imputation methods. Panel A presents the comparison of Recovered R&D, Zero R&D, and positive R&D firm characteristics. Panel B presents the comparison of Recovered R&D with different imputation methods. *Recovered R&D* is the recovered R&D expenditure as reported in 10-K filings, its t-stat presents its difference from 0. *Imputed R&D (Industry Avg.)* is the average industry expenditure (two-digit) for the observations that are recovered, *Imputed R&D MI (Full Sample)* is the multiply imputed R&D using only the Lasso variables: ln(total assets), stock liquidity, and industry patent intensity, by industry (two-digit) for the complete sample, *Imputed R&D MI (Sub Sample)* is the multiply imputed R&D using the same MI model on the restated R&D sub sample and the industry and size matched peers. “Diff.” is the difference between Recovered R&D and imputed R&D. t-stat. represent the t-statistic for the difference between Recovered R&D an imputed R&D. Panel C shows the rank correlation among text-based innovation, patent applications (USPTO only and PATSTAT), R&D expenditure, and multiple imputation. *, **, and *** represent significance at the 10%, 5%, and 1% levels, respectively.

Panel A. Recovered R&D and R&D Firms

	Recovered R&D	Zero R&D	Diff.	Positive R&D	Diff.
R&D (\$ value)	6.69	0.00	6.69***	112.76	-106.07***
R&D Expenditure	0.87	0.00	0.87	0.34	0.53
Ln(Total Assets)	3.60	4.84	-1.24***	4.69	-1.09***
ROA	-3.11	-2.81	-0.30	-0.70	-2.41
PPE	0.20	0.30	-0.10***	0.18	0.02***
Sales growth	15.31	1.54	13.77	1.49	13.82
Capex	0.06	0.06	-0.01	0.05	0.01
Leverage	2.61	6.74	-4.13***	1.81	0.80

Panel B. Comparison with Imputation

Variable	Mean	St. Dev.	R&D	Diff.	t-stat.
Recovered R&D	6.91	24.24			8.84
Imputed R&D (Industry Avg.)	77.86	92.13	6.91	-70.95	-23.19
Imputed R&D MI (Full Sample)	6.36	242.75	6.91	0.55	0.07
Imputed R&D MI (Sub Sample)	8.66	245.55	6.91	-1.74	-0.22

Panel C. Rank Correlation MI and Text-based Innovation

	Text-based Innovation	Text-based Negative Innovation
Patent USPTO	0.22***	0.17***
Patent PATSTAT	0.21***	0.15***
R&D	0.26***	0.29***
Imputed R&D MI Full Sample	0.30***	0.27***

Table 6
Simulation Based on the Empirical Distribution from Compustat Data

This table provides the evaluation statistics, bias (relative bias over true parameter) and root mean squared error (RMSE) for the simulation based on the empirical distribution from Compustat (US) data, as described in section 5.1. The empirical distribution is from the panel of 783 firms with non-missing information for all variables except R&D. The methods evaluated are listwise deletion (LD), imputation with zero (ImpZero), imputation with industry mean, two-digit SIC code (ImpMean), inverse probability weighting (IPW), Heckman procedure (Heckman), and multiple imputation (MI). MI uses total assets, stock liquidity, industry patent intensity identified using Lasso analysis in the regression and is estimated using MCMC with 200 iterations for convergence. The regressions for imputation with zero and industry mean include a dummy variable for the imputed observations. Absolute average represents the average of the absolute bias across all variables. We present results for three missingness mechanisms: missing completely at random (MCAR) in Panel A, missing at random (MAR) in Panel B, and missing not at random (MNAR) in Panel C. Variable definitions are presented in Table A1. We generate missingness R&D for 50 and 70% of the sample. We conduct 500 simulations.

		Missing 70%						Missing 50%					
		LD	Imp Zero	Imp Mean	IPW	Heckman	MI	LD	Imp Zero	Imp Mean	IPW	Heckman	MI
<i>Panel A. MCAR</i>													
Bias	R&D	0.84	-0.58	-0.48	0.76	0.80	-0.04	0.67	-0.52	-0.47	0.73	0.68	-0.19
	Ln(Total Assets)	1.42	-0.23	-0.22	18.00	18.11	0.01	0.86	-0.34	-0.33	18.00	18.11	-0.18
	Tobin's Q	1.18	0.15	0.13	0.40	0.10	0.07	0.73	0.10	0.09	0.35	0.08	0.03
	Leverage	0.24	-0.02	-0.02	0.82	0.75	0.03	0.20	-0.02	-0.01	0.81	0.69	0.02
	ROA	0.44	-0.01	-0.01	1.23	1.19	0.15	0.39	-0.01	-0.01	1.20	1.04	0.09
	<i>Avg. Abs. Bias</i>		<i>0.83</i>	<i>0.20</i>	<i>0.17</i>	<i>4.24</i>	<i>4.19</i>	<i>0.06</i>	<i>0.57</i>	<i>0.19</i>	<i>0.18</i>	<i>4.22</i>	<i>4.12</i>
RMSE	R&D	0.01	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	Ln(Total Assets)	0.02	0.01	0.01	0.19	0.19	0.01	0.02	0.01	0.01	0.19	0.19	0.01
	Tobin's Q	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
	Leverage	0.08	0.03	0.03	0.15	0.13	0.03	0.06	0.03	0.03	0.14	0.12	0.03
	ROA	0.13	0.04	0.04	0.22	0.20	0.05	0.10	0.04	0.04	0.19	0.17	0.05

Panel B. MAR

Bias	R&D	1.09	-0.66	-0.54	4.60	-19.52	-0.17	0.84	-0.38	-0.26	3.63	-23.51	0.00
	Ln(Total Assets)	1.43	-0.32	-0.28	62.02	12.88	-0.03	1.00	-0.18	-0.16	62.21	26.46	0.05
	Tobin's Q	1.25	0.24	0.24	2.68	-4.33	0.12	0.93	0.06	0.04	1.40	-17.88	0.02
	Leverage	0.42	-0.04	-0.04	3.57	-0.67	-0.06	0.22	-0.03	-0.03	2.67	-0.71	0.00
	ROA	0.86	-0.04	-0.04	5.10	1.89	0.09	0.43	-0.04	-0.04	4.49	5.53	0.06
	<i>Avg. Abs. Bias</i>	<i>1.01</i>	<i>0.26</i>	<i>0.23</i>	<i>15.59</i>	<i>7.86</i>	<i>0.09</i>	<i>0.69</i>	<i>0.14</i>	<i>0.11</i>	<i>14.88</i>	<i>14.82</i>	<i>0.03</i>
RMSE	R&D	0.01	0.00	0.00	0.02	0.33	0.00	0.00	0.00	0.00	0.01	0.30	0.00
	Ln(Total Assets)	0.02	0.01	0.01	0.64	0.94	0.01	0.02	0.01	0.01	0.64	0.90	0.01
	Tobin's Q	0.02	0.00	0.00	0.02	1.19	0.01	0.01	0.01	0.01	0.02	1.06	0.01
	Leverage	0.12	0.04	0.04	0.61	2.97	0.04	0.06	0.04	0.04	0.46	2.90	0.03
	ROA	0.22	0.05	0.05	0.78	4.25	0.06	0.11	0.05	0.05	0.68	4.32	0.05

