



RIETI Discussion Paper Series 17-E-068

Forecasting Firm Performance with Machine Learning: Evidence from Japanese firm-level data

MIYAKAWA Daisuke

Hitotsubashi University

MIYAUCHI Yuhei

MIT

Christian PEREZ

Carnegie Mellon University



Research Institute of Economy, Trade & Industry, IAA

The Research Institute of Economy, Trade and Industry

<http://www.rieti.go.jp/en/>

Forecasting Firm Performance with Machine Learning:
Evidence from Japanese firm-level data*

MIYAKAWA Daisuke

Hitotsubashi University

MIYAUCHI Yuhei

MIT

Christian PEREZ

Carnegie Mellon University

Abstract

The goal of this paper is to forecast future firm performance with machine learning techniques. Using data on over one million Japanese firms with supply-chain linkage information provided by a credit reporting agency, we show high performance in the prediction of exit, sales growth, and profit growth. In particular, our constructed proxies far outperform the credit score assigned by the credit reporting agency based on a detailed survey and interviews of firms. Against such baseline score, our models are able to ex-ante identify 16% of exiting firms (baseline: 11%), 25% of firms experiencing growth in sales (baseline: 8%), and 22% of firms exhibiting positive profit growth (baseline: 13%). The proof of concept of this paper provides practical usage of machine learning methods in firm performance prediction.

Keywords: Machine learning, Big data, Prediction, Firm exit, Firm growth

JEL classification: G31; L25

RIETI Discussion Papers Series aims at widely disseminating research results in the form of professional papers, thereby stimulating lively discussion. The views expressed in the papers are solely those of the author(s), and neither represent those of the organization to which the author(s) belong(s) nor the Research Institute of Economy, Trade and Industry.

*Miyakawa: Graduate School of International Corporate Strategy, Hitotsubashi University, dmiyakawa@ics.hit-u.ac.jp; Miyauchi: Department of Economics, MIT, miyauchi@mit.edu; Perez: School of Computer Science, Carnegie Mellon University, cperez1@andrew.cmu.edu. This research is conducted as a part of the Research Institute of Economy, Trade and Industry (RIETI) research project (Study on Corporate Finance and Firm Dynamics). We thank Hiroshi Ohashi, Miho Takizawa, Iichiro Uesugi, Kosuke Uetake, Yasutora Watanabe, Makoto Yano, and the seminar participants at RIETI for helpful suggestions. We also thank Katsuhiko Komatsu for his excellent research assistant work. Miyakawa gratefully acknowledges financial supports from the Grant-in-Aid for Scientific Research No. 16K03736 JSPS and the grant-in-aid from Zengin Foundation for Studies on Economics and Finance. Miyauchi acknowledges the Nakajima Foundation for the financial support for their Ph.D scholarship.

1 Introduction

Prediction of firms' future performance is a central mandate for many stakeholders. First and foremost, it is a crucial activity for for-profit activity of banks, investors, and supply chain management. It is also crucial from a policy perspective. For example, modern banking regulation (e.g., Basel) requires banks to construct their internal model for evaluating client firms' credit worthiness, which reflects the estimates of client firms' future performance.

In the practical process of examining and predicting firms' future performance, credit reporting agencies play a crucial role. Credit reporting agencies are entities that collect and survey firms and provide the information for commercial purposes. Examples of credit reporting agencies include Dunn and Bradstreet in the US, Experian in European countries, and Tokyo Shoko Research Ltd. (TSR) in Japan. In addition to providing raw information such as financial statements, they typically create a score that summarizes the overall performance of the firm. These scores are typically constructed from both observable firm characteristics and financial statements (i.e. "hard" information) and in-depth interviews based on owner characteristics, reputation, and growth opportunity (i.e. "soft" information). The score is used for various purposes; e.g. evaluating the credit worthiness of client firms, screening on transaction partners, and understanding overall market environment.

Traditionally, credit reporting agencies have relied on their own (often confidential) algorithm to construct the scores. For example, the score by TSR in Japan is the summation of (i) the ability of owner (max: 20 points) based on the business attitude, experience, their asset condition, (ii) the growth possibility (max: 25 points) based on past sales growth, the growth of profit, and the characteristics of products, and (iii) stability (max: 45 points) based on firm age, stated-capital, financial statement information, room of collateral provision, real and financial transaction relationships and (iv) reputation (max 10 points) based on the level of disclosure and overall reputation, with further detail not disclosed. Although this set of information is intuitive, it is not immediately clear whether these particular variables or weights are optimal for the construction of a score to predict the future performance of firms.

A recent revolution of machine learning techniques opens up a scope to tackle such a problem possibly more accurately, systematically and in a non-arbitrary manner. Machine learning is the study of efficient and accurate prediction using models which summarize potential sets of predictors. It is used in different contexts such as the prediction of crime in a specific area, mechanical failure in a plant, and weather forecasts.

The goal of this paper is to apply machine learning techniques to predict various future firm performance measures (i.e., firm exit, sales growth, and profit growth) and compare their predictive power with the score assigned by the credit reporting agency. Toward this end, we utilize the firm-level data from TSR in Japan, which consist of a subset of firm characteristics, supply chain linkage information, as well as its score assigned by the TSR, of nearly all firms that TSR covers in 2006, 2011 and 2014.

We find that, although the score constructed by TSR has reasonably good predictive performance, particularly for exit prediction, combining firm-level characteristics to the score with machine learning far out-performs the score. In particular, against such a baseline score, our models are able to ex-ante identify 16% of exiting firms (baseline: 11%), 25% of firms experiencing growth in sales (baseline: 8%), and 22% of firms exhibiting positive profit growth

(baseline: 13%). These results suggest the usefulness of machine learning methods in firm performance prediction.

The rest of the paper proceeds as follows. Section 2 describes the data we use for our analysis. Section 3 explains the empirical methodology and Section 4 presents and discusses the results. Section 5 concludes.

2 Background and Data

2.1 Tokyo Shoko Research

Throughout the paper, we use the datasets provided by TSR (Tokyo Shoko Research Ltd.), one of the largest credit reporting agencies in Japan. TSR is a private company operating in the areas of credit research, publishing, and database distribution. The central product TSR provides is the unsolicited-basis company report accounting for the performance of each targeted firm, which they sell to a variety of clients including banks, security firms, non-financial enterprises, and governmental organizations.

A typical report consists of more than ten pages and includes firms' basic characteristics and financial statement information. The clients of TSR purchase the reports for various reasons; e.g. evaluating the credit worthiness of client firms, screening on transaction partners, and understanding the overall market environment.

Among the items reported in the company report, a proxy computed by TSR to summarize the performance of firms, which we call as "*fscore*", is provided. We will describe this score in detail in the following section.

2.2 Data

In this section, we will go over the data we use in the present study. All data is obtained from TSR through the support of Research Institute of Economy, Trade and Industry (RIETI), which is a governmental research institute affiliated with Japanese Ministry of Economy, Trade and Industry.

