# Discussion Papers

## An Evaluation of Early Warning Models for Systemic Banking Crises: Does Machine Learning Improve Predictions?

Johannes Beutel, Sophia List, Gregor von Schweinitz

## Authors

**Johannes Beutel**
Deutsche Bundesbank
E-mail: johannes.beutel@bundesbank.de

**Sophia List**
Deutsche Bundesbank
E-mail: sophia.list@bundesbank.de

**Gregor von Schweinitz**
Halle Institute for Economic Research (IWH) –
Member of the Leibniz Association,
Department of Macroeconomics, and
Leipzig University, Institute for Theoretical
Economics
E-mail: gregorvon.schweinitz@iwh-halle.de
Tel +49 345 7753 744

Comments and suggestions on the methods
and results presented are welcome.

IWH Discussion Papers are indexed in
RePEc-EconPapers and in ECONIS.

## Editor

# An Evaluation of Early Warning Models for Systemic Banking Crises: Does Machine Learning Improve Predictions?*

## Abstract

This paper compares the out-of-sample predictive performance of different early warning models for systemic banking crises using a sample of advanced economies covering the past 45 years. We compare a benchmark logit approach to several machine learning approaches recently proposed in the literature. We find that while machine learning methods often attain a very high in-sample fit, they are outperformed by the logit approach in recursive out-of-sample evaluations. This result is robust to the choice of performance measure, crisis definition, preference parameter, and sample length, as well as to using different sets of variables and data transformations. Thus, our paper suggests that further enhancements to machine learning early warning models are needed before they are able to offer a substantial value-added for predicting systemic banking crises. Conventional logit models appear to use the available information already fairly effciently, and would for instance have been able to predict the 2007/2008 financial crisis out-of-sample for many countries. In line with economic intuition, these models identify credit expansions, asset price booms and external imbalances as key predictors of systemic banking crises.

*Keywords: early warning system, logit, machine learning, systemic banking crises*

*JEL classification: C35, C53, G01*

# 1 Introduction

The global financial crisis has spurred a new wave of research on the importance of a stable financial system for macroeconomic stability. New early warning models for financial crises have been developed and are being employed by central banks to monitor the stability of the financial system and to guide macroprudential policy (see, for example European Central Bank, 2010, 2017; Drehmann and Juselius, 2014). Given the high costs associated with financial crises, it is important to understand the circumstances under which countries are likely to experience them and to provide accurate early warning signals of these events. Failing to activate macroprudential policy tools in time might lead to large costs for taxpayers, policymakers, and society as a whole, while issuing false alarms might lead to costly over-regulation of the financial system.[1]

Recently, early warning models that rely on machine learning methods have been proposed as an alternative to the traditionally employed methods in this field, such as the signaling approach (e.g. Kaminsky and Reinhart, 1999; Knedlik and von Schweinitz, 2012) and discrete choice (probit or logit) models (e.g. Frankel and Rose, 1996; Lo Duca and Peltonen, 2013). For instance, Alessi and Detken (2018) as well as Tanaka, Kinkyo, and Hamori (2016) have argued that random forests may improve early warning predictions in comparison to the logit model and the signaling approach. Holopainen and Sarlin (2017) have extended this argument to at least four other machine learning methods, namely artificial neural networks, support vector machines, k-nearest-neighbors, and decision trees.

Using a comprehensive dataset encompassing systemic banking crises for 15 advanced economies over the past 45 years, we compare the out-of-sample prediction accuracies of the logit model to four machine learning methods employed in the existing literature (random forest, support vector machines, k-nearest neighbors, and decision trees). We come to an interesting and perhaps surprising conclusion: simple logit models systematically outperform all machine learning methods considered under a large variety of circumstances. In particular, we show that, while machine learning methods are able to achieve near perfect in-sample fit, they perform worse than the logit model in recursive out-of-sample prediction, and often even worse than a naïve benchmark. This result is remarkable, as it cautions against the use of machine learning methods whose impressive in-sample performance may backfire in the context of actual out-of-sample forecasting situations.

We subject our key result to a variety of tests. First, we document the superiority of logit models for different combinations of leading indicator variables as well as for different measures of prediction accuracy. Second, we perform standard robustness checks, such as different data transformations, crisis databases, estimation periods, and parameterizations. Finally, we propose a bootstrap as a uniform approach to account for estimation uncertainty, allowing us to establish statistically significant differences in performance between methods. Moreover, we seek to determine ex ante optimal hyperparameters for machine learning methods using a computationally intensive re-sampling procedure (a specific variant of cross-validation). Even with this considerable effort, machine learning methods still generate out-of-sample predictions which are inferior to those of the logit

---

[1]The costs of financial crises are documented, for instance, in Jordà, Schularick, and Taylor (2011), and Laeven and Valencia (2012). An overview of internationally employed macroprudential policy tools can be found in Lim, Costa, Columba, Kongsamut, Otani, Saiyid, Wezel, and Wu (2011), Cerutti, Claessens, and Laeven (2017) or Claessens (2015).

model.