Panel C. MNAR

Bias	R&D	0.89	-0.55	-0.49	4.81	-30.91	-0.21	0.64	-0.55	-0.52	3.61	-16.35	-0.15
	Ln(Total Assets)	1.31	-0.18	-0.18	61.87	-2.98	0.09	0.96	-0.26	-0.25	62.06	1.27	-0.04
	Tobin's Q	1.57	0.40	0.40	3.17	-33.98	0.29	0.89	0.31	0.30	2.07	-7.09	0.26
	Leverage	0.36	-0.08	-0.08	3.51	-5.10	-0.10	0.20	-0.06	-0.06	2.63	-1.56	-0.03
	ROA	0.83	-0.07	-0.07	5.06	-0.46	0.05	0.45	-0.06	-0.06	4.21	-0.17	0.05
	<i>Avg. Abs. Bias</i>	<i>0.99</i>	<i>0.25</i>	<i>0.24</i>	<i>15.68</i>	<i>14.69</i>	<i>0.15</i>	<i>0.63</i>	<i>0.25</i>	<i>0.24</i>	<i>14.92</i>	<i>5.29</i>	<i>0.10</i>
RMSE	R&D	0.01	0.00	0.00	0.02	0.38	0.00	0.00	0.00	0.00	0.01	0.35	0.00
	Ln(Total Assets)	0.02	0.01	0.01	0.64	0.96	0.01	0.02	0.01	0.01	0.64	0.92	0.01
	Tobin's Q	0.02	0.01	0.01	0.03	1.29	0.01	0.01	0.01	0.01	0.02	1.17	0.01
	Leverage	0.11	0.04	0.04	0.61	2.87	0.04	0.06	0.04	0.04	0.45	2.62	0.04
	ROA	0.21	0.04	0.04	0.78	4.34	0.05	0.11	0.05	0.05	0.64	4.84	0.05

Table 7
Imputation Effect on Empirical Inference

This table replicates the results in Fama and French (2002) using different imputation methods and two-way fixed effects. We present the results of a contemporaneous regression with two-way fixed effects: $\frac{L_t}{A_t} = \beta_0 + \beta_1 \frac{V_t}{A_t} + \beta_2 \frac{ET_t}{A_t} + \beta_3 \frac{Dp_t}{A_t} + \beta_4 RDD_t + \beta_5 \frac{RD_t}{A_t} + \beta_6 \ln(A_t) + e_t$. “ImpZero” presents the result for the sample with imputation with zero and an indicator variable, “LD” presents the results for listwise deletion, “MI” presents the results for multiple imputation implemented using all the variables in the regression in the imputation, “MI Lasso” presents the results for multiple imputation implemented using all the variables in the regression and the Lasso variables stock liquidity and industry patent intensity in the imputation, “Pseudo RD” presents the result using pseudo R&D as an explanatory variable, and “An. Innov.” presents the results for the analyst coverage based innovation variable (Bellstam et al., 2020). The dependent variable is book leverage $\frac{L_t}{A_t}$ at time T . $\frac{V_t}{A_t}$ is the market to book ratio, $\frac{ET_t}{A_t}$ is earnings before interest and taxes as a proportion of total assets, $\frac{Dp_t}{A_t}$ is depreciation as a proportion of total assets, $\frac{RD_t}{A_t}$ is the R&D expenses as a proportion of total assets, RDD_t is an indicator variable equal to 1 if R&D expenditure is missing and has been imputed with zero, and zero otherwise, $Pseudo\ R\&D_t$ is an indicator variable equal to 1 if a firm applies for a patent in PATSTAT and has no reported R&D, and zero otherwise, $An.\ Innov._t$ is the firm analyst-based innovation measure from (Bellstam et al., 2020), and $\ln(A_t)$ is the natural logarithm of total assets. Non-dividend payers include firms that do not pay dividend in year $T-1$. Panel A presents the results for the dividend paying firms and Panel B for the non-dividend paying firms. The sample period is 1965-1999. Standard errors are double clustered.

Panel A. Dividend Payer Firms

Variable	Imp Zero (1)	LD (2)	MI (3)	MI Lasso (4)	Pseudo R&D (5)	Text-based Innov. (6)
Intercept	0.305*** (22.62)	0.344*** (19.83)	0.366*** (56.52)	0.368*** (55.24)	0.300*** (22.13)	0.246*** (3.94)
$\frac{V_t}{A_t}$	-0.001 (-0.15)	-0.001 (-0.47)	0.001 (0.40)	0.001 (0.60)	0.000 (-0.10)	-0.006 (-1.29)
$\frac{ET_t}{A_t}$	-0.158** (-1.99)	-0.215** (-2.66)	-0.184** (-2.07)	-0.192 (-2.17)	-0.157 (-1.99)	-0.628 (-3.91)
$\frac{Dp_t}{A_t}$	-1.076*** (-6.12)	-0.059 (-0.30)	-1.057** (-10.67)	-1.048*** (-10.56)	-1.049*** (-6.05)	-0.797** (-2.77)
RDD_t	0.070*** (11.96)				0.075*** (12.35)	
$\frac{RD_t}{A_t}$	-0.290*** (-2.71)	-0.435*** (-4.54)	0.081*** (3.91)	0.033 (1.48)	-0.290*** (-2.72)	
<i>Pseudo R&D</i>					-0.098*** (-12.43)	
<i>An. Innov.</i>						-0.009* (-1.78)
$\ln(A_t)$	0.041*** (29.95)	0.029*** (13.61)	0.038*** (96.36)	0.038*** (95.10)	0.042*** (30.14)	0.048*** (6.89)

Panel B. Non-dividend Payer Firms

Variable	Zero (1)	Delete (2)	MI (3)	MI Lasso (4)	Pseudo R&D (5)	Text-based Innov. (6)
Intercept	0.325*** (4.70)	0.394*** (20.07)	0.376*** (6.74)	0.381*** (7.05)	0.323*** (4.66)	0.242 (1.11)
$\frac{V_t}{A_t}$	0.027 (1.32)	-0.004*** (-3.24)	0.028** (2.24)	0.029*** (2.30)	0.027 (1.32)	-0.008 (-1.40)
$\frac{ET_t}{A_t}$	-0.517*** (-3.15)	-0.301*** (-4.77)	-0.139 (-0.54)	-0.136 (-0.52)	-0.517*** (-3.14)	-0.404 (-1.56)
$\frac{Dp_t}{A_t}$	0.691 (1.29)	1.984*** (7.96)	0.636 (1.66)	0.651 (1.70)	0.692 (1.29)	1.725 (1.84)
<i>RDD</i>	0.079*** (4.61)				0.082*** (4.73)	
$\frac{RD_t}{A_t}$	-0.702*** (-2.83)	-0.335*** (-3.33)	0.955*** (3.45)	0.962*** (3.43)	-0.701*** (-2.83)	
<i>Pseudo R&D</i>					-0.134*** (-6.03)	
<i>An. Innov.</i>						-0.095*** (-5.50)
$\ln(A_t)$	0.032*** (4.54)	0.013** (2.60)	0.022*** (4.98)	0.024*** (5.35)	0.033*** (4.62)	0.042 (1.60)

Table 8
New Products, Patents, and Imputation Methods

The table presents an analysis of new product announcements, patents, and imputation methods for firms with different combinations of products and patents. Panel A presents the innovation (R&D and patents) and new product characteristics. New product announcement data is from (Mukherjee et al., 2016) and patents are based on USPTO data. New products include the average returns for all new product announcements and major new products includes the number of new products in the 75th percentile of returns. Panel B presents the comparison of single and multiple imputation methods for patents with different product announcements. Single imputation includes: imputation with zero (*Impute 0*), and imputation with the two-digit industry average (*Impute Industry Average*). MI M1 presents the multiply imputed USPTO patents using the Lasso variables: ln(total assets), stock liquidity, R&D stock, and industry patent intensity from PATSTAT patents and MI M2 is the same as M1 with the addition of ROA, PPE, capital expenditure, sales growth, and leverage. Column (1) presents the information for the subsample with no USPTO patents and with product announcements; Column (2) presents the information for the subsample with no USPTO patents and no product announcements; Column (3) presents the information for firms with USPTO patents and product announcements; and Column (4) presents the information for firms with USPTO patents and no new product announcements. *Diff.* presents the difference between Columns (1) and (4), and *t-stat.* present the t-statistic for the difference.