2.2.1 Overview

Our main data source is a panel of Japanese firm data accounting for firms' basic performance information (e.g., sales, number of employees, stated capital, and profit) as well as basic characteristics including company owner characteristics, precise geographic location, firm age, etc., for 2006, 2011, and 2014. For some variables (e.g. sales, profit, dividend), the data records the information in the preceding year. The data records around 800,000 firms in 2006 and reaches nearly 1,700,000 firms in 2014.

In addition to the firm-level characteristics, the dataset also includes linked firm-firm pair-level data accounting for firms' supply chain network. As discussed in, for example, Acemoglu et al. (2015), which suggests firm-level shocks are transmitted through a network of interconnections in the economy, it is reasonable to presume that this supply chain network information has predictive power.

2.2.2 Firm Performance Indicators

We consider three firm performance indicators to be predicted: firm exit, sales growth, and profit growth. Each outcome variable is defined for two time intervals; from 2006 to 2011, and

2011 to 2014. We use information from 2006 to predict outcomes defined for 2006 to 2011, and information from 2011 to predict outcomes for 2011 to 2014.

Firm Exit

We define firm exit in any subsequent panel periods if firms exited the market for the following four reasons reported by TSR: bankruptcies, closure, dissolution, and suspension.

Sales Growth

To characterize firms which exhibit high sales growth relative to other firms in the same industry, we prepare a dummy variable that takes 1 if the sales growth in the subsequent panel periods exceeds the average plus one standard deviation within the same 2-digit industry.

Profit Growth

We prepare a dummy variable that takes 1 if the profit growth in the subsequent panel periods exceeds the average plus one standard deviation within the same 2-digit industry. In this analysis, we restrict our data to firms which realize positive profits.

2.2.3 Predictors

We use four categories of predictors to predict the firm performance measures described in the previous subsection: (1) the score constructed by TSR (*fscore*), (2) firms' own characteristics, (3) geography and industry-related variables, and (4) supply-chain network related variables. Here we overview the variables categorized in each group below. The Appendix 1 describes the full list of variables.

(1) Solvency Score

The score (*fscore*) takes values between 0 and 100. The number is computed as the sum of the four sub-scores accounting for (i) the ability of owner (max: 20 points) based on the business attitude, experience, their asset condition, and so on, (ii) the growth possibility (max: 25 points) based on past sales growth, the growth of profit, the characteristics of products, and so on, (iii) stability (max: 45 points) based on firm age, stated-capital, financial statement information, room of collateral provision, real and financial transaction relationships, and so on, and (iv) reputation (max 10 points) based on the level of disclosure and overall reputation. We only have access to the *fscore*, but not the decomposition of each component. The variable should reflect some of the predictors we use in this analysis, but not necessarily all since we do not have the full information that TSR obtains (e.g. detailed financial statements and "soft" information from the interviews).

TSR guidelines provide the following categorization of *fscore* ranges: a. caution required (scores 29 and under), b. medium caution required (scores between 30 and 49), c. little caution required (scores between 50 and 64), d. no specific concern (scores between 65 and 79), and e. no concern at all for scores 80 and above. We should note that while the categorization exists, the score is highly concentrated around 50 (as shown in Table 1). This implies that there could be large room for the score to be improved for the purpose of accurate prediction of firm performance.

(2) *Own-Firm Characteristics*

As predictors accounting for firms' own characteristics, we use the information stored in financial statements and firms' attributes. The former consists of firm size measured by their sales and its change, their profit (loss or not) and its change, the number of employees, their stated capital, and the status of dividend payment and its change. The latter consists of firm age, owner age, the number of establishments, and their listed status.

(3) *Industry and Geographic Information*

As predictors accounting for the industry and area which firms belong to, we set up the following two groups of variables. First, we construct the two variables measuring the average sales growth of firms located in the same city as the targeted (i.e., we compute a score for) firms as well as the average sales growth of firms belonging to the same industry classified in the 2-digit level. Second, we also employ dummy variables representing the 2-digit industry classification as well as dummy variables representing the large classification used in Japan Standard Industry Classification.

(4) *Supply-Chain Linkage Information*

As predictors accounting for the supply chain network, we construct the following two groups of variables. First, we compute widely used network metrics for each firm by using the supply chain network information. The metrics consist of degree centrality, eigenvector centrality, egonet eigenvalue, co-transaction, and the number of direct (i.e., customers and suppliers) and indirect (i.e., suppliers' suppliers, suppliers' customers, customers' customers, and customers' suppliers) transaction partners. Second, we construct a large number of variables accounting for the characteristics of transaction partners. To summarize this information, we employ an average, maximum, minimum and the sum of various attributes associated with transaction partners. Note that while the network metrics cover both the direct and indirect transaction partners, the transaction partners' characteristics only consider the direct transaction partners.

Regarding the extant studies using the same dataset as ours and focusing on the economic implication of supply chain network structure, Fu and Ogura (2017) reports that the sensitivity of lending rate with respect to the score is lower for the firms more closely connected to other firms. Their finding implies that position in the supply chain network represents additional information to firms' own characteristics. Also, focusing on the economic implication of the exogenous change in transaction partners' characteristics, Calvalho et al. (2016) finds that the damage due to a natural disaster to transaction partners causes a reduction in the sales of firms transacting with the damaged firms. Their finding implies that the characteristics of transaction partners contain information useful for the purpose of evaluating the performance of firms. These extant studies share a motivation to examine the economic implication of the information embedded in the supply chain network with the present study. Nonetheless, the largest difference between the present study and them lies in the fact that we aim at incorporating such information for predicting firm performance. Using a machine learning method, which has not been extensively used in the economics literature, we explicitly compare the relative importance of those additional variables.

3 Method

We utilize state-of-the-art machine learning methods for developing our prediction model. Our particular problem of predicting relatively rare firm exit events (which occur with 7%

probability) falls in to the class of “imbalanced label prediction” tasks in computer science. Following the literature, we apply weighted random forest, a minority-class oversampling method.

3.1 Weighted Random Forest

Random forests models aggregate many individual decision tree models, each trained on a randomly selected sample from the training data. Particularly for predicting rare events, Chen et al. (2004) develop an extension of random forest, called weighted random forest. Intuitively, the method weighs data corresponding to minority event (e.g. exit) much more heavily than that corresponding to the majority event (e.g. non-exit).

3.2 Variable Selection

For the purpose of increasing the out-of-sample predictive performance of our models, we utilize Lasso variable selection for removing noisy variables not contributing to predictive performance. We select variables for each outcome variable.

The Appendix 2 describes the list of variables selected through Lasso. The three panels of Table 1 show the summary statistics of all the selected variables in the three categories (i.e., firms’ own characteristics, geography- and industry-specific variables, and network variables).

3.3 Measuring Prediction Performance

In our baseline exercise, we train models with the realization of outcome variables from 2006 to 2011 using the information available at 2006, and conduct out-of-sample prediction of the realization of outcome variables from 2011 to 2014 using the information available at 2011.