We suggest an explanation for this result and compare our findings to other studies using machine learning methods for predicting financial crises. Machine learning methods typically contain a much larger number of parameters than the logit model and are able to flexibly approximate a large space of functions. This allows them to fit in-sample data quite closely, but, at the same time, entails the risk of an overfit, and as a consequence weak out-of-sample performance. We provide empirical and theoretical arguments to show that this risk appears to materialize in the early warning context.

We complement our horse race with a detailed discussion of the best forecasting model (a logit model with 10 predictor variables). First, we provide insights into the economic variables driving the predictions of the best model. This illustrates the interpretability of logit models, e.g. in terms of coefficient signs and marginal effects, as an additional advantage relative to machine learning models. Second, we show that estimated coefficients of the logit model across time and prediction performance across different forecasting setups are remarkably stable. Finally, we discuss the role of prediction uncertainty for policymaking and show that error rates can be further reduced if one is willing to focus exclusively on significant signals.

Our paper is related to the new wave of research on early warning models spurred by the global financial crisis of 2008, as, for instance, in Alessi and Detken (2011); Rose and Spiegel (2012); Gourinchas and Obstfeld (2012); Lo Duca and Peltonen (2013); Drehmann and Juselius (2014). These papers construct different early warning models but do not consider machine learning methods or horse races between different methods. Going further, random forests are used by Alessi and Detken (2018) in early warning models of systemic banking crises at the country level, and by Tanaka et al. (2016) to predict failures at the level of individual banks. However, neither paper performs a systematic out-of-sample comparison of the random forest relative to other methods. By contrast, Holopainen and Sarlin (2017) run out-of-sample comparisons of several methods. Yet, they do so on a dataset containing a relatively small number of crisis episodes. We build on their pioneering work, but refine it in several important ways, namely regarding our careful construction of datasets and robustness checks, as well as our bootstrap and hyperparameter selection schemes, taking into account cross-sectional and serial dependence structures. We show how our out-of-sample results differ with respect to their paper and provide a potential explanation for this difference.

Overall, we add to this strand of the literature by conducting extensive out-of-sample model evaluations on a sample of 15 advanced economies covering 22 systemic banking crises over the period 1970-2016. Using this comprehensive sample and a variety of techniques enables us to refine the conclusions of earlier studies. We complement recent assessments of machine learning methods in other fields, for instance, regarding the closely related task of predicting civil wars in political science (Neunhoeffer and Sternberg, 2018). We thereby seek to contribute to a realistic assessment of the strengths and limitations of the various methods, and to stimulate further research in this area.

# 2 Methodology

## 2.1 The Early Warning Setup

In line with the recent early warning literature (see Drehmann and Juselius, 2014; Alessi and Detken, 2018; Holopainen and Sarlin, 2017) we estimate the probability of a financial crisis starting between the next 5 to 12 quarters (conditional on not already being in an acute crisis period)based on a set of potential early warning indicators. Details on this window forecasting approach and the resulting definition of the dependent variable for the estimations are given in Appendix A.1.

To inform decision-making, the estimated probabilities may be mapped into binary signals. Using a threshold $\tau$, the signal is set to one if the probability exceeds $\tau$ and to zero otherwise. These signals or their absence will ex post turn out as right or wrong, and can be classified into true positives, false positives, true negatives, or false negatives as indicated in Table 1. False negatives FN (also called type-1 errors) are observations where no signal is given during an early warning window (missed crises), while false positives FP (type-2 errors) result from observations where a signal is given outside of an early warning window (false alarms). A higher classification threshold $\tau$ implies fewer signals, reducing both true and false positives. Converting probabilities into signals via a threshold thus entails a trade-off between type-1 errors (missed crises) and type-2 errors (false alarms). Selecting the optimal threshold generally depends on the loss function of the forecast user. The current standard is to choose the classification threshold $\tau$ such that it maximizes the relative usefulness function of Alessi and Detken (2011), which weighs errors (as a share of the respective actual class) by a parameter $\mu$ representing the forecast user's preferences.[2]

Table 1: A contingency matrix.

|  |  | Actual class $C$ | |
| --- | --- | --- | --- |
|  |  | Pre-crisis period | Tranquil period |
| Predicted class $S$ | Signal | Correct call *True positive (TP)* | False alarm *False positive (FP)* |
|  | No signal | Missed crisis *False negative (FN)* | Correct silence *True negative (TN)* |

*Note:* This contingency matrix follows Holopainen and Sarlin (2017).

## 2.2 Estimation Methods

We employ the following methods for estimating crisis probabilities: Logistic regression, k-nearest neighbors, decision trees, random forests and support vector machines. Appendix A.2 provides an overview of these methods as well as details on their implementation. The selection of methods follows the previous literature (Berg and Pattillo, 1999; Bussière and

---

[2]In an out-of-sample experiment, it is important to choose the threshold using in-sample information only and not by maximizing out-of-sample performance, which is a priori unknown to the forecaster.

Fratzscher, 2006; Alessi and Detken, 2018; Holopainen and Sarlin, 2017).[3]

While binary choice models such as the logit are standard tools in the early warning literature, machine learning methods are sometimes thought to allow for stronger non-linearities and more flexible distributional assumptions, which might be beneficial when forecasting extreme events such as systemic banking crises.[4] We have a panel dataset with observations for several countries at different points in time. In order to treat methods uniformly and keep the setup parsimonious, we estimate each method on the same pooled sample of observations (pooling observations in the cross-section and time dimension).