Panel A. Patents and New Products

	No Patents + New Products	No Patents + No New Products	Patents + New Products	Patents + No New Products	Diff.	t-stat.
	(1)	(2)	(3)	(4)	(5) = (1)-(4)	
Obs. (% of Sample)	5.71%	75.58%	3.10%	13.83%		
R&D (\$ mio)	85.72	32.05	285.20	148.39		
R&D (% of report)	56.79%	39.12%	91.38%	85.31%		
Number of Patents	0.00	0.00	47.87	32.07		
New products	0.07	0.00	0.13	0.00	0.07	40.18
Major New Products	0.79	0.00	1.60	0.00	0.79	42.02

Panel B: Compare Imputation Method: Imputed Patents

	No Patents + New Products	No Patents + No New Products	Patents + New Products	Patents + No New Products	Diff.	t-stat
	(1)	(2)	(3)	(4)	(5) = (1)-(4)	
Single Imputation						
Impute 0	0.00	0.00	47.87	32.07	-32.07	-21.87
Impute Industry Average	24.54	28.74	47.87	32.07	-7.53	-4.00
Multiple Imputation						
Imputed Patents MI M1	29.38	11.98	47.87	32.07	-2.68	-1.26
Imputed Patents MI M2	30.50	10.28	47.87	32.07	-1.57	-0.73

Figure 1
Ranking Innovative Firms

The figure shows the rank of innovative firms for S&P 500 using R&D and Text-Based Analysis. The figure shows innovation ranks based on the R&D expenditure imputed with zero for unreported R&D (*Imp Zero*), Text-based innovation measure (*Text-based*), and Multiple Imputation (*MI*). Firms with reported R&D are shown in light orange, and those with unreported R&D as red.

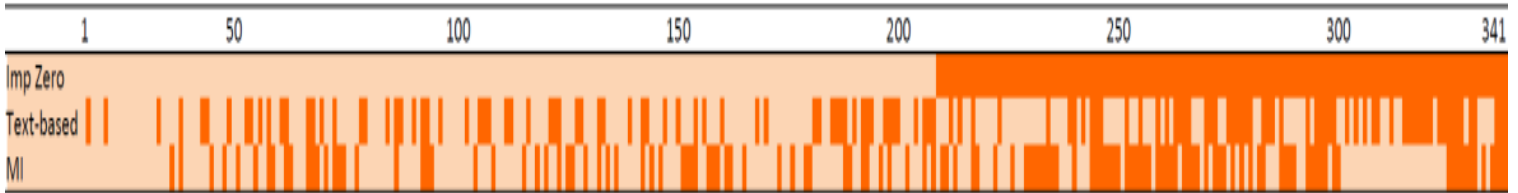


Figure 2
Bias From Deleting Firms Without Reported Innovation

This figure presents the bias of the R&D coefficient for listwise deletion (LD) and multiple imputation (MI) across different missingness levels. The simulation is based on the empirical distribution of the panel of 783 firms with non-missing information for all variables except R&D. MI uses all the variables in the regression in Section 5.1 and is estimated using MCMC with 200 iterations for convergence. We present results for data missing at random (MAR). We conduct 500 simulations.

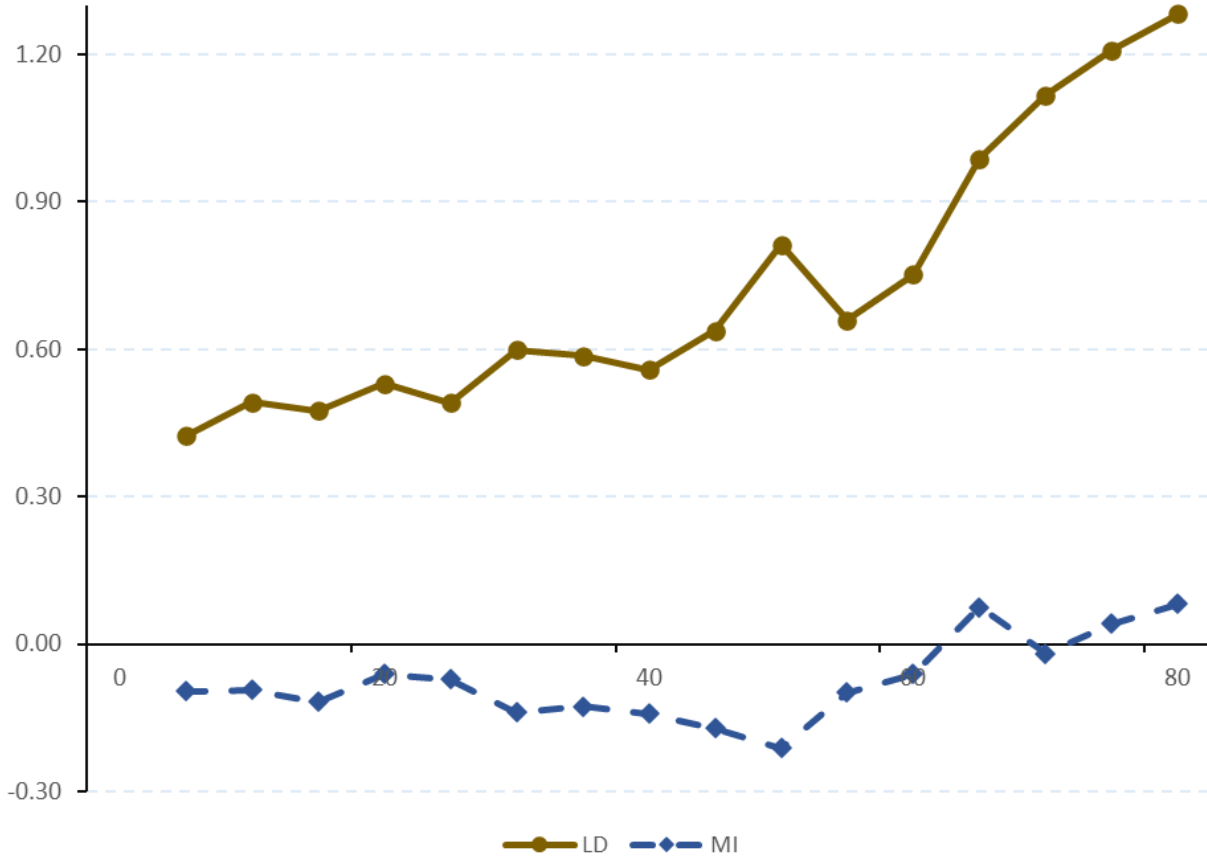
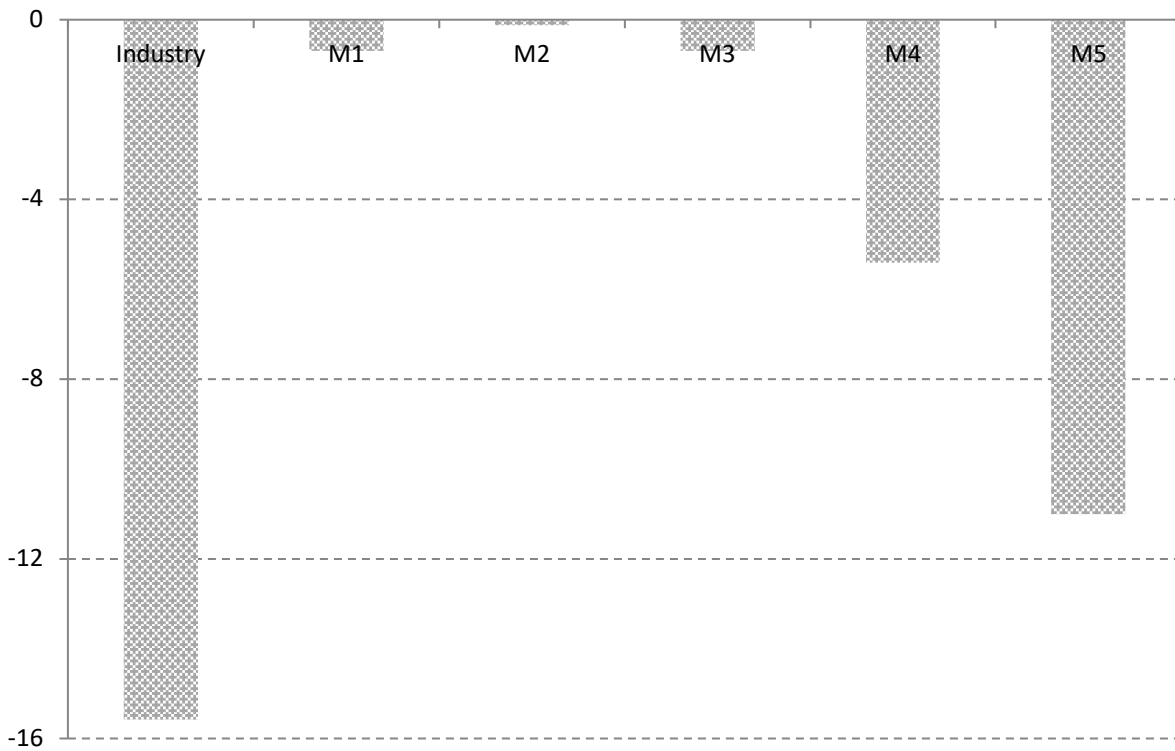