We utilize the ROC curve to evaluate the predictive performance of the model. Our tasks of binary exit, growth, and profit growth classification require the setting of thresholds for which predicted probabilities surpassing this level will indicate a positive binary outcome. Given a fixed model, the ROC curve plots the true and false positive rates corresponding to the varying of this threshold value. Without any predictors (i.e. random guess), the curve should trace the 45-degree line, and curves closer to the top-left corner are desirable (maximize true positive rate and minimize false positive rate). With this motivation, it is conventional to also summarize the ROC curve by the area under the curve, called AUC.

In addition to evaluating performance over different threshold values to observe the trade-off between true and false positive rates, we find it useful to compare true positive rates fixing the threshold value corresponding to the realized proportion of events (i.e. exit, high sales growth and profit growth). This threshold choice provides a practical exposition of this model comparison.

4 Results

Figure 1 shows the four ROC curves in the case of exit prediction. We evaluate four models. The first model uses only *fscore* while the second model additionally contains firms’ own characteristics selected as predictors. The third model further includes the geography- and industry-specific variables selected as predictors. The fourth model includes all the selected variables including the ones in network variables. All ROC curves are based on the out-of-

sample prediction of the exit from 2011 to 2014, using the exit from 2006 to 2011 as a training data.

We observe that while the model solely using *fscore* as the predictor performs well as it is located well-above the 45 degree line, it can be also confirmed that our constructed proxy outperforms *fscore* in terms of predictive power for firm exit. As the three lines corresponding to the model 2 to 4 are largely overlapping, we can also confirm that own firm characteristics explain most of the additional predictive power while geographic and supply-chain related information does not add much predictive power.

Figure 2 and Figure 3 show the ROC curves of firm sales and profit growth. Two points are noteworthy: first, the predictive power of *fscore* in the case of sales prediction is low compared to exit (Figure 1). This might reflect the fact that *fscore* is designed to be used as an early warning indicator. Second, partly reflecting the poor performance of *fscore* for sales growth prediction, the gain in the predictive power obtained from adding variables to the model is larger than exit (Figure 1).

Table 2 summarizes the performance of our prediction model. In the first column, we show the size of area under the ROC curve (AUC) for each curve, which confirms our finding based on the figures we showed above. In the second column, we show the true positive rate fixing the threshold value corresponding to the realized proportion of events (i.e. exit, high sales growth and profit growth). We can see that, against the prediction solely on *fscore*, our models are able to ex-ante identify 16% of exiting firms (baseline: 11%), 25% of firms experiencing growth in sales (baseline: 8%), and 22% of firms exhibiting positive profit growth (baseline: 13%). These results again confirm a practical importance of machine learning methods in firm performance prediction.

These results provide an important insight. Namely, although *fscore* is constructed from various firm-level information as well as detailed survey conducted by TSR investigators, there is still a large room for the score to be improved in terms of the predictive power for firm exit. As the computation algorithm for *fscore* is confidential, we cannot identify exactly what generates such a difference. One potential factor leading to the current result is that our constructed proxy might use information not used for the computation of *fscore*. Another possibility is that the weight we assign to each variable is chosen in a better way than for *fscore*. In either case, the obtained result suggests the usefulness of using machine learning methods for the purpose of predicting firm performance.

5 Conclusion and Future Work

In this paper, we apply machine learning techniques to over a million Japanese firms' data to predict future firm performance. First, for each of exit, sales growth and profit growth, our constructed proxies out-perform the credit score assigned by the credit reporting agency based on a detailed survey and interviews of firms. In particular, the predictions of sales and profit growth are improved largely. Second, the predictors chosen by a machine learning method vary depending on the choice of the performance measure we predict. Third, own firm characteristics explain most of the predictive power while geographic and supply-chain related information does not add much predictive power. The proof of concept of this paper provides a practical usage of machine learning methods in firm performance prediction.

The analysis in the present study could be expanded toward various directions. First, we are planning to implement the “true” out-of-sample test using more recent firm performance data (e.g., exit after 2014). This additional analysis allows us to rigorously test our models by committing ourselves to only information currently available to us. Second, we can use more detailed network-related and accounting information to further improve our models. This is inspired by recent studies such as Acemoglu et al. (2015) which discuss the economic implications of geographical and supply chain network information.

References

Acemoglu, D., U. Akcigit, and W. Kerr. 2015. Networks and the Macroeconomy: An Empirical Exploration. Working paper.

Carvalho, V. M., M. Nirei, Y. U. Saito, A. Tahbaz-Salehi. 2016. Supply Chain Disruptions: Evidence from the Great East Japan Earthquake. Working paper.

Fu, J., and Y. Ogura. 2017. Product Network Connectivity and Information for Loan Pricing. RIETI Discussion Paper Series 17-E-028.

Chen C, Liaw A, Breiman L. 2004. Using Random Forest to Learn Imbalanced Data. Technical Report 666 Statistics Department of Univeristy of California at Berkley.