The machine learning methods come with hyperparameters that have to be set exogenously prior to estimation. For example, the k-nearest neighbor method has a single exogenous hyperparameter, k, determining the number of neighbors to consider. Following the standard in the literature (James, Witten, Hastie, and Tibshirani, 2013; Murphy, 2012), hyperparameters for all methods are chosen such that they optimize a performance criterion (e.g. relative usefulness) in a cross-validation exercise. We have implemented a fairly sophisticated cross-validation algorithm taking into account the cross-sectional and serial correlation present in our dataset (see appendix section A.4 for details). As a consequence, cases where the performance of machine learning methods falls short of the logit approach cannot be easily attributed to sub-optimal hyperparameters, but appear to be more deeply rooted in the given model. A list of optimized hyperparameters can be found in the Appendix (see table A.2).

We deliberately do not use cross-validation to evaluate models. As Neunhoeffer and Sternberg (2018) show, using cross-validation both for hyperparameter selection and model evaluation may lead to serious over-estimation of (machine learning) model performance. Thus, we use cross-validation only for hyperparameter selection and perform a classic out-of-sample prediction experiment to evaluate models. To ensure a strict separation between in-sample and out-of-sample data, our cross-validation routine uses information before the start of the out-of-sample window only (i.e. data before 2005Q3).

## 2.3 Evaluating Predictions

For every observation, the early warning models estimate the probability of a crisis starting in the following five to twelve quarters. Based on this probability a binary signal can be derived. The performance of an early warning model can therefore be evaluated either

---

[3]We have also experimented with artificial neural networks (ANN), which – set up properly – need many (random) initial guesses for ANN weights to solve a global optimization problem (Shalev-Shwartz and Ben-David, 2014, Chap. 20.6). Consequently, when combining ANN with the re-sampling required by our hyperparameter selection and bootstrap schemes, computation time of ANN becomes extraordinarily high compared to the other methods. As neither our tentative results, nor the results by Holopainen and Sarlin (2017) suggest significant gains of (single hidden layer) ANN relative to the other methods, we do not include them in the analysis presented in this paper. Recent advances in the literature on deep neural networks point to a potentially interesting avenue for future research (LeCun, Bengio, and Hinton, 2015), although overfitting may pose challenges for these models as well.

[4]Other commonly mentioned advantages of machine learning methods relate to potential benefits on large datasets ("big data"), and their ability to deal with a large number of potentially relevant variables. In the context of early warning models however, datasets typically contain only a limited number of observations. Moreover, the amount of variables can generally still be challenging for some machine learning methods (see discussion of methods' properties in appendix section A.2 and empirical results for our early warning application 4.2).

with respect to signals or probabilities. We employ three different performance measures that are standard in the literature, namely the above-mentioned relative usefulness ($U_r$), the area under the curve (AUC) and the Brier probability score (BPS).

The relative usefulness ($U_r$) as a function of the preference parameter $\mu$ sets the loss of misspecification, $L(\mu) = \mu \frac{FN}{FN+TP} + (1-\mu) \frac{FP}{FP+TN}$, in relation to the loss of a naïve decision rule, $\min(\mu, 1 - \mu)$, resulting from either always or never signaling a crisis depending on the preference parameter (Alessi and Detken, 2011):[5]

$$U_r(\mu) = 1 - \frac{L(\mu)}{\min(\mu, 1 - \mu)} \tag{1}$$

This implies a maximum relative usefulness of one, if $L(\mu) = 0$, and a usefulness of zero, if $L(\mu) = \min(\mu, 1 - \mu)$. A usefulness above (below) zero therefore means that the model is more (less) informative than the naïve decision rule. We use a standard choice of $\mu = 0.5$ for our baseline results, thus weighting the two types of errors equally.

In contrast to relative usefulness, the two other performance measures do not rely on an additional preference parameter. The Brier probability score (Brier, 1950; Diebold and Rudebusch, 1989; Knedlik and von Schweinitz, 2012) operates directly on probabilities instead of signals. It is simply given by the mean of the squared differences between predicted probabilities and actual outcomes (i.e. a special case of mean squared forecast error for binary dependent variables). By contrast, the area under the (receiver-operator characteristic) curve (AUC or AUROC) does operate on signals, but aggregates type-1 errors and type-2 errors over *all* possible classification thresholds $\tau$ (Janes, Longton, and Pepe, 2009; Drehmann and Juselius, 2014). The AUC can take on values between 0 and 1, with 0 being a misleading, 0.5 an uninformative and 1 a perfect set of forecasts.