Figure 3
Imputation of Missing Patents

This figure presents the t-statistics of the comparison of years with only non-USPTO patents for US firms with different imputation methods. Industry is the difference between non-USPTO patents and industry patents defined as the average industry expenditure for the observations the year with non-USPTO patents only, M1 is the multiply imputed non-USPTO patents using $\ln(\text{total assets})$, ROA, PPE, capital expenditure, sales growth, and leverage by industry (two-digit), M2 is the multiply imputed non-USPTO patents using the same model as M1 without sales growth and leverage as conditioning information, M3 is the multiply imputed non-USPTO patents using the same model as M1 with the addition of R&D expenditure as conditioning information, M4 is the multiply imputed non-USPTO patents using the Lasso variables $\ln(\text{total assets})$, stock liquidity, industry patent intensity, and stock R&D, M5 is the multiply imputed non-USPTO patents combining models M1 and M5.



Appendix

Table A1
Variable Definitions

This table shows the variable definitions.

Variable Names	Variable Definitions	Code
R&D Expenditure	R&D expenditure divided by total assets	XRD/AT
Report R&D	Indicator variable: 1 if a firm reported zero or positive R&D expenditure; 0 otherwise	
PPE	Net property, plant, and equipment divided by total assets	PPENT/AT
Tobin's Q	Tobin's Q, measured as market value of equity divided by total assets	MKTVAL/AT
Leverage	Total liabilities divided by total assets	LT/AT
Ln(Total Assets)	Natural log of total assets	Ln(AT)
Capital Expenditure	Capital expenditure divided by total assets	CAPX/AT
ROA	EBIT divided by total assets	EBIT/AT
Sales Growth	Annual sales growth	$(\text{Sale}_t - \text{Sale}_{t-1}) / \text{Sale}_{t-1}$
HH Index	Herfindahl industry concentration index	
No. of Patent Applications	Total number of patent applications	
No. of Patents Granted	Total number of patents granted	
Citations	Total number of citations per patent	
Liquidity	Yearly sum of daily trading volume in USD	PRC*VOL (for US-Stocks), PRCCD*CS HTRD*Exchange Rate (for non-US stocks)
Patent Intensity	Number of PATSTAT patents per total assets for industry, using two-digit SIC across countries, unless specified otherwise	

Table A2
Simulation Based on the Empirical Distribution from Compustat Data

This table provides the evaluation statistics, bias (relative bias over true parameter) and root mean squared error (RMSE) for the simulation based on the empirical distribution from Compustat (US) data, as described in section 5.1. Bias presents the average of the absolute bias across all five variables and RMSE presents the average RMSE across the five variables. The empirical distribution is from the panel of 783 firms with non-missing information for all variables except R&D. The methods evaluated are listwise deletion (LD), imputation with zero (ImpZero), imputation with industry mean, two-digit SIC code (ImpMean), inverse probability weighting (IPW), Heckman procedure (Heckman), and multiple imputation (MI). The regressions for imputation with zero and industry mean include a dummy variable for the imputed observations. MI is spec uses all the variables in the regression and is estimated using MCMC with 200 iterations for convergence. We present results for three missingness mechanisms: missing completely at random (MCAR), missing at random (MAR), and missing not at random (MNAR). Panel A presents the results for the missingness regression which includes the lasso variables. Panel B presents the results with the MI specification in Panel A as well as includes the Lasso variables in the Sales growth regression. Panel C presents the Double Lasso results. Variable definitions are presented in Table A1. We generate missingness R&D for 50 and 70% of the sample. We conduct 500 simulations.

		Missing 70%						Missing 50%					
		LD	Imp. Zero	Imp. Mean	IPW	Heckman	MI	LD	Imp. Zero	Imp. Mean	IPW	Heckman	MI

Panel A. Missingness Regression with Q, A, V, and PI

MCAR	Bias	0.80	0.24	0.22	3.69	3.67	0.11	0.63	0.16	0.13	3.42	3.36	0.05
	RMSE	0.05	0.02	0.02	0.09	0.09	0.02	0.03	0.02	0.02	0.07	0.07	0.02
MAR	Bias	1.02	0.13	0.12	20.22	118.27	0.10	0.63	0.17	0.15	17.43	100.29	0.07
	RMSE	0.07	0.02	0.02	0.53	2.94	0.02	0.04	0.02	0.02	0.45	2.87	0.02
MNAR	Bias	0.98	0.13	0.12	11.96	67.10	0.11	0.58	0.18	0.15	10.18	59.42	0.08
	RMSE	0.07	0.02	0.02	0.26	2.27	0.02	0.03	0.02	0.02	0.25	2.05	0.02

Panel B. Missingness Regression with Q, A, V, and PI and Sales growth regression with V and PI

MCAR	Bias	0.86	0.26	0.24	3.74	3.75	0.17	0.60	0.23	0.15	3.52	3.48	0.08
	RMSE	0.05	0.02	0.02	0.09	0.09	0.02	0.04	0.02	0.02	0.07	0.08	0.02
MAR	Bias	1.10	0.25	0.23	18.96	100.69	0.09	0.61	0.18	0.17	16.26	98.45	0.06
	RMSE	0.08	0.02	0.02	0.48	2.73	0.02	0.04	0.02	0.02	0.39	2.50	0.02
MNAR	Bias	0.99	0.13	0.11	11.49	61.56	0.14	0.60	0.17	0.14	9.93	52.77	0.05
	RMSE	0.07	0.02	0.02	0.29	1.86	0.02	0.03	0.02	0.02	0.24	1.76	0.02

Panel C. Double Lasso

MCAR	Bias	0.77	0.24	0.22	3.77	3.74	0.08	0.60	0.21	0.19	3.59	3.47	0.09
	RMSE	0.05	0.02	0.02	0.09	0.09	0.02	0.04	0.02	0.02	0.07	0.07	0.02
MAR	Bias	0.66	0.18	0.15	15.51	5.76	0.05	0.48	0.19	0.17	16.46	4.50	0.09
	RMSE	0.03	0.02	0.02	0.45	0.40	0.02	0.02	0.02	0.02	0.43	0.35	0.02
MNAR	Bias	0.63	0.24	0.20	10.00	3.58	0.07	0.49	0.20	0.18	10.08	3.34	0.06
	RMSE	0.03	0.02	0.02	0.29	0.24	0.02	0.02	0.02	0.02	0.25	0.18	0.02

Table A3
Simulation Based on Simulated Data

This table provides the evaluation statistics, bias (relative bias over true parameter) and root mean squared error (RMSE) for the simulation based on simulated data, as described in section 5.2. The methods evaluated are listwise deletion (LD), imputation with zero (ImpZero), imputation with industry mean, two-digit SIC code (ImpMean), inverse probability weighting (IPW), Heckman procedure (Heckman), and multiple imputation (MI). MI is estimated using MCMC with 200 iterations for convergence. The regressions for imputation with zero and industry mean include a dummy variable for the imputed observations. We present results for three missingness mechanisms: missing completely at random (MCAR) in Panel A, missing at random (MAR) in Panel B, and missing not at random (MNAR) in Panel C. We generate missingness in x_1 for 50 and 70% of the sample. We conduct 500 simulations.