Figures and Tables

Figure 1. ROC curve for Exit Prediction

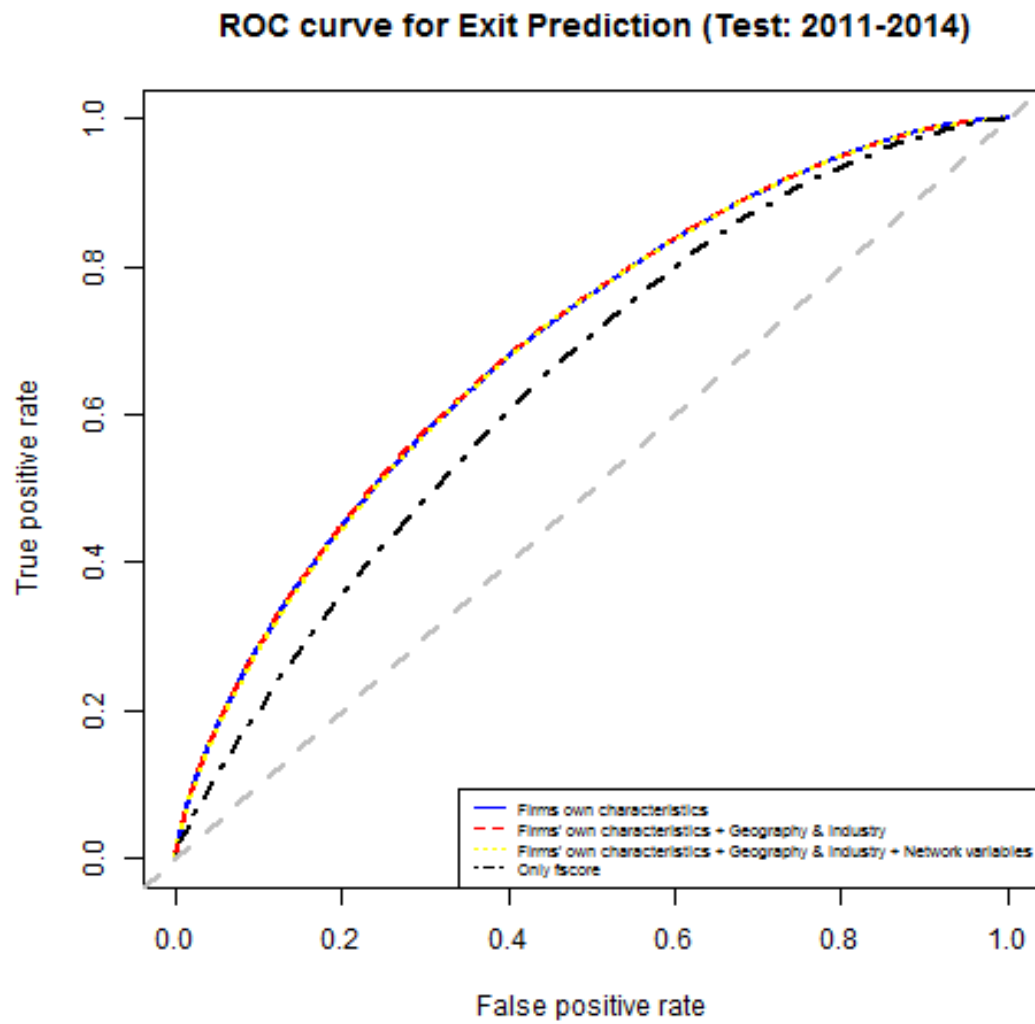


Figure 2. ROC curve for Sales Growth Prediction

ROC curve for Growth Prediction (Test: 2011-2014)