As usual in forecasting, we are predominantly interested in the out-of-sample performance of different models. Thus, we split the panel dataset into two distinct parts: estimations are performed on an in-sample part (the *training sample*), while predictions and performance evaluations are derived on an out-of-sample part (the *test sample*). For comparability with previous findings, we follow Holopainen and Sarlin (2017) in setting our out-of-sample window to the period between 2005Q3 and 2016Q4. This leads to a good balance between observations available for estimation and for evaluating predictions, with approximately half of the pre-crisis observations contained in the in-sample part and half of the pre-crisis observations contained in the out-of-sample part. In most of the paper, we focus on recursive out-of-sample estimations where we predict the crisis probability quarter-by-quarter between 2005Q3 and 2016Q4 based on the information that was available in each respective quarter.[6] The performance measures are then based on the recursive predictions for the out-of-sample part of the dataset. If, instead, we are interested in in-sample performance, we use the same dataset for estimation and performance evaluation, i.e. we set the test and training sample equal to the full sample.

---

[5]An alternative usefulness function proposed by Sarlin (2013) was shown to be equivalent under a constant unconditional crisis probability (Sarlin and von Schweinitz, 2017).

[6]The definition of the early warning window is forward-looking. In order to account for that, all observations where $\bar{C}$ is yet unknown given information at time $t$ have to be excluded from the training sample. That is, for a forecast made in 2006Q1 we can only estimate the model on observations until 2003Q1 (unless a crisis occurs between 2003Q1 and 2006Q1, in which case the realization of $\bar{C}$ for some additional periods is known).

## 2.4 Bootstrap

Several of our estimation methods do not readily come with measures of estimation uncertainty. Moreover, even if such measures can be derived, they are conditional on very different (distributional) assumptions for different methods, making a comparison difficult. We solve this problem by bootstrapping, which provides a straightforward approach for calculating measures of estimation uncertainty under identical assumptions for all estimation methods. This allows us to test whether differences between model performances are statistically significant.

Bootstrapped measures of estimation uncertainty can be derived from the dispersion of estimates across random variations of the original dataset. These measures of estimation uncertainty are conditional on the statistical properties of the bootstrap datasets. Therefore, it is important to construct bootstrap datasets such that they preserve important statistical properties of the original dataset that are likely to affect the precision of estimates. In our case, autocorrelation and cross-sectional correlation are strong features of the data. Based on El-Shagi, Knedlik, and von Schweinitz (2013) and Holopainen and Sarlin (2017), we use a panel-block-bootstrap to account for these properties, as described in Appendix A.3 in more detail.

# 3 Data

## 3.1 Crisis Variable

We use the database for systemic banking crises established by the European System of Central Banks (ESCB) and the European Systemic Risk Board (ESRB) covering European countries from 1970 to 2016 (Lo Duca, Koban, Basten, Bengtsson, Klaus, Kusmierczyk, Lang, Detken, and Peltonen, 2017). This latest database refines previous crisis databases, both with respect to the identification of events and their timing. Crises are identified by the following two-step procedure. In a first step, "systemic financial stress events" are identified using the quantitative methodology of Duprey, Klaus, and Peltonen (2015). These financial stress events together with additional crises identified in previous databases (Laeven and Valencia, 2012; Babeckỳ, Havránek, Matějů, Rusnák, Šmídková, and Vašíček, 2014; Detken, Weeken, Alessi, Bonfim, Boucinha, Frontczak, Giordana, Giese, Jahn, Kakes, Klaus, Lang, Puzanova, and Welz, 2014) form a list of potential crisis events. In the second step, this list of potential crises is checked against a set of qualitative criteria defining systemic financial crises (see Lo Duca et al. (2017) for details).

Following Drehmann and Juselius (2014), we focus on systemic banking crises with at least partially domestic origins.[7] Furthermore, we expand the coverage of the crisis database to include two additional (non-European) advanced countries with important crisis experience, namely Japan and the United States.[8] As a result, our dataset covers

---

[7]Focusing on crises with at least partially domestic origins makes sense, as our modeling framework (where domestic variables determine the crisis probability of each country) does not allow for cross-country spillover effects. That is, we know a priori that these events are largely unforeseeable given the present modeling framework.

[8]For these countries, we use the crisis episodes identified by Laeven and Valencia (2012) adapting start

all of the "big five" crises identified by Reinhart and Rogoff (2008). The full list of crisis episodes used in our analysis after taking into account the availability of the explanatory variables may be found in Table B.1 in the Appendix. It includes 19 crises for European countries (of which 11 take place before 2008) as well as three crisis events in the United States and Japan (of which two take place before 2008). The majority of countries are included for a time period starting in the early to mid-1970s until the beginning of 2016.

As a robustness check we also run an estimation with crisis dates taken from the well-known Laeven and Valencia (2012) database. Their database has the advantage of being somewhat more agnostic in its definition of crisis events. Yet, the most recent European crisis database, which we use for our core results, provides a more comprehensive and more precise account of the crises in the European countries of our sample.

## 3.2  Potential Early Warning Indicators

Our collection of explanatory variables aims at capturing a wide range of sources of vulnerabilities for the banking sector. The channels we focus on are (i) asset prices, (ii) credit developments, (iii) the macroeconomic environment, as well as (iv) external and global imbalances. In the rest of this section, we take a closer look at these four channels before explaining transformations of the employed early warning indicators and our ultimate model specifications.