		Missing 70%						Missing 50%					
		LD	Imp Zero	Imp Mean	IPW	Heck man	MI	LD	Imp Zero	Imp Mean	IPW	Heck man	MI
<i>Panel A. MCAR</i>													
Bias	θ_1	0.00	-0.19	-0.19	0.00	0.00	-0.01	0.00	-0.13	-0.13	0.00	0.00	-0.01
	θ_2	0.01	0.28	0.28	0.01	0.01	0.01	0.00	0.19	0.19	0.00	0.00	0.00
RMSE	θ_1	0.11	0.21	0.21	0.08	0.11	0.09	0.06	0.07	0.07	0.06	0.06	0.05
	θ_2	0.11	0.29	0.29	0.08	0.11	0.09	0.06	0.10	0.10	0.07	0.06	0.05
<i>Panel B. MAR</i>													
Bias	θ_1	-0.15	-0.23	-0.23	-0.16	-0.08	-0.08	-0.11	-0.16	-0.16	-0.11	-0.09	-0.05
	θ_2	-0.12	0.12	0.12	-0.12	-0.08	-0.05	-0.08	0.04	0.04	-0.07	-0.06	-0.04
RMSE	θ_1	0.17	0.24	0.24	0.18	0.17	0.10	0.13	0.17	0.17	0.13	0.12	0.08
	θ_2	0.15	0.13	0.13	0.15	0.16	0.09	0.07	0.06	0.06	0.07	0.07	0.05
<i>Panel C. MNAR</i>													
Bias	θ_1	-0.17	-0.28	-0.28	-0.19	-0.13	-0.10	-0.13	-0.19	-0.19	-0.13	-0.11	-0.05
	θ_2	-0.16	0.14	0.14	-0.15	-0.13	-0.08	-0.11	0.04	0.04	-0.11	-0.10	-0.07
RMSE	θ_1	0.19	0.29	0.29	0.20	0.17	0.12	0.14	0.20	0.20	0.14	0.13	0.07
	θ_2	0.17	0.12	0.12	0.17	0.16	0.10	0.13	0.06	0.06	0.13	0.12	0.08

Internet Appendix

Internet Appendix I. *Handling Missing Data*

This appendix provides a summary of the missing data problem and discusses several popular econometric approaches to handling missing data that are considered in this paper. With partially observed data, we can rarely be sure of the mechanism leading to such missing data. Therefore, we highlight some approaches to analyzing missing data under different mechanisms, which helps to establish inference robustness in the face of uncertainty about the missingness mechanism. In particular, we consider listwise deletion, deterministic imputation, inverse probability weighting, Heckman correction, and multiple imputation. For exposition simplicity (as in the main body of the paper), we consider the case where only one explanatory variable contains missing observations. Let y_i be the dependent variable and z_i be the explanatory variables with missingness. We have the linear relation:

$$y_i = \alpha + \theta z_i + \varepsilon_i, \quad i = 1, \dots, N. \quad (\text{IA1})$$

Let s_i be a selection indicator where $s_i = 1$ when z_i is not missing and firm i is included in the regression. Otherwise, when $s_i = 0$ firm i is deleted from the data. The validity of solutions to this problem depends on the missing mechanism, thus we first present the three missing mechanisms.

1. Missing completely at random (MCAR):

$$P(s = 0|y, z, x) = P(s = 0).$$

This means that the missing probability does not depend on any random variables.

2. Missing at random (MAR): The probability of missing can be formulated by:

$$P(s = 0|y, z, x) = P(s = 0|x).$$

In other words, the probability of missingness only depends on the set of *observed* variables x , but not on the missing variable itself nor on unobservables.

3. Missing not at random (MNAR): the missing mechanism is neither MAR nor MCAR. For example, the missing mechanism depends on the value of z itself, or on unobserved variables, e.g., high-income individuals do not participate in surveys related to income.

Effects of Listwise Deletion

Listwise deletion only uses a subsample of observations, deleting those that contain missing values in the z -variable.¹ This leads to estimating the following regression using the subsample of the data:

¹ We consider the univariate setup for simplicity. There might be other covariates of interest that drive the outcome variable but including them in the regression does not change the problem of deletion.

$$y_i = s_i\alpha + \theta s_i z_i + s_i \varepsilon_i, \quad (\text{IA2})$$

where $s_i z_i$ is now the explanatory variable and $s_i \varepsilon_i$ is the error term. The OLS (ordinary least squares) estimator is unbiased if $E(s_i \varepsilon_i z_i) = 0$, which can be implied by $E(\varepsilon_i | z_i, s_i) = 0$. We can see that if MCAR holds and z_i is exogenous, then $E(\varepsilon_i | z_i, s_i) = E(\varepsilon_i | z_i) = 0$. Thus, deletion can lead to consistent estimates in the case of MCAR. However, if selection is driven by observed or even unobserved variables as in MAR and MNAR cases, $E(\varepsilon_i | z_i, s_i) \neq 0$ in general because ε_i can be correlated with s_i even if one controls for z_i , leading to biased estimates produced by deletion.

Deterministic Imputation

Another popular approach used in empirical studies is to impute the missing observations using various methods, and then treat the resulting data as given for further analysis. Frequently used deterministic imputation employs, e.g., zero, overall average, average from “similar” observations, or fitted values based on some pre-specified models. The validity of this method obviously depends on whether the specified imputation models are correct. If the imputation model perfectly coincides with the missing mechanism, then the resulting estimate using the imputed sample is consistent. On the contrary, misspecification of the imputation models can lead to potentially highly biased estimates because of the distortion of the variance-covariance matrices. For example, in our case, the missing R&D clearly does not equate to zero R&D in general (see results in Sections 3.2.2 and 3.3), and thus imputation using zeros leads to biased estimates when conditioning on R&D as an explanatory variable.

Inverse Probability Weighting

Inverse probability weighting assigns different weights to observed data points depending on their probability of being observed. Thus, the computation of IPW requires researchers to know the probability of being observed. Consider the case of MAR, where the probability of missing (or equivalently being observed) only depends on a set of observed variables x . Denote $p(x) \equiv P(s = 1 | x) = P(s = 1 | y, x, z)$, then we can solve the missing data problem by:

$$\min_{\alpha, \theta} \sum_{i=1}^N \left(\frac{s_i}{p(x_i)} \right) (y_i - \alpha - \theta z_i)^2.$$

In practice, $p(x)$ is often unknown except in some special cases, and thus we need to estimate it. To this end, we can regress the selection indicator s on x using flexible binary choice models, such as logit or probit, or even nonparametric models, and obtain the estimated selection probability (or alternatively called the propensity score) $\hat{p}(x)$.