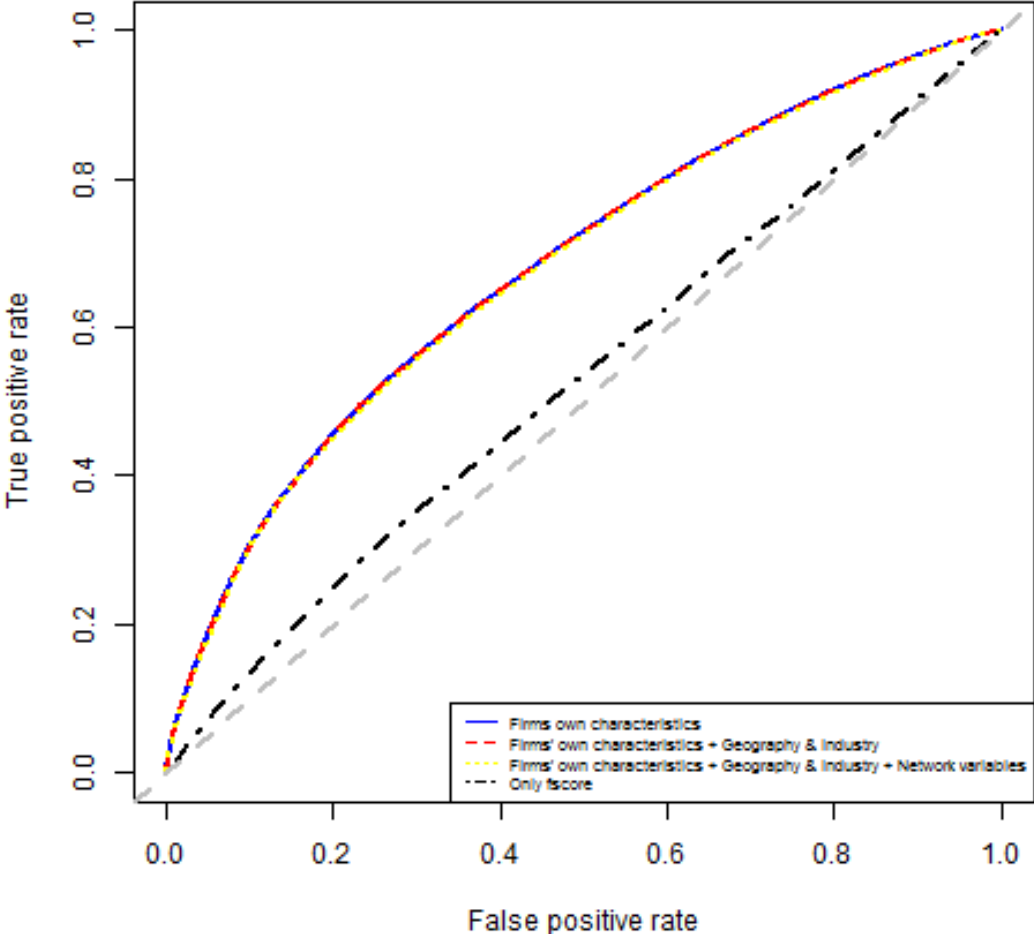


Figure 3. ROC curve for Profit Growth Prediction

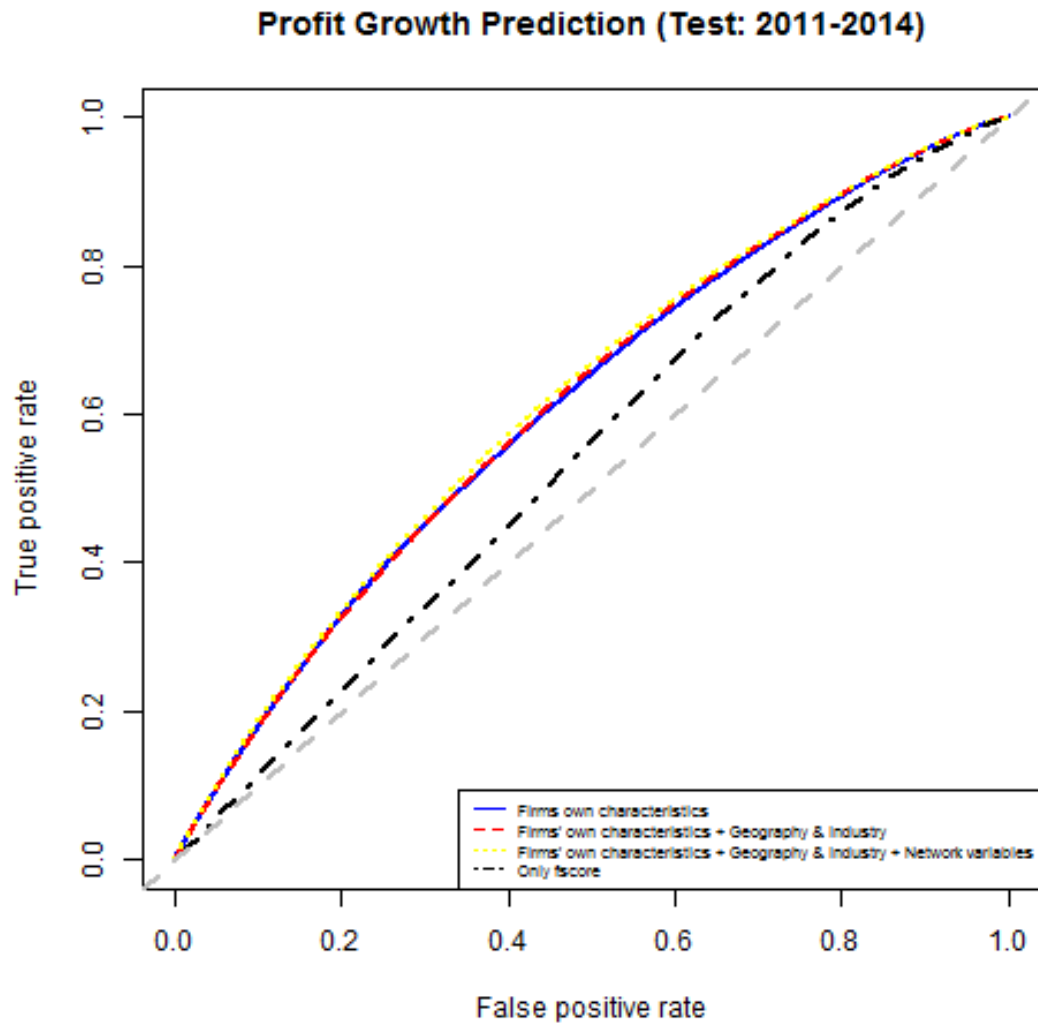


Table 1. Summary Statistics of LASSO selected variables

(a) Firms' own characteristics

	count	mean	standard deviation	min	p50	max
<i>exit</i>	1,763,750	0.073	0.261	0	0	1
<i>growth_sales</i>	1,507,839	0.085	0.278	0	0	1
<i>growth_profit</i>	595,547	0.131	0.338	0	0	1
<i>fscore</i>	1,759,351	48.294	6.010	0	48	91
<i>logsales</i>	1,763,750	11.911	1.661	0	11.775	23.837
<i>logsalesdif</i>	1,763,750	-0.029	0.340	-11.522	-0.002	10.784
<i>div</i>	1,763,750	0.028	0.166	0	0	1
<i>divdif</i>	1,763,750	-0.002	0.104	-1	0	1
<i>deficit</i>	1,261,802	0.216	0.412	0	0	1
<i>deficitdif</i>	1,184,895	-0.007	0.460	-1	0	1
<i>firmage</i>	1,607,580	29.168	15.481	0	27	136
<i>ceoage</i>	1,604,065	60.042	10.760	4	61	114
<i>logemp</i>	1,750,424	2.032	1.293	0	1.792	12.476
<i>logcap</i>	1,572,942	9.297	1.196	0	9.210	22.757
<i>listed</i>	1,763,750	0.004	0.064	0	0	1
<i>single</i>	1,333,470	0.225	0.418	0	0	1
N	1,763,750					

(b) Geography- and Industry-specific variables

	count	mean	sd	min	p50	max
<i>avg_salesgrowth_city</i>	1,763,750	-0.012	0.014	-0.297	-0.014	0.112
<i>avg_salesgrowth_industry_2d</i>	1,763,750	-0.014	0.014	-0.043	-0.015	0.201
N	1,763,750					

(b) Geography- and Industry-specific variables

	count	mean	sd	min	p50	max
nsup	1,419,014	5.288	31.644	1	3	7227
ncus	1,294,124	5.944	32.285	1	3	10734
nsup_supmean	1,366,741	138.142	412.791	1	22	7475
nsup_cusmean	1,294,124	332.976	818.372	1	40.5	7926
ncus_supmean	1,419,014	339.283	838.762	1	61	12057
ncus_cusmean	1,209,472	145.222	434.513	1	16.5	12057
logsalesdif_supmean	1,402,188	-0.013	0.211	-7.237	-0.010	10.341
logsalesdif_cusmean	1,244,852	-0.021	0.232	-7.896	-0.014	10.341
deficit_supmean	1,349,303	0.186	0.289	0	0	1
deficit_cusmean	1,202,640	0.185	0.280	0	0	1
deficit_supmin	1,349,303	0.071	0.257	0	0	1
deficitdif_cusmean	1,189,406	-0.029	0.338	-1	0	1
listed_supmin	1,419,014	0.037	0.190	0	0	1
listed_cusmean	1,294,124	0.159	0.276	0	0	1
div_cusmean	1,254,725	0.150	0.295	0	0	1
logsalesdif_cusmax	1,244,852	0.173	0.446	-7.896	0.083	10.784
logsalesdif_cusmin	1,244,852	-0.213	0.411	-10.973	-0.118	10.341
logsalesdif_supmin	1,402,188	-0.172	0.356	-10.439	-0.095	10.341
deficitdif_cusmax	1,189,406	0.246	0.519	-1	0	1
div_supmax	1,414,127	0.271	0.445	0	0	1
div_cusmax	1,254,725	0.261	0.439	0	0	1
divdif_supmean	1,414,127	-0.018	0.144	-1	0	1
divdif_supmin	1,414,127	-0.107	0.320	-1	0	1
listed_cusmax	1,294,124	0.360	0.480	0	0	1
listed_cusmin	1,294,124	0.056	0.229	0	0	1
single_supmean	1,346,608	0.146	0.256	0	0	1
N	1,728,428					

Table 2. Main Results

	AUC	True Positive Rate With Fixed Threshold ¹
(a) Exit Prediction		
Probit: F-Score Only	0.65	11%
Weighted R.F.: Adding Own-Firm	0.7	16%
Weighted R.F.: Adding Geo/Industry	0.7	16%
Weighted R.F.: Adding Network	0.7	16%
(b) Sales Growth Prediction		
Probit: F-Score Only	0.53	8%
Weighted R.F.: Adding Own-Firm	0.68	25%
Weighted R.F.: Adding Geo/Industry	0.68	25%
Weighted R.F.: Adding Network	0.68	25%
(c) Profit Growth Prediction		
Probit: F-Score Only	0.55	13%
Weighted R.F.: Adding Own-Firm	0.61	22%
Weighted R.F.: Adding Geo/Industry	0.61	22%
Weighted R.F.: Adding Network	0.62	22%

¹Estimates from fixing the binary-classification threshold on predicted probabilities corresponding to the realized proportions of the events, as described in Section 3.3. The resulting model comparison represents a practical application of our approach.

Appendix 1: Full list of variables

In the appendix, we provide a list of all the variables we use as inputs for variable selection. There are three categories of variables and 200 variables in total which we include in our prediction model through Lasso.

