

Marcus Baulig

Machine Learning in Statistics

Forecasting Applications in
Corporate Finance

Machine Learning in Statistics – Forecasting Applications in Corporate Finance

Dissertation

zur Erlangung des Grades eines Doktors der Wirtschaftswissenschaft

(doctor rerum politicarum)

der Fakultät für Empirische Humanwissenschaften und

Wirtschaftswissenschaft der Universität des Saarlandes

vorgelegt von

Marcus Baulig

Tag der Disputation: 14. Juni 2019
Dekan: Univ.-Prof. Dr. Stefan Strohmeier
Erstberichterstatter: Univ.-Prof. Dr. Ashok Kaul
Zweitberichterstatter: Univ.-Prof. Dr. Dieter Hess

Berichte aus der Statistik

Marcus Baulig

Machine Learning in Statistics

Forecasting Applications in Corporate Finance

Shaker Verlag
Düren 2019

Bibliographic information published by the Deutsche Nationalbibliothek

The Deutsche Nationalbibliothek lists this publication in the Deutsche Nationalbibliografie; detailed bibliographic data are available in the Internet at <http://dnb.d-nb.de>.

Zugl.: Saarbrücken, Univ., Diss., 2019

Copyright Shaker Verlag 2019

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, without the prior permission of the publishers.

Printed in Germany.

ISBN 978-3-8440-6820-7

ISSN 1619-0963

Shaker Verlag GmbH • Am Langen Graben 15a • 52353 Düren

Phone: 0049/2421/99011-0 • Telefax: 0049/2421/99011-9

Internet: www.shaker.de • e-mail: info@shaker.de

Abstract

Given the growing awareness of machine learning outside of computer science in both academia and business in recent years, I examine two corporate financing forecasting problems and show how machine learning models compare to established literature models. The forecasting of future corporate bankruptcies serves vicariously as a classification problem and the forecasting of future corporate earnings as a regression problem, which are both traditionally approached with econometric techniques such as logistic and linear regression.

I forecast bankruptcies and earnings using a diverse set of machine learning models—ranging from subset selection models to highly flexible Boosting models—and combine these models to stacked ensembles. Ensemble learning can potentially outstrip the performance of individual machine learning models and to date has not been investigated for bankruptcy and earnings forecasts. Besides the focus on high performing models, I create highly interpretable logistic and linear regression models with the most predictive variables assessed over all machine learning models.

For both forecast problems, stacked ensembles show in their optimal calibration the best forecast performance, exceeding established literature models by at least 5% in terms of standard evaluation criteria. Individual machine learning models achieve at least an improvement of 4% and my self-generated regression models at least 1%.

Despite the already capable performance of established literature models in the field of corporate finance already, the utilization of machine learning successfully enriches the way researchers can transform data into forecasts. The engineering of new data with machine learning methods has the potential to be an even more influential approach for future researchers.

Acknowledgments

I would like to express my special gratitude to Ashok Kaul, my first advisor, for the numerous academic and professional opportunities he has provided that have pointed the way ahead for me. Among those opportunities, I wish to highlight my time abroad in China, the US, and South Africa, all of which have found their way into this thesis. I am also very grateful to Dieter Hess, my second advisor, for his outstanding supervision from the very beginning of my PhD studies.

My colleagues Christian Agethen, Nathalie Neu-Yanders, and Manuel Schieler provided me not only with inspiring academic ideas, but also with a very pleasurable atmosphere at the chair. I truly appreciate you all and the time we have spent together. A special thanks also goes to Andrew Yanders for his proofreading of my thesis.

Finally, I would like to thank everyone who has accompanied me on a personal level these last years: most of all my parents, Felix Malzahn, Matthias Schindler, and Jani Coetzee.