### 3.2.1  Sources of Vulnerabilities and Corresponding Indicators

**Asset prices:**  Historically, banking crises have often been preceded by asset price booms. Banking crises associated with house price booms and busts, could, for example, not only be observed during the global financial crisis of 2008, but also in a number of industrial countries in the late 1970s to early 1990s, such as in Spain, Sweden, Norway, Finland, and Japan (Reinhart and Rogoff, 2008, 2009). We therefore include *house prices* and *equity prices* to capture booms and busts in asset prices.

**Credit developments:**  High private sector indebtedness poses risks to the financial system when asset price booms are debt-financed, asset prices decrease and borrowers are unable to repay their debt (Kindleberger and Aliber, 2005; Jordà, Schularick, and Taylor, 2015). As a consequence of decreasing asset values, banks may be forced to deleverage, in particular when market liquidity is low and banks relying mainly on short-term funding face a liquidity mismatch (Brunnermeier and Oehmke, 2013; Brunnermeier, 2009). Deleveraging may induce a credit crunch and potentially lead to a recession. The effects of losses in asset values may be amplified by fire sales and may spill over to other assets as these are sold to meet regulatory and internal standards, such as capital and liquidity ratios. Moreover, bank runs may occur when the net worth of banks decreases and depositors lose confidence in the affected institutions (Allen and Gale, 2007). To capture risks related to high private sector indebtedness, we use *total credit to the private non-financial sector relative to GDP* as an indicator of how far credit developments are in line with real economic developments.

---

and end dates such that they are consistent with the definition in our core database.

**Macroeconomic environment:** Closely related to credit and asset prices are real economic developments. On the one hand, rapid economic growth may increase risk appetite, asset prices and credit growth (Drehmann, Borio, and Tsatsaronis, 2011; Kindleberger and Aliber, 2005; Minsky, 1982). On the other hand, real economic downturns may lead to repayment difficulties on the borrower side inducing asset price declines and financial sector difficulties (Allen and Gale, 2007). To capture real economic developments we include *GDP*, *gross fixed capital formation relative to GDP* and *inflation*. Furthermore, we include *three-month interbank interest rates*, as banks and investors may take on excessive risks when interest rates are low and, hence, low-risk assets are less attractive (Maddaloni and Peydró, 2011; Allen and Gale, 2007; Rajan, 2005). Conversely, an abrupt increase in interest rates may put pressure on banks as well (Minsky, 1982).

**External and global imbalances:** The external sector played a prominent role in the first seminal contributions to the early warning literature (Frankel and Rose, 1996; Kaminsky and Reinhart, 1999). These papers tended to focus more on balance-of-payment crises than on systemic banking crises. However, both types of crises may occur jointly and often reinforce each other as "twin crises" (Kaminsky and Reinhart, 1999). While classic balance-of-payment crises may be less of a concern for the countries considered in this paper, external imbalances may still add to vulnerabilities. Similarly to the reasoning on credit expansion and asset prices, large capital inflows from abroad may support asset price booms and induce a reversal in asset prices when these inflows decline or stop (Kaminsky and Reinhart, 1999; Calvo, 1998). Hence, we include the *real effective exchange rate* and the *current account balance relative to GDP*. Furthermore, global shocks may affect the domestic banking system through various channels of contagion, such as financial sector interconnectedness and trade links (Kaminsky and Reinhart, 2000). We therefore add *oil prices* as an indicator for global developments.

Any list of potential early warning indicators is naturally incomplete. Yet, for the purpose of comparing predictions across methods, this is not the key point, as long as the same variables are used across all methods. Furthermore, it turns out that data availability is a key issue. While several additional variables would have been plausible predictors on economic grounds, these variables are not available for a long enough time span and/or not available for all countries in our sample. In addition, lack of comparability across countries can be an issue for some variables. For instance, while both theoretical and empirical arguments for the inclusion of a debt service ratio variable can be made (e.g. Drehmann and Juselius, 2014; Drehmann, Juselius, and Korinek, 2017), the extent to which this variable would truncate the sample outweighs its potential benefit in our case.[9] Another important class of variables that we cannot include are those based on bank balance sheet data, where availability in the time series is even much more restricted than for the debt

---

[9]Proprietary debt service ratio data from the BIS is available starting at the earliest in 1980 (public debt service ratio data from 1999). We compared the availability of the debt service ratio by country with our sample of crises and early warning indicators described in table B.1. Including the debt service ratio starting from 1980 would exclude five of the 13 crises episodes prior to the financial crisis of 2007/2008 and one crisis episode starting in 2008 (as BIS debt service ratio data is not available for Ireland). In total, a substantial amount of around 400 out of 1801 total observations would be excluded from our dataset.

service ratio. Nevertheless, we will see in the results section that the variables we were able to include do have substantial explanatory power for predicting banking crises. Thus, for the purpose of comparing different prediction methods (as opposed, say, to the question of finding the most important early warning indicator(s) as, for instance, in Drehmann and Juselius, 2014), having a sufficient number of observations in the sample appears to outweigh the benefits of using a complete set of all potentially important early warning indicators. Indeed, a robustness check using a shorter sample length shows that reducing the amount of observations available for estimation substantially reduces out-of-sample prediction performance (see section 4.3).