Variable name	Definition of the variable	Year 2006-2011	Year 2011-2014
<u>Variable category: TSR score</u>			
<i>fscore</i>		2005	2010
<u>Variable category: Firms' own characteristics</u>			
<i>logsales</i>	Log of sales	2005	2010
<i>logsalesdif</i>	Log difference of sales	2004-05	2009-10
<i>div</i>	Dummy for dividend payment	2005	2010
<i>div_pre</i>	Dummy for dividend payment	2004	2010
<i>divdif</i>	Difference of <i>div</i>	2004-05	2009-10
<i>deficit</i>	Dummy for loss	2005	2010
<i>deficitdif</i>	Difference of <i>deficit</i>	2004-05	2009-10
<i>logemp</i>	Log of number of employees	2005	2010
<i>firmage</i>	Firm age	2005	2010
<i>CEOage</i>	CEO age	2005	2010
<i>logcap</i>	Log of stated capital	2005	2010
<i>single</i>	Dummy for single establishment	2005	2010
<i>listed</i>	Dummy for listed firm	2005	2010
<u>Variable category: Geography-specific variables & Industry-specific variables</u>			
<i>avg_salesgrowth_city</i>	Average sales growth of the firms located in the same city	2005	2010
<i>avg_salesgrowth_industry_2d</i>	Average sales growth of the firms belonging to the same 2-digit industry	2005	2010

<i>ind_code_2d</i>	Dummy variable associated with 2-digit industry code	2005	2010
<i>ind_code_L</i>	Dummy variable associated with a large classification used in Japan Standard Industry Classification	2005	2010
<i>logsalesdif2</i>	Log difference of sales	2004-05	2009-10
<u>Variable category: Network characteristics associated with supply chain network & Transaction partners' information</u>			
<i>d centrality</i>	Degree centrality	2005	2010
<i>ev centrality</i>	Eigen vector centrality	2005	2010
<i>egonet_evalue</i>	Egonet eigenvalue	2005	2010
<i>cotransaction</i>	Co-transaction	2005	2010
<i>nsup</i>	Number of direct suppliers	2005	2010
<i>ncus</i>	Number of direct customers	2005	2010
<i>ncus_supmean</i>	Average number of suppliers' customers	2005	2010
<i>ncus_supmax</i>	Maximum number of suppliers' customers	2005	2010
<i>ncus_supmin</i>	Minimum number of suppliers' customers	2005	2010
<i>ncus_supsum</i>	Total number of suppliers' customers	2005	2010
<i>nsup_supmean</i>	Average number of suppliers' suppliers	2005	2010
<i>nsup_supmax</i>	Maximum number of suppliers' suppliers	2005	2010
<i>nsup_supmin</i>	Minimum number of suppliers' suppliers	2005	2010
<i>nsup_supsum</i>	Total number of suppliers' suppliers	2005	2010
<i>single_supmean</i>	Average of <i>single</i> over direct suppliers	2005	2010
<i>single_supmax</i>	Maximum of <i>single</i> over direct suppliers	2005	2010

<i>single_supmin</i>	Minimum of <i>single</i> over direct suppliers	2005	2010
<i>single_supsum</i>	Sum of <i>single</i> over direct suppliers	2005	2010
<i>logcap_supmean</i>	Average of <i>logcap</i> over direct suppliers	2005	2010
<i>logcap_supmax</i>	Maximum of <i>logcap</i> over direct suppliers	2005	2010
<i>logcap_supmin</i>	Minimum of <i>logcap</i> over direct suppliers	2005	2010
<i>logcap_supsum</i>	Sum of <i>logcap</i> over direct suppliers	2005	2010
<i>CEOage_supmean</i>	Average of <i>CEOage</i> over direct suppliers	2005	2010
<i>CEOage_supmax</i>	Maximum of <i>CEOage</i> over direct suppliers	2005	2010
<i>CEOage_supmin</i>	Minimum of <i>CEOage</i> over direct suppliers	2005	2010
<i>CEOage_supsum</i>	Sum of <i>CEOage</i> over direct suppliers	2005	2010
<i>logprofitdif_2_supmean</i>	Average log difference of profit among suppliers	2003-04	
<i>logprofitdif_2_supmax</i>	Maximum of log difference of profit among suppliers	2003-04	
<i>logprofitdif_2_supmin</i>	Minimum of log difference of profit among suppliers	2003-04	
<i>logprofitdif_2_supsum</i>	Sum of log difference of profit among suppliers	2003-04	
<i>logprofitdif_supmean</i>	Average log difference of profit among suppliers	2004-05	2009-10
<i>logprofitdif_supmax</i>	Maximum of log difference of profit among suppliers	2004-05	2009-10
<i>logprofitdif_supmin</i>	Minimum of log difference of profit among suppliers	2004-05	2009-10
<i>logprofitdif_supsum</i>	Sum of log difference of profit among suppliers	2004-05	2009-10

<i>deficitdif_pre_supmean</i>	Average of <i>deficitdif</i> among suppliers	2003-04	
<i>deficitdif_pre_supmax</i>	Maximum of <i>deficitdif</i> among suppliers	2003-04	
<i>deficitdif_pre_supmin</i>	Minnimum of <i>deficitdif</i> among suppliers	2003-04	
<i>deficitdif_pre_supsum</i>	Sum of <i>deficitdif</i> among suppliers	2003-04	
<i>deficitdif_supmean</i>	Average of <i>deficitdif</i> among suppliers	2004-05	2009-10
<i>deficitdif_supmax</i>	Maximum of <i>deficitdif</i> among suppliers	2004-05	2009-10
<i>deficitdif_supmin</i>	Minimum of <i>deficitdif</i> among suppliers	2004-05	2009-10
<i>deficitdif_supsum</i>	Sum of <i>deficitdif</i> among suppliers	2004-05	2009-10
<i>deficit_supmean</i>	Average of <i>deficit</i> among suppliers	2005	2010
<i>deficit_supmax</i>	Maximum of <i>deficit</i> among suppliers	2005	2010
<i>deficit_supmin</i>	Minimum of <i>deficit</i> among suppliers	2005	2010
<i>deficit_supsum</i>	Sum of <i>deficit</i> among suppliers	2005	2010
<i>logemp_supmean</i>	Average <i>logemp</i> among suppliers	2005	2010
<i>logemp_supmax</i>	Maximum of <i>logemp</i> among suppliers	2005	2010
<i>logemp_supmin</i>	Minimum of <i>logemp</i> among suppliers	2005	2010
<i>logemp_supsum</i>	Sum of <i>logemp</i> among suppliers	2005	2010
<i>divdif_pre_supmean</i>	Average <i>divdif</i> among suppliers	2003-04	
<i>divdif_pre_supmax</i>	Maximum of <i>divdif</i> among suppliers	2003-04	
<i>divdif_pre_supmin</i>	Minimum of <i>divdif</i> among suppliers	2003-04	