Heckman Correction for Selection Bias

We know from (IA2) that the OLS estimator $\hat{\theta}$ is biased because $E(y_i | s_i = 1, z_i) = \alpha + \theta z_i + E(\varepsilon_i | z_i, s_i = 1)$, and $E(\varepsilon_i | z_i, s_i = 1) \neq 0$ in general. Heckman’s method assumes that the missing mechanism is determined by the following model:

$$s_i^* = \beta x_i + \eta_i, \quad i = 1, \dots, N, \quad (\text{IA3})$$

where s_i^* is the latent variable associated with s_i , i.e. $s_i = 1$ if $s_i^* > 0$ and $s_i = 0$ if $s_i^* \leq 0$. Further, assume that the error terms in (IA3) is normally distributed with variance σ_η^2 and correlated with ε_i in (IA1), and their covariance is ρ ; x and z are both exogeneous. The Heckman procedure approximates the “omitted variable” ($\varepsilon_i|z_i, s_i = 1$) by its consistent estimate and includes this proxy in the regression to correct for the bias. In particular, based on the joint distribution of η_i and ε_i , one could write $E(\varepsilon_i|z_i, s_i = 1) = \sigma_\eta \rho \lambda(x_i \beta) = \gamma \lambda(x_i \beta)$, where $\lambda(x_i \beta)$ is the inverse Mills ratio defined by:

$$\lambda(x_i \beta) = \frac{\phi(-x_i \beta)}{1 - \Phi(-x_i \beta)} = \frac{\phi(x_i \beta)}{\Phi(x_i \beta)}.$$

Then we can rewrite the conditional expectation of y_i given x_i and selection into the sample as:

$$E(y_i|x_i, s_i = 1) = \alpha + \theta z_i + \gamma \lambda(x_i \beta).$$

This leads to Heckman’s two-step procedure.

Step 1: Estimate a probit regression $P(s_i = 1|x_i) = \Phi(x_i \beta)$ using all N observations and obtain the estimate $\hat{\beta}$. Then compute the inverse Mills ratio $\lambda(x_i \hat{\beta})$.

Step 2: Estimate the regression $y_i = \alpha + \theta z_i + \gamma \lambda(x_i \hat{\beta})$ using OLS.

The estimates $\hat{\alpha}$, $\hat{\theta}$, and $\hat{\gamma}$ are consistent when x correctly includes all of the selection variables. The validity of Heckman’s procedure also heavily relies on the distributional assumptions of the two errors, η_i and ε_i . For example, the deviation from the normality assumption of η_i may negatively affect the performance of the Heckman’s procedure. Since γ captures the covariance between η_i and ε_i and a nonzero correlation implies selection bias, we can test whether selection is exogenous (or equivalently MCAR) by testing whether $\hat{\gamma} = 0$. For more extensions of Heckman’s procedure, see Wooldridge (2002, Chapter 17).

Multiple Imputation

Multiple imputation (MI) is essentially an iterative version of stochastic imputation, which aims at explicitly modeling the uncertainty/variability ignored by the deterministic imputation procedures. Instead of imputing in a single value, multiple imputation uses the (joint) distribution of the observed data to estimate the parameters of interest multiple times to capture the uncertainty/variability in this imputation procedure. A general multiple imputation procedure consists of three steps:

Step 1. Imputation: Impute the missing data with their estimates and create a complete sample. Repeat this process multiple times.

Step 2. Estimation: For each complete sample, estimate the parameters of interest.

Step 3. Pooling: Combine the parameter estimates obtained from each completed data set.

The imputation method should be chosen depending on the type of variables with missing observations and the pattern of missingness. For example, MI with multivariate normal regressions can be applied to impute one or more continuous variables of arbitrary missing-value patterns; MI with chained equations employs a separate conditional distribution for each imputed variable, and is often used to impute a variable with finite and discrete support (e.g., binary, multinomial, or count variable). We illustrate the MI with multivariate normal regressions (MI_MVN). As all MI methods, MI with multivariate normal regressions analyses the data in three steps: imputation, estimation, and pooling. We discuss the three steps in turn.

First, MI_MVN imputes the missing observations using data augmentation. In this case, we assume that the variable containing missing observations z is related with a set of (completely) observed variables x by:

$$z_i = \delta' x_i + v_i, \quad i = 1, \dots, N,$$

where $v_i \sim N(0, \sigma_v^2)$. Denote $w_i = (z_i, x_i)$. Data augmentation in this case is essentially an iterative Markov chain Monte Carlo (MCMC) procedure that iterates between two (sub-)steps, a replacement step and posterior step.

- Replacement step: We replace the missing values of z_i with draws from the conditional posterior distribution of z_i given observed variables and the values of model parameters in this iteration. Particularly, for each iteration t , we can replace the missing observations by:

$$z_i^{(t)} \sim P\left(z_i \mid x_i, \delta^{(t-1)}, \sigma_v^{(t-1)}\right), \quad \text{for } i \in \{i \mid s_i = 1\}.$$

- Posterior step: We draw the new values of model parameters from their conditional posterior distribution given the observed data and imputed data from the previous replacement step.

$$\sigma_v^{(t)} \sim P\left(\sigma_v \mid x_i, z_i^{(t)}\right), \quad \text{and} \quad \delta^{(t)} \sim P\left(\delta \mid x_i, z_i^{(t)}, \sigma_v^{(t)}\right),$$

where $z_i^{(t)}$ is the imputed value from iteration t if it is missing and the original value if non-missing.

The conditional posterior distributions above are jointly determined from the prior distribution for the model parameter $P(\delta, \sigma_v)$, e.g., uniform, Jeffreys, or ridge, and the assumed normal distribution of the data. These two steps (replacement and posterior) are iterated until a specified number of iterations or there is numerical convergence.

Second, we estimate the regression of interest (IA1) with the imputed (pseudo-complete) data set using various approaches, e.g., OLS, IV. Since the imputation is conducted for multiple times, sD times, we obtain multiple estimates for the same regression parameter θ .

Third, we combine/pool the estimates (coefficients and standard errors) across all imputed datasets and obtain a single statistic for each parameter. The final estimated slope coefficient $\hat{\theta}$ is simply an arithmetic mean of the corresponding estimate obtained from each of the imputed data. The variance of $\hat{\theta}$ is obtained by the total variance formula and is written by the average estimated variance of coefficient estimates across D imputed datasets plus the sample variance of coefficient estimates based on D imputations.

A major advantage of multiple imputation over deterministic imputation is that the final statistics appropriately reflect the uncertainty caused by imputation. If the joint normality is a reasonable assumption and the specification of \boldsymbol{x} is correct (i.e. MAR), MI_MVN produces consistent estimates. In practice, a safe strategy is to include all observables in \boldsymbol{x} including \boldsymbol{y} to better approximate the posterior distribution.

Internet Appendix II. *Simulation Extensions*

A. *Simulation of the with Generated Data*

We extend the benchmark simulation design by considering alternative specifications of the covariance matrix for the generated variables. We report results for the case of MAR with 70% missing observations in Table IA4 in Internet Appendix III, to conserve space. The ranking of the methods remains similar when considering MNAR and different levels of missingness.

First, we consider how the correlation between errors influences the performance of the methods. We increase the correlation between η and ε to 0.6. When the correlation is higher, the bias and RMSE in θ_1 under all six methods deteriorates (relative to benchmark case reported in Table A3 Panel B). Heckman and MI continue to have among the lowest biasness and RMSE in θ_1 , while deterministic imputation and MI exhibit the lowest biasness and RMSE in θ_2 . Next, we consider how the correlation between the selection variables and variables of interest influences the estimates. We increase the correlation between x_1 and z_1 to 0.6, and we find that an increase in correlation increases the biasness in θ_1 under LD and IPW but not Heckman; while that of ImpZero, ImpMean and MI improved. On the other hand, the biasness in θ_2 improved for all six methods.