Contents

List of Figures	vii
List of Tables	ix
List of Abbreviations	xi
1 Preface	1
2 Overview of Machine Learning Methods in Statistics	6
2.1 Introduction	6
2.2 General Concepts at Machine Learning	10
2.2.1 Trade-offs at Machine Learning	10
2.2.2 Cross-Validation	13
2.3 Classification and Regression Techniques	14
2.3.1 Discriminant Analysis	15
2.3.2 Logistic Regression	17
2.3.3 Linear Regression	17
2.4 Machine Learning Methods	18
2.4.1 Subset Selection: Best Subset and Stepwise Selection	18
2.4.2 Shrinkage Methods: Ridge Regression, Lasso and Elastic Net	20
2.4.3 Spline-Based Methods: Multivariate Adaptive Regression Splines and Generalized Additive Models	24
2.4.4 K-Nearest Neighbors	27
2.4.5 Support Vector Machines	29
2.4.6 Neural Networks	32
2.4.7 Tree-Based Methods: Decision Trees, Boosting, Random Forests	34
2.5 Combining Machine Learning Methods: Ensemble Learning	39
2.6 Conclusion	41

3	Classification Problem: Bankruptcy Forecast	43
3.1	Introduction	43
3.2	Bankruptcy Forecast Literature Review	44
3.2.1	Statistical Bankruptcy Forecast Models	45
3.2.2	Theoretical Bankruptcy Forecast Models	46
3.2.3	Machine Learning Bankruptcy Forecast Studies	47
3.2.4	Methodological Issues at Bankruptcy Research	48
3.3	Data and Method	49
3.4	Empirical Analysis	57
3.4.1	Bankruptcy Forecast with Machine Learning Methods	57
3.4.2	Stacked Ensemble Learning	63
3.4.3	Variable Importance Analysis	66
3.4.4	Accounting vs. Market Variables	72
3.4.5	Extended Forecast Horizons	72
3.5	Conclusion	76
4	Regression Problem: Earnings Forecast	80
4.1	Introduction	80
4.2	Earnings Forecast Literature Review	82
4.2.1	Econometric Earnings Forecast Models	82
4.2.2	Machine Learning Studies	85
4.3	Data and Method	85
4.4	Empirical Analysis	90
4.4.1	Earnings Forecast with Machine Learning Methods	90
4.4.2	Stacked Ensemble Learning	94
4.4.3	Variable Importance Analysis	95
4.4.4	Accounting vs. Market Variables	102
4.4.5	Extended Forecast Horizons	104
4.5	Conclusion	107
5	Conclusion and Outlook	109
	Appendices	112
	Bibliography	120

List of Figures

Figure 2.1:	Overview of Implemented Machine Learning Methods	8
Figure 2.2:	Trade-off between Interpretability and Flexibility	11
Figure 2.3:	Optimal Fit of a Model	12
Figure 2.4:	Concept of 5-fold Cross-Validation	14
Figure 2.5:	Linear Discriminant Analysis	16
Figure 2.6:	Shrinkage Parameter λ	21
Figure 2.7:	Variable Selection at Shrinkage Methods	23
Figure 2.8:	Regression Splines	25
Figure 2.9:	K-Nearest Neighbors	28
Figure 2.10:	Support Vector Classifier and Support Vector Machine	30
Figure 2.11:	Neural Network Architecture	32
Figure 2.12:	Decision Tree	36
Figure 2.13:	Ensemble Learning	41
Figure 3.1:	Receiver Operating Characteristics (ROC) Curve	56
Figure 3.2:	Variable Importance of Bankruptcy Predictors	67
Figure 4.1:	Variable Importance of Earnings Predictors	97