### 3.2.2 Data Transformations

Several of our potential predictor variables naturally contain a time trend (the exception being inflation, money market rates and current account to GDP), which needs to be removed prior to estimation. Different approaches for filtering out the trend could be considered, and are an active area of research. Yet, our emphasis is on comparing machine learning and logit models on given identical datasets. We therefore focus on two of the most frequently employed approaches in the early warning literature: A Hodrick-Prescott (HP) filtering approach for our benchmark results and a growth rates approach for robustness.[10]

In the HP filtering approach, we transform early warning indicators into gaps by calculating deviations from the trend computed by a one-sided HP filter. Using a one-sided filter ensures that the information set at every point in time does not contain future information.[11] For variables such as total-credit-to-GDP ratio, real residential real estate prices and real equity prices we take into account recent evidence on lower frequency financial cycles, as documented, for instance, in Drehmann, Borio, and Tsatsaronis (2012), and Schüler, Hiebert, and Peltonen (2015). Thus, for these variables we take the value of $\lambda = 400'000$ often employed for early warning models (Drehmann and Juselius, 2014), corresponding to financial cycles being roughly four times as long as business cycles, which is broadly in line with the findings in the aforementioned literature on financial cycles. Moreover, this ensures that the total credit-to-GDP gap used in our analysis is in line with the definitions of the Basel Committee on Banking Supervision (BCBS) used for Basel III and for setting countercyclical capital buffers (Drehmann and Juselius, 2014; Basel Committee on Banking Supervision, 2010). For typical business cycle variables such as real GDP, gross fixed capital formation-to-GDP, and the real oil price we use a standard HP filter smoothing parameter of $\lambda = 1'600$. In the case of the real effective exchange rate, we also use $\lambda = 400'000$. The reason for this is that real effective exchange rate imbalances have been found to be extremely persistent, especially since the introduction of the Euro made adjustments via nominal exchange rate movements impossible (El-Shagi, Lindner, and von Schweinitz, 2016).[12]

---

[10] To remove extreme outliers, we furthermore winsorize the data at the 1%- and 99%-quantile.

[11] For the first $k$ observations in every country, we need to apply a two-sided filter instead of the recursive one-sided version, given that the filter needs a certain minimum number of observations to compute a trend. We set $k$ to ten years in the case of $\lambda = 400'000$ and to four years in the case of $\lambda = 1'600$.

[12] We also follow Drehmann and Juselius (2014) in calculating relative gaps (i.e. the deviations from trend normalized by the trend) for certain HP filtered variables, which can be useful to improve the comparability of gaps across time and countries. Relative gaps are used for real equity prices, real

Robustness checks are performed by transforming the variables into growth rates. In line with the reasoning for applying different HP filter smoothing parameters for capturing business cycles and financial cycles, we also use two different growth rate horizons. Our business cycle variables are transformed into four-quarter growth rates, while our financial cycle variables are transformed into 16-quarter growth rates.

A detailed description of all variables and their transformations may be found in Table B.2 in the Appendix. A list of non-transformed original data with the corresponding sources is documented in Table B.3. Summary statistics of the transformed (standardized) predictor variables are displayed in Table 2 for the HP filter and in Table B.4 for the growth rate transformation. When comparing the means of the selected indicators in the pre-crisis and tranquil periods, we note that the difference is particularly pronounced for the credit-to-GDP gap, the residential real estate price gap, the gross fixed capital formation-to-GDP gap and the current account balance relative to GDP. The first three of these indicators are, on average, higher and the current account balance to GDP is, on average, lower during pre-crisis periods. The volatility of these indicators is similar across pre-crisis and tranquil periods. Similar findings are obtained when using the growth rates transformation.

### 3.2.3 Specifications

For our model specifications, we use four different combinations of the described variables (also referred to as datasets). Dataset (iv) uses all of the available variables. In addition, we specify smaller models using subsets of variables in order to illustrate the relative performance of methods across datasets of varying complexity and information content. While reduced information content should generally reduce models' predictive ability, this may in some cases be offset by the gains from estimating less complex models. Datasets (i)-(iii) are mutually exclusive selections of indicators based on the different sources of vulnerabilities: (i) asset prices and credit developments, (ii) macroeconomic environment, and (iii) external and global imbalances. A list of the variables used in each dataset can be found in Table B.5. In order to guarantee comparability across the different datasets, we use the same sample for all datasets.

These a priori specified, economically motivated datasets have been chosen to allow for an economic interpretation of the information contained in each dataset. Moreover, by limiting ourselves to a priori specified, economically motivated variables and transformations, we seek to limit potential problems of data-mining. As Inoue and Kilian (2005) explain, when trying many variables and specifications, one is likely to find "spurious rejections of the no-predictability null and thus overfitting relative to the true model". Thus, avoiding data-mining is important for a realistic assessment of the true out-of-sample performance of early warning models. This is especially critical for early-warning models as policy-relevant analysis tools, as an overestimation of their accuracy might lead to a wrong sense of security.

In a similar vein, we have refrained from adding formal variable selection procedures to our analysis as another layer of complexity. Since different variable selection procedures could be applied, some of which specific to the different methods, we feel that adding such procedures might obscure the sources of differences in performance across methods.

---

residential real estate prices, the real oil price, real GDP and the real effective exchange rate.