<i>divdif_pre_supsum</i>	Sum of <i>divdif</i> among suppliers	2003-04	
<i>divdif_supmean</i>	Average <i>divdif</i> among suppliers	2004-05	2009-10
<i>divdif_supmax</i>	Maximum of <i>divdif</i> among suppliers	2004-05	2009-10
<i>divdif_supmin</i>	Minimum of <i>divdif</i> among suppliers	2004-05	2009-10
<i>divdif_supsum</i>	Sum of <i>divdif</i> among suppliers	2004-05	2009-10
<i>logsalesdif_2_supmean</i>	Average of <i>logsalesdif</i> among suppliers	2003-04	
<i>logsalesdif_2_supmax</i>	Maximum of <i>logsalesdif</i> among suppliers	2003-04	
<i>logsalesdif_2_supmin</i>	Minimum of <i>logsalesdif</i> among suppliers	2003-04	
<i>logsalesdif_2_supsum</i>	Sum of <i>logsalesdif</i> among suppliers	2003-04	
<i>logsalesdif_supmean</i>	Average of <i>logsalesdif</i> among suppliers	2004-05	2009-10
<i>logsalesdif_supmax</i>	Maximum of <i>logsalesdif</i> among suppliers	2004-05	2009-10
<i>logsalesdif_supmin</i>	Minimum of <i>logsalesdif</i> among suppliers	2004-05	2009-10
<i>logsalesdif_supsum</i>	Sum of <i>logsalesdif</i> among suppliers	2004-05	2009-10
<i>logsales_supmean</i>	Average of <i>logsales</i> among suppliers	2005	2010
<i>logsales_supmax</i>	Maximum of <i>logsales</i> among suppliers	2005	2010
<i>logsales_supmin</i>	Minimum of <i>logsales</i> among suppliers	2005	2010
<i>logsales_supsum</i>	Sum of <i>logsales</i> among suppliers	2005	2010
<i>firmage_supmean</i>	Average <i>firmage</i> among suppliers	2005	2010

<i>firmage_supmax</i>	Maximum of <i>firmage</i> among suppliers	2005	2010
<i>firmage_supmin</i>	Minimum of <i>firmage</i> among suppliers	2005	2010
<i>firmage_supsum</i>	Sum of <i>firmage</i> among suppliers	2005	2010
<i>CEOfemale_supmean</i>	Average <i>CEOfemale</i> among suppliers	2005	2010
<i>CEOfemale_supmax</i>	Maximum of <i>CEOfemale</i> among suppliers	2005	2010
<i>CEOfemale_supmin</i>	Minimum of <i>CEOfemale</i> among suppliers	2005	2010
<i>CEOfemale_supsum</i>	Sum of <i>CEOfemale</i> among suppliers	2005	2010
<i>div_supmean</i>	Average <i>div</i> among suppliers	2005	2010
<i>div_supmax</i>	Maximum of <i>div</i> among suppliers	2005	2010
<i>div_supmin</i>	Minimum of <i>div</i> among suppliers	2005	2010
<i>div_supsum</i>	Sum of <i>div</i> among suppliers	2005	2010
<i>listed_supmean</i>	Average <i>listed</i> among suppliers	2005	2010
<i>listed_supmax</i>	Maximum of <i>listed</i> among suppliers	2005	2010
<i>listed_supmin</i>	Minimum of <i>listed</i> among suppliers	2005	2010
<i>listed_supsum</i>	Sum of <i>listed</i> among suppliers	2005	2010
<i>fscore_supmean</i>	Average <i>fscore</i> among suppliers	2005	2010
<i>fscore_supmax</i>	Maximum of <i>fscore</i> among suppliers	2005	2010
<i>fscore_supmin</i>	Minimum of <i>fscore</i> among suppliers	2005	2010
<i>fscore_supsum</i>	Sum of <i>fscore</i> among suppliers	2005	2010
<i>ncus_cusmean</i>	Average number of customer' customers	2005	2010
<i>ncus_cusmax</i>	Maximum number of customer' customers	2005	2010

<i>ncus_cusmin</i>	Minimum number of customer' customers	2005	2010
<i>ncus_cussum</i>	Total number of customer' customers	2005	2010
<i>nsup_cusmean</i>	Average number of customer' suppliers	2005	2010
<i>nsup_cusmax</i>	Maximum number of customer' suppliers	2005	2010
<i>nsup_cusmin</i>	Minimum number of customer' suppliers	2005	2010
<i>nsup_cussum</i>	Total number of customer' suppliers	2005	2010
<i>positive_growth_cusmean</i>	Average <i>positive_growth</i> among customers	2004-05	2009-10
<i>positive_growth_cusmax</i>	Maximum <i>positive_growth</i> among customers	2004-05	2009-10
<i>positive_growth_cusmin</i>	Minimum <i>positive_growth</i> among customers	2004-05	2009-10
<i>positive_growth_cussum</i>	Sum of <i>positive_growth</i> among customers	2004-05	2009-10
<i>single_cusmean</i>	Average of <i>single</i> over direct customers	2005	2010
<i>single_cusmax</i>	Maximum of <i>single</i> over direct customers	2005	2010
<i>single_cusmin</i>	Minimum of <i>single</i> over direct customers	2005	2010
<i>single_cussum</i>	Sum of <i>single</i> over direct customers	2005	2010
<i>logcap_cusmean</i>	Average of <i>logcap</i> over direct customers	2005	2010
<i>logcap_cusmax</i>	Maximum of <i>logcap</i> over direct customers	2005	2010
<i>logcap_cusmin</i>	Minimum of <i>logcap</i> over direct customers	2005	2010
<i>logcap_cussum</i>	Sum of <i>logcap</i> over direct customers	2005	2010

<i>CEOage_cusmean</i>	Average of <i>CEOage</i> over direct customers	2005	2010
<i>CEOage_cusmax</i>	Maximum of <i>CEOage</i> over direct customers	2005	2010
<i>CEOage_cusmin</i>	Minimum of <i>CEOage</i> over direct customers	2005	2010
<i>CEOage_cussum</i>	Sum of <i>CEOage</i> over direct customers	2005	2010
<i>logprofitdif_2_cusmean</i>	Average log difference of profit among customers	2003-04	2010
<i>logprofitdif_2_cusmax</i>	Maximum of log difference of profit among customers	2003-04	2010
<i>logprofitdif_2_cusmin</i>	Minimum of log difference of profit among customers	2003-04	2010
<i>logprofitdif_2_cussum</i>	Sum of log difference of profit among customers	2003-04	2010
<i>logprofitdif_cusmean</i>	Average log difference of profit among customers	2004-05	2009-10
<i>logprofitdif_cusmax</i>	Maximum of log difference of profit among customers	2004-05	2009-10
<i>logprofitdif_cusmin</i>	Minimum of log difference of profit among customers	2004-05	2009-10
<i>logprofitdif_cussum</i>	Sum of log difference of profit among customers	2004-05	2009-10
<i>deficitdif_pre_cusmean</i>	Average of <i>deficitdif</i> among customers	2003-04	
<i>deficitdif_pre_cusmax</i>	Maximum of <i>deficitdif</i> among customers	2003-04	
<i>deficitdif_pre_cusmin</i>	Minimum of <i>deficitdif</i> among customers	2003-04	
<i>deficitdif_pre_cussum</i>	Sum of <i>deficitdif</i> among customers	2003-04	
<i>deficitdif_cusmean</i>	Average of <i>deficitdif</i> among customers	2004-05	2009-10
<i>deficitdif_cusmax</i>	Maximum of <i>deficitdif</i> among customers	2004-05	2009-10