Finally, we allow for correlation between the observed selection variables and the error in the main regression. We set the correlation between x_1 and ε to 0.4, but we generate η independently from ε to avoid direct endogeneity in the selection equation. In this case, even though the two errors are uncorrelated, the correlation between selection variables and the error term in the main regression also significantly biases θ_1 estimates under LD, ImpZero, ImpMean, and IPW. Heckman also performs poorly in this setting, because it is derived based on the joint distribution of ε and η , but does not consider the correlation between ε and other observables. The θ_2 biasness deteriorates significantly under LD, IPW and Heckman methods, while ImpZero and ImpMean exhibit smaller bias than under the baseline simulation. Overall, MI's performance improves and produces the lowest bias in θ_1 and θ_2 estimates among all methods. Its performance is even better than in the benchmark case because the correlation between errors is zero, leading to more accurate stochastic imputation based on the joint distribution with the selection variables.

B *Patent Empirical Data-Based Simulation*

Patents and R&D expenditures may have different determinants and missingness levels. To understand the properties of the different methods for handling missing data in the patent setting, we replicate the empirical distribution-based simulation, with the USPTO patent data distribution. The empirical distribution is derived from a panel of 783 firms with non-missing information for all variables except USPTO patents for the period 1992 to 2012. We follow the same simulation procedure as described in Section 5.1. We analyze the case of 70% missing data, as Table 1 shows that patents exhibit very large levels of missingness. In addition, we only show the results for MAR and MNAR, since the analyses in Table 2 and Table 4 show that patent data is not *missing completely at random*. Table IA5 in Internet Appendix III presents the results of the simulation based on the patent empirical distribution. Under MAR, IPW and Heckman generate the highest biasness in coefficient estimates relative to both imputation and deletion. Focusing on MNAR, deterministic imputation and multiple imputation both perform better than listwise deletion, IPW and Heckman approaches.

Internet Appendix III. Tables and Figures

Table IA1
Relaxing Firm Constraints

This table replicates the results of Table 1 Panel B on the difference between the univariate comparisons of the sample data with the full sample and innovation-variable based sample. “Full Sample” uses all available observations, “Report R&D” includes only observations that report R&D, “Report Patent” includes only observations that patent applications in any patent office, “R&D and Patent” includes only observations that have positive R&D and patent filings in the USPTO. Panel A only includes countries with more than 1,000 listed firms in the sample, Panel B only firms from Industrial and Commercial Machinery (SIC 35) and Chemical and Allied Products (SIC 28) industries, and Panel C excludes small firms, i.e., firms that have total assets smaller than the 10th percentile of the total assets in the country sample.

Panel A. Countries with more than 1,000 Listed Firms

	Full Sample (1)	Report R&D (2)	Report Patent (3)	R&D and Patent (4)	Differences		
					(5) = ((1)-(2))/(1)	(6) = ((1)-(3))/(1)	(7) = ((1)-(4))/(1)
Ln(Total Assets)	6.92	7.40	7.73	7.46	-7%***	-12%***	-8%***
PPE	0.29	0.25	0.23	0.20	14%***	21%***	31%***
Tobin's Q	1.29	1.38	1.65	1.88	-7%***	-28%***	-46%***
Leverage	0.80	0.53	0.56	0.48	34%***	30%***	40%***
Capital Expenditure	0.06	0.05	0.05	0.05	17%***	17%***	17%***
ROA	0.01	0.00	-0.01	-0.03	100%***	200%***	400%***
Sales Growth	0.26	0.23	0.27	0.30	12%***	-4%	-15%***

Panel B. SIC 25 and 38 Industries

	Full Sample (1)	Report R&D (2)	Report Patent (3)	R&D and Patent (4)	Differences		
					(5) = ((1)-(2))/(1)	(6) = ((1)-(3))/(1)	(7) = ((1)-(4))/(1)
Ln(Total Assets)	5.76	5.99	6.02	6.17	-4%***	-5%***	-7%***
PPE	0.19	0.17	0.16	0.16	11%***	16%***	16%***
Tobin's Q	1.69	1.73	1.93	2.01	-2%	-14%***	-19%***
Leverage	0.56	0.52	0.44	0.40	7%	21%*	29%***
Capital Expenditure	0.05	0.04	0.04	0.04	20%***	20%***	20%***
ROA	-0.02	-0.03	-0.03	-0.02	-50%***	-50%**	0%
Sales Growth	0.26	0.25	0.24	0.25	4%	8%	4%

Panel C. Excluding Small Firms

	Full Sample (1)	Report R&D (2)	Report Patent (3)	R&D and Patent (4)	Differences		
					(5) = ((1)-(2))/(1)	(6) = ((1)-(3))/(1)	(7) = ((1)-(4))/(1)
Ln(Total Assets)	7.11	7.64	7.87	7.76	-7%***	-11%***	-9%***
PPE	0.29	0.25	0.24	0.21	14%***	17%***	28%***
Tobin's Q	1.57	1.47	1.72	1.78	6%	-10%*	-13%***
Leverage	0.53	0.49	0.49	0.45	8%***	8%***	15%***
Capital Expenditure	0.06	0.05	0.05	0.05	17%***	17%***	17%***
ROA	0.03	0.03	0.02	0.01	0%	33%***	67%***
Sales Growth	0.26	0.22	0.26	0.30	15%***	0%	-15%***

Table IA2
Predictability of Unreported Innovation with Lasso

The table presents the OLS regression results for predictability of unreported innovation using only Lasso variables. Columns (1)-(4) *World* present the results for all countries in the sample. Columns (5)-(7) present the results for the *US* only. Panel A presents the results for unreported R&D, and Panel B presents the results for non-USPTO patent seeking firms. Standard errors are double clustered at firm and time level. T-statistics are presented in brackets. Variable definitions are presented in Table A1. *, **, and *** represent significance at the 10%, 5%, and 1% levels, respectively. Adj. R² is the adjusted R².

Panel A. Unreported R&D

	World				US		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
Ln(Total Assets)	-0.020*** (-10.20)	-0.017*** (-8.31)	-0.012*** (-7.86)	-0.009*** (-3.54)	0.040*** (13.28)	0.004 (1.73)	-0.014*** (-4.40)
Stock Liquidity	-0.007*** (-8.92)	-0.007*** (-9.99)	-0.005*** (-14.33)	-0.001*** (-3.47)	-0.008*** (-13.50)	-0.003*** (-6.23)	-0.001*** (-3.27)
Patent Intensity	-604.300*** (-17.42)	-1.310 (-0.11)	13.370 (1.06)	38.122** (2.02)	-700.176*** (-21.33)	-19.498*** (-3.08)	-17.045* (-1.77)
Country FE			Yes				
Industry FE		Yes	Yes			Yes	
Year FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Firm FE				Yes			Yes
N	300,634	300,634	300,634	328,734.0	77,982	77,982	76,944
Adj. R ²	0.13	0.23	0.38	0.80	0.23	0.53	0.93

Panel B. Non-USPTO Patent Seeking Firms

Ln(Total Assets)	-0.012*** (-11.62)	-0.009*** (-9.41)	-0.020*** (-19.28)	-0.006*** (-6.88)	-0.000 (-0.08)	-0.023*** (-9.67)	-0.013*** (-4.03)
Stock Liquidity	-0.006*** (-20.77)	-0.007*** (-21.86)	-0.004*** (-17.74)	-0.001*** (-5.11)	-0.007*** (-14.73)	-0.005*** (-12.24)	-0.001*** (-3.55)
Patent Intensity	-0.000** (-2.09)	-0.000** (-2.10)	-0.000** (-2.22)	-0.000 (-1.14)	-0.000*** (-3.22)	-0.000** (-2.60)	0.000 (1.01)
R&D Stock	-368.647*** (-17.34)	-30.191** (-2.57)	-25.181*** (-3.23)	-4.871 (-0.90)	-554.522*** (-18.37)	-47.131*** (-3.72)	-0.239 (-0.02)
Country FE			Yes				
Industry FE		Yes	Yes			Yes	
Year FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Firm FE				Yes			Yes
N	327,997	327,997	327,997	326,067	77,958	77,958	76,926
R ²	0.09	0.15	0.24	0.76	0.15	0.32	0.78