Instead, our focus is on comparing the methods on identical sets of explanatory variables, which are selected a priori on economic grounds. The results we obtain, with logit outperforming machine learning methods on virtually *every* given dataset (see next section), appear to some extent orthogonal to the issue of variable selection. We therefore leave the issue of formal variable selection in the context of method comparison for future research.

Table 2: Summary statistics: Gap dataset

| Variable name | Pre-crisis periods | | | | | Tranquil periods | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | mean | sd | min | max | obs | mean | sd | min | max | obs |
| Total credit-to-GDP gap | 0.76 | 1.03 | -1.16 | 3.19 | 171 | -0.08 | 0.96 | -3.06 | 3.19 | 1608 |
| Real residential real estate price gap | 0.58 | 1.06 | -1.36 | 2.48 | 171 | -0.06 | 0.97 | -2.58 | 2.48 | 1608 |
| Real equity price gap | 0.23 | 0.74 | -1.73 | 1.68 | 171 | -0.02 | 1.02 | -1.96 | 3.61 | 1608 |
| Real GDP gap | 0.10 | 0.94 | -3.15 | 2.20 | 171 | -0.01 | 1.01 | -4.81 | 3.71 | 1608 |
| Inflation rate | -0.08 | 0.87 | -1.24 | 3.64 | 171 | 0.01 | 1.01 | -1.38 | 4.83 | 1608 |
| Gross fixed capital formation-to-GDP gap | 0.32 | 1.02 | -2.46 | 2.97 | 171 | -0.03 | 0.99 | -4.63 | 3.19 | 1608 |
| Real 3-month money market rate | 0.06 | 1.10 | -3.29 | 2.32 | 171 | -0.01 | 0.99 | -3.29 | 2.61 | 1608 |
| Current account-to-GDP ratio | -0.73 | 1.08 | -3.50 | 2.17 | 171 | 0.08 | 0.96 | -3.50 | 2.75 | 1608 |
| Real effective exchange rate gap | 0.20 | 0.85 | -2.08 | 2.72 | 171 | -0.02 | 1.01 | -2.61 | 2.72 | 1608 |
| Real oil price gap | 0.26 | 1.00 | -2.56 | 2.34 | 171 | -0.03 | 1.00 | -2.56 | 3.13 | 1608 |

*Note*: Data have been standardized using the unconditional mean and standard deviation across all periods. Since data are winsorized at the 1% and 99% level, minimum (maximum) values may be the same in pre-crisis and tranquil periods.

# 4 Results

## 4.1 In-sample Predictive Performance across Methods

Table 3 reports the *in-sample* relative usefulness of the five different methods (in rows) for four different sets of explanatory variables (in columns). The best performance on each dataset is indicated in bold. Significance stars (obtained from our bootstrap procedure) indicate whether the respective usefulness is significantly below that of the best-performing method on the same dataset.

In line with the literature (Alessi and Detken, 2018; Holopainen and Sarlin, 2017; Tanaka et al., 2016), machine learning methods such as knn and random forest always attain substantially higher in-sample relative usefulness than the corresponding logit model. Random forest achieves the highest in-sample relative usefulness on dataset 1 (credit/asset prices), dataset 3 (external imbalances) and 4 (all variables), while knn has the best in-sample performance on dataset 2 (macroeconomic environment).[13] The inferiority of the logit model's in-sample performance relative to the best model on every dataset is significant. As Table C.1 shows, these findings are robust to using alternative measures of prediction performance. Moreover, the fit of knn and random forest is often close to perfect, that is a value close to one for relative usefulness and area under the curve (AUC) is achieved.[14]

Table 3: *In-sample* relative usefulness

|  | (1) Credit/Asset Prices | (2) Macro | (3) External | (4) All |
|---|---|---|---|---|
| logit | 0.347*** | 0.202*** | 0.390*** | 0.511*** |
|  | [0.088,0.457] | [0.114,0.306] | [0.271,0.461] | [0.294,0.655] |
| trees | 0.432*** | 0.000*** | 0.000*** | 0.674*** |
|  | [0.172,0.576] | [-0.099,0.135] | [-0.207,0.147] | [0.406,0.883] |
| knn | 0.693*** | **0.965** | 0.685*** | 0.955 |
|  | [0.419,0.825] | [0.693,1.000] | [0.470,0.832] | [0.660,1.000] |
| rf | **0.959** | 0.956 | **1.000** | **0.990** |
|  | [0.609,1.000] | [0.683,1.000] | [0.702,1.000] | [0.672,1.000] |
| svm | 0.327*** | 0.512*** | 0.230*** | 0.930 |
|  | [0.080,0.460] | [0.340,0.676] | [0.015,0.428] | [0.566,1.000] |

*Note*: Highest usefulness on each dataset in bold. Stars indicate whether the respective usefulness is significantly below the best performance on the same dataset (***/**/* for the 1%/5%/10% level). Numbers in brackets indicate 90% confidence bands.

---

[13]See Table B.5 for a list of the four sets of explanatory variables (also referred to as datasets).

[14]Table C.1 also shows what is behind the zero usefulness of trees on dataset 2 and 3, which is due to a degenerate tree that always predicts a crisis. This happens because the available variables are not very informative relative to the required tree complexity. The optimal tree (according to the tree's internal cost function) then consists of the unconditional forecast. By definition, relative usefulness and AUC will be zero and 0.5, respectively.