<i>deficitdif_cusmin</i>	Minimum of <i>deficitdif</i> among customers	2004-05	2009-10
<i>deficitdif_cussum</i>	Sum of <i>deficitdif</i> among customers	2004-05	2009-10
<i>deficit_cusmean</i>	Average of <i>deficit</i> among customers	2005	2010
<i>deficit_cusmax</i>	Maximum of <i>deficit</i> among customers	2005	2010
<i>deficit_cusmin</i>	Minimum of <i>deficit</i> among customers	2005	2010
<i>deficit_cussum</i>	Sum of <i>deficit</i> among customers	2005	2010
<i>logemp_cusmean</i>	Average <i>logemp</i> among customers	2005	2010
<i>logemp_cusmax</i>	Maximum of <i>logemp</i> among customers	2005	2010
<i>logemp_cusmin</i>	Minimum of <i>logemp</i> among customers	2005	2010
<i>logemp_cussum</i>	Sum of <i>logemp</i> among customers	2005	2010
<i>divdif_pre_cusmean</i>	Average <i>divdif</i> among customers	2003-04	
<i>divdif_pre_cusmax</i>	Maximum of <i>divdif</i> among customers	2003-04	
<i>divdif_pre_cusmin</i>	Minimum of <i>divdif</i> among customers	2003-04	
<i>divdif_pre_cussum</i>	Sum of <i>divdif</i> among customers	2003-04	
<i>divdif_cusmean</i>	Average <i>divdif</i> among customers	2004-05	2009-10
<i>divdif_cusmax</i>	Maximum of <i>divdif</i> among customers	2004-05	2009-10
<i>divdif_cusmin</i>	Minimum of <i>divdif</i> among customers	2004-05	2009-10
<i>divdif_cussum</i>	Sum of <i>divdif</i> among customers	2004-05	2009-10
<i>logsalesdif_2_cusmean</i>	Average of <i>logsalesdif</i> among customers	2003-04	

<i>logsalesdif_2_cusmax</i>	Maximum of <i>logsalesdif</i> among customers	2003-04	
<i>logsalesdif_2_cusmin</i>	Minimum of <i>logsalesdif</i> among customers	2003-04	
<i>logsalesdif_2_cussum</i>	Sum of <i>logsalesdif</i> among customers	2003-04	
<i>logsalesdif_cusmean</i>	Average of <i>logsalesdif</i> among customers	2004-05	2009-10
<i>logsalesdif_cusmax</i>	Maximum of <i>logsalesdif</i> among customers	2004-05	2009-10
<i>logsalesdif_cusmin</i>	Minimum of <i>logsalesdif</i> among customers	2004-05	2009-10
<i>logsalesdif_cussum</i>	Sum of <i>logsalesdif</i> among customers	2004-05	2009-10
<i>logsales_cusmean</i>	Average of <i>logsales</i> among customers	2005	2010
<i>logsales_cusmax</i>	Maximum of <i>logsales</i> among customers	2005	2010
<i>logsales_cusmin</i>	Minimum of <i>logsales</i> among customers	2005	2010
<i>logsales_cussum</i>	Sum of <i>logsales</i> among customers	2005	2010
<i>firmage_cusmean</i>	Average <i>firmage</i> among customers	2005	2010
<i>firmage_cusmax</i>	Maximum of <i>firmage</i> among customers	2005	2010
<i>firmage_cusmin</i>	Minimum of <i>firmage</i> among customers	2005	2010
<i>firmage_cussum</i>	Sum of <i>firmage</i> among customers	2005	2010
<i>CEOmale_cusmean</i>	Average <i>CEOmale</i> among customers	2005	2010
<i>CEOmale_cusmax</i>	Maximum of <i>CEOmale</i> among customers	2005	2010
<i>CEOmale_cusmin</i>	Minimum of <i>CEOmale</i> among customers	2005	2010

<i>CEOfemale_cussum</i>	Sum of <i>CEOfemale</i> among customers	2005	2010
<i>div_cusmean</i>	Average <i>div</i> among customers	2005	2010
<i>div_cusmax</i>	Maximum of <i>div</i> among customers	2005	2010
<i>div_cusmin</i>	Minimum of <i>div</i> among customers	2005	2010
<i>div_cussum</i>	Sum of <i>div</i> among customers	2005	2010
<i>listed_cusmean</i>	Average <i>listed</i> among customers	2005	2010
<i>listed_cusmax</i>	Maximum of <i>listed</i> among customers	2005	2010
<i>listed_cusmin</i>	Minimum of <i>listed</i> among customers	2005	2010
<i>listed_cussum</i>	Sum of <i>listed</i> among customers	2005	2010
<i>fscore_cusmean</i>	Average <i>fscore</i> among customers	2005	2010
<i>fscore_cusmax</i>	Maximum of <i>fscore</i> among customers	2005	2010
<i>fscore_cusmin</i>	Minimum of <i>fscore</i> among customers	2005	2010
<i>fscore_cussum</i>	Sum of <i>fscore</i> among customers	2005	2010
<i>fscore_ssmean</i>	Average <i>fscore</i> among suppliers' suppliers	2005	2010
<i>fscore_ccmean</i>	Average <i>fscore</i> among customers' customers	2005	2010
<i>fscore_scmean</i>	Average <i>fscore</i> among suppliers' customers	2005	2010
<i>fscore_csmean</i>	Average <i>fscore</i> among customers' suppliers	2005	2010

Appendix 2: List of selected variables

In the appendix, we provide a list of the variables selected through Lasso.

<p>Variable name</p> <p><u>Variable category: Firms' own characteristics for exit, sales growth, and profit growth</u></p> <p><i>logsales, logsalesdif, div, divdif, deficit, difcitdif, firmage, CEOage, logemp, logcap, listed, single.</i></p>
<p><u>Variable category: Geography-specific variables & Industry-specific variables for exit, sales growth, and profit growth</u></p> <p><i>avg_salesgrowth_city, avg_salesgrowth_industry_2d.</i></p>
<p><u>Variable category: Network characteristics associated with supply chain network & Transaction partners' information for exit</u></p> <p><i>nsup, ncus, nsup_supmean, nsup_cusmean, ncus_supmean, ncus_cusmean, logsalesdif_supmean, logsalesdif_cusmean.</i></p>
<p><u>Variable category: Network characteristics associated with supply chain network & Transaction partners' information for sales growth</u></p> <p><i>logsalesdif_supmean, deficit_supmean, deficit_cusmean, deficitdif_supmin, deficitdif_cusmean, listed_supmin, listed_cusmean, div_cusmean.</i></p>
<p><u>Variable category: Network characteristics associated with supply chain network & Transaction partners' information for profit growth</u></p> <p><i>deficit_supmean, deficit_cusmean, listed_supmin, and listed_cusmean, logsalesdif_cusmax, logsalesdif_cusmin, logsalesdif_supmin, deficitdif_cusmax, div_supmax, div_cusmax, divdif_supmean, divdif_supmin, listed_cusmax, listed_cusmin, single_supmean.</i></p>