Table IA3
Recovered R&D and Imputed Unreported R&D

The table presents the comparison of Recovered R&D with different imputation methods. Panel A presents the result for the MI estimated for the complete sample, Panel B present the results for the multiple imputation on the restated R&D sub sample and the industry and size matched peers. Panel C shows the correlation between text-based innovation measures and the full sample multiple imputation. *R&D* is the recovered R&D expenditure as reported in 10-K filings, *MI R&D M1* is the multiply imputed R&D using ln(total assets), ROA, PPE, sales growth and leverage, by industry (two-digit), *MI R&D M2* is the multiply imputed R&D using the same model as M1 with the addition of the lagged R&D as conditioning information. *MI R&D M3* is the multiply imputed R&D using the same model as M1 with the addition of the Lasso variables: stock liquidity and industry patent intensity, *MI R&D M4* is the multiply imputed R&D using the same model as M2 with the addition of the Lasso variables: stock liquidity and industry patent intensity. “Diff.” is the difference between Recovered R&D and imputed R&D. T-stats represent the t-statistic for the difference between Recovered R&D an imputed R&D. *, **, and *** represent significance at the 10%, 5%, and 1% levels, respectively.

MEAN	STD	RD	Diff.	t-stat	
<i>Panel A. R&D Full Sample</i>					
MI R&D M1	17.19	334.17	6.91	-10.28	-0.97
MI R&D M2	12.48	252.47	6.91	-5.56	-0.69
MI R&D M3	17.95	334.12	6.91	-11.04	-1.04
MI R&D M4	15.10	246.92	6.91	-8.19	-1.05
<i>Panel B. R&D Sub Sample</i>					
MI R&D M1	17.20	340.49	6.91	-10.29	-0.93
MI R&D M2	14.14	257.20	6.91	-7.23	-0.87
MI R&D M3	19.09	337.75	6.91	-12.17	-1.12
MI R&D M4	15.99	256.29	6.91	-9.07	-1.10
<i>Panel C. Correlation with Text-based Innovation</i>					
	Text-based Innovation	Text-based Negative Innovation			
MI R&D M1 Full Sample	0.30***	0.28***			
MI R&D M2 Full Sample	0.29***	0.26***			
MI R&D M3 Full Sample	0.30***	0.28***			
MI R&D M4 Full Sample	0.30***	0.26***			

Table IA4
Robustness of Simulation based on Simulated Data

This table provides robustness for the simulation based on simulated data, as described in section 5.2.3. The methods evaluated are listwise deletion (LD), imputation with zero (ImpZero), imputation with industry mean, two-digit SIC code (ImpMean), inverse probability weighting (IPW), Heckman procedure (Heckman), and multiple imputation (MI). The regressions for imputations using zero and the industry mean also include a dummy variable to denote the imputed observations. We present the bias (relative bias over true parameter) and root mean squared error (RMSE) for MAR for 70% missingness. In Panel A the correlation between η and ε is 0.6, in Panel B the correlation between z_1 and x_1 is 0.6, in Panel C the correlation between z_1 and ε is 0.4, but η is generated independently from ε to avoid direct endogeneity in the selection equation.

		LD	Imp Zero	Imp Mean	IPW	Heckman	MI
<i>Panel A. Large Correlation between Errors</i>							
Bias	θ_1	-0.20	-0.29	-0.29	-0.21	-0.11	-0.10
	θ_2	-0.18	0.08	0.08	-0.17	-0.13	-0.09
RMSE	θ_1	0.21	0.30	0.30	0.23	0.18	0.11
	θ_2	0.19	0.10	0.10	0.19	0.17	0.11
<i>Panel B. Large Correlation between z_1 and x_1</i>							
Bias	θ_1	-0.16	-0.22	-0.22	-0.19	-0.08	-0.06
	θ_2	-0.11	0.10	0.10	-0.11	-0.07	-0.03
RMSE	θ_1	0.20	0.25	0.25	0.22	0.20	0.11
	θ_2	0.16	0.12	0.12	0.15	0.16	0.09
<i>Panel C. Correlation between z_1 and ε, no Correlation between Errors</i>							
Bias	θ_1	-0.27	-0.36	-0.36	-0.31	-0.12	-0.01
	θ_2	-0.23	0.05	0.05	-0.21	-0.17	0.02
RMSE	θ_1	0.28	0.37	0.37	0.32	0.26	0.07
	θ_2	0.24	0.08	0.08	0.22	0.26	0.09

Table IA5
Patent Simulation Based on the Empirical Distribution of Data

This table provides the evaluation statistics, bias (relative bias over true parameter) and root mean squared error (RMSE) for the simulation based on the empirical distribution from Compustat (US) and USPTO data. The empirical distribution comes from the panel of 783 firms with non-missing information for all variables except USPTO patents. The methods evaluated are listwise deletion (LD), imputation with zero (ImpZero), imputation with industry mean, two-digit SIC code (ImpMean), inverse probability weighting (IPW), Heckman procedure (Heckman), and multiple imputation (MI). MI uses all the variables in sample and is estimated using MCMC with 200 iterations for convergence. The regressions for imputation with zero and industry mean include a dummy variable for the imputed observations. Absolute average represents the average of the absolute bias across all variables. Variable definitions are presented in Table A1. We present results for two missingness mechanisms: missing at random (MAR) in Panel A and missing not at random (MNAR) in Panel B. We generate missingness in patents for 70% of the sample. We conduct 500 simulations.

		LD	Imp Zero	Imp Mean	IPW	Heckman	MI
<i>Panel A. MAR</i>							
Bias	Patent	-14.77	-2.38	-2.88	16.86	-48.13	0.34
	Ln(Total Assets)	21.85	1.64	3.33	62.00	-15.40	2.51
	Tobin's Q	-7.95	1.47	1.13	-190.58	-238.28	1.23
	Leverage	-1.94	0.02	0.04	9.08	9.84	0.07
	ROA	0.41	0.11	0.10	3.95	2.66	0.11
	<i>Avg. Abs. Bias</i>	<i>9.38</i>	<i>1.13</i>	<i>1.50</i>	<i>56.50</i>	<i>62.86</i>	<i>0.85</i>
RMSE	Patent	0.00	0.00	0.00	0.00	0.00	0.00
	Ln(Total Assets)	0.06	0.02	0.03	0.12	0.09	0.03
	Tobin's Q	0.03	0.02	0.02	0.41	0.52	0.02
	Leverage	0.31	0.12	0.13	0.91	1.07	0.12
	ROA	0.51	0.32	0.33	2.96	2.42	0.32
	<i>Panel B. MNAR</i>						
Bias	Patent	-15.14	-0.22	0.05	10.25	-23.99	1.26
	Ln(Total Assets)	19.45	-1.00	-1.22	52.67	2.04	0.06
	Tobin's Q	-7.58	-0.34	-0.23	-106.94	-136.07	-0.39
	Leverage	-1.84	0.03	0.02	4.07	3.65	0.06
	ROA	0.43	0.01	-0.01	2.23	1.37	0.01
	<i>Avg. Abs. Bias</i>	<i>8.89</i>	<i>0.32</i>	<i>0.31</i>	<i>35.23</i>	<i>33.42</i>	<i>0.36</i>
RMSE	Patent	0.00	0.00	0.00	0.00	0.00	0.00
	Ln(Total Assets)	0.06	0.02	0.03	0.10	0.07	0.03
	Tobin's Q	0.03	0.02	0.02	0.23	0.30	0.02
	Leverage	0.29	0.13	0.13	0.42	0.57	0.13
	ROA	0.51	0.30	0.33	1.69	1.67	0.30