List of Tables

Table 3.1:	Mean Comparison of Bankruptcy Predictors by Bankruptcy Status	51
Table 3.2:	Performance Evaluation of Bankruptcy Forecast using Machine Learning Methods	58
Table 3.3:	Performance Evaluation of Bankruptcy Forecast with Extreme Gradient Boosting per Year	60
Table 3.4:	Correlations of Machine Learning Bankruptcy Forecast Models . . .	63
Table 3.5:	Performance Evaluation of Bankruptcy Forecast using Ensemble Learning Methods	65
Table 3.6:	Performance Evaluation of Bankruptcy Forecast with Variable Importance Model	69
Table 3.7:	Parameter Estimates of the Bankruptcy Variable Importance Model	70
Table 3.8:	Performance Evaluation of Bankruptcy Forecast using Five Most Important Variables	71
Table 3.9:	Performance Evaluation of Bankruptcy Forecast using Accounting Variables	73
Table 3.10:	Performance Evaluation of Bankruptcy Forecast with Two-Year-Ahead Forecast Horizon ($t+2$)	74
Table 3.11:	Performance Evaluation of Bankruptcy Forecast with Three-Year-Ahead Forecast Horizon ($t+3$)	75
Table 4.1:	Summary Statistics of Earnings Predictors	87
Table 4.2:	Performance Evaluation of Earnings Forecast using Machine Learning Methods	91
Table 4.3:	Performance Evaluation of Earnings Forecast with Random Forests per Year	93
Table 4.4:	Correlations of Machine Learning Earnings Forecast Models	94
Table 4.5:	Performance Evaluation of Earnings Forecast using Ensemble Learning Methods	96
Table 4.6:	Performance Evaluation of Earnings Forecast with Variable Importance Models	99
Table 4.7:	Parameter Estimates of the Earnings Variable Importance Models .	101

Table 4.8:	Performance Evaluation of Earnings Forecast using Four Most Important Variables	102
Table 4.9:	Performance Evaluation of Earnings Forecast using Accounting Variables	103
Table 4.10:	Performance Evaluation of Earnings Forecast with Two-Year-Ahead Forecast Horizon ($t+2$)	105
Table 4.11:	Performance Evaluation of Earnings Forecast with Three-Year-Ahead Forecast Horizon ($t+3$)	106
Table A.1:	Summary of Machine Learning Methods	113
Table A.2:	Bankruptcies per Year	117
Table A.3:	Performance Evaluation of Bankruptcy Forecast with Jones et al. (2017) Variables	118
Table A.4:	Performance Evaluation of Bankruptcy Forecast with a Random 70%/30% Data Split	119

List of Abbreviations

AI	Artificial Intelligence
AMEX	American Stock Exchange
ARIMA	Autoregressive Integrated Moving Average (Model)
AUC	Area under the (ROC) Curve
BOOST.GEN	Generalized Boosting
BOOST.XG	Extreme Gradient Boosting
BS	Breir-Score
BSS	Best Subset Selection
BSTEP	Backward Stepwise Selection
CV	Cross-validation
CRSP	Center for Research in Security Prices
D	Cross-deviance
Diff.	Difference
ENET	Elastic Net Regression
EP	Earnings Persistence Model (Li and Mohanram, 2014)
FPR	False Positive Rate
FSTEP	Forward Stepwise Selection
GAM	Generalized Additive Models
GDP	Gross Domestic Product
HVZ	Hou et al. (2012)
iid	Independent and identically distributed
Ind.-Spec.	Industry-Specific
KNN	K-Nearest Neighbors
Lasso	Least Absolute Shrinkage and Selection Operator (Lasso Regression)
LCI	Lower Confidence Interval
LDA	Linear Discriminant Analysis
LOG	Logistic Regression
LOOCV	Leave-One-Out Cross-Validation
LM	Li and Mohanram (2014)

MARS	Multiple Adaptive Regression Splines
MSE	Mean Squared Error
NASDAQ	National Association of Securities Dealers Automated Quotations
NN	Neural Networks
NYSE	New York Stock Exchange
OLS	Ordinary Least Squares Regression
PS	Per Share
RF	Random Forests
Ridge	Ridge Regression
RIM	Residual Income Model (Li and Mohanram, 2014)
ROC	Receiver Operating Characteristics (Curve)
RSS	Residual Sum of Squares
SVM	Support Vector Machine
TPR	True Positive Rate
TREE	Decision Tree
UCI	Upper Confidence Interval
VAR	Vector Autoregressive (Model)
VIM	Variable Importance Model
VIMnl	Variable Importance Model (non-linear)
WRDS	Wharton Research Data Service