## 4.2 Out-of-sample Predictive Performance across Methods

In line with the standard in the empirical literature, the focus of our evaluation is on recursive out-of-sample performance rather than in-sample performance. For every point in time from 2005Q3 until the end of our sample in 2016Q4, we estimate the model recursively, strictly using only data until that time.[15] We thus obtain predictions for approximately 300 out-of-sample observations, including 60 pre-crisis periods. These predictions are then used to calculate out-of-sample performance measures.

Table 4: *Out-of-sample* relative usefulness

|       | (1) Credit/Asset Prices | (2) Macro | (3) External | (4) All |
|-------|-------------------------|-----------|--------------|---------|
| logit | **0.368** | -0.236* | **0.438** | **0.605** |
|       | [0.237,0.487] | [-0.394,-0.074] | [0.34,0.532] | [0.481,0.727] |
| trees | 0.201** | -0.605*** | 0.390*** | 0.126*** |
|       | [0.084,0.329] | [-0.703,-0.499] | [0.27,0.516] | [0.007,0.248] |
| knn   | 0.247* | -0.087 | 0.293** | -0.062*** |
|       | [0.123,0.374] | [-0.147,-0.023] | [0.182,0.399] | [-0.142,0.024] |
| rf    | 0.246* | -0.274*** | 0.117*** | -0.003*** |
|       | [0.134,0.358] | [-0.343,-0.194] | [0.036,0.203] | [-0.126,0.124] |
| svm   | 0.243* | **-0.082** | 0.060*** | -0.186*** |
|       | [0.118,0.355] | [-0.198,0.039] | [-0.11,0.215] | [-0.305,-0.063] |

*Note*: Highest usefulness on each dataset in bold. Stars indicate if the respective usefulness is significantly below the best performance on the same dataset (***/**/* for the 1%/5%/10% level). Numbers in brackets indicate 90% confidence bands.

Table 4 shows the out-of-sample relative usefulness of five different methods (in rows) for four different sets of explanatory variables (in columns). The table shows that the logit model almost always outperforms the machine learning methods. The only exception is the dataset based on macroeconomic variables. However, in that case the relative usefulness is negative for all methods. That is, a naïve forecast would be better than using any of the considered models. Table C.2 shows that the superiority of the logit model is confirmed by the other two performance measures, namely the area under the curve (AUC) and Brier probability score (BPS) - again, with two exceptions on the macroeconomic dataset for AUC. Performance of all machine learning methods on dataset 1 (credit and asset prices), dataset 3 (external variables) and dataset 4 (all variables) is significantly worse than that of the logit model using any of the three performance measures.

## 4.3 Robustness

In this section, we assess the degree to which our results are robust to four key variations of the modeling setup. These variations concern the choice of preference parameter, data transformation, sample length, and the crisis database.

---

[15]As explained in the methodology section, we follow (Holopainen and Sarlin, 2017) in choosing 2005Q3 as starting point for the recursive out-of-sample evaluation.

**Preference parameter:** First of all, we make sure that our results do not hinge on the choice of loss function preference parameter (see equation (1)). While AUC and BPS are preference-independent measures, the preference parameter enters into the computation of relative usefulness via the loss function, which is used to evaluate forecasts and to compute optimal thresholds. Our benchmark value for the preference parameter, $\mu = 0.5$, represents a balanced trade-off between type-1 (missed crises) and type-2 errors (false alarms). It is easily conceivable that missing a crisis may be more costly than issuing a false alarm. In this case, more weight should be given to type-1 errors. Therefore, we conduct a robustness check for a preference parameter of $\mu = 0.6$ which assigns slightly more weight to type-1 errors. We re-estimate all models including hyperparameters given this new preference parameter and report results in Table C.3 in the Appendix.[16] We find that relative usefulness drops for all methods. However, while the logit model still has a substantially positive relative usefulness on all datasets, machine learning methods seem to deteriorate more strongly, for instance on datasets (4) and (1). As a result, logit continues to outperform machine learning methods on all datasets except dataset (2), and the logit with all variables remains the "best" overall model.

**Data transformation:** As a second robustness test, we check the extent to which our results hinge on the choice of using HP filter gaps for removing the trend in our explanatory variables. To this end, we replace the HP filter gaps by simple growth rates, where lags used to compute growth rates differentiate between business cycle and financial cycle variables (for details see Section 3.2.2). The results of this robustness check are shown in Table C.4. Looking at relative usefulness, we see that some of the machine learning models perform better when using growth rates rather than HP filter gaps. By contrast, all four logit specifications have somewhat lower relative usefulness and higher BPS than before. Knn actually outperforms logit in terms of relative usefulness and BPS on dataset 1, but not on datasets 3 and 4. Across all models considered, the logit model with all variables (logit.4) remains the best model according to relative usefulness and BPS.

**Sample length:** A third important test concerns the robustness of our results relative to variations of our sample. Specifically, we consider a robustness check, where we cut off the first ten years of our dataset, amounting to approximately one-sixth of our observations. Financial repression during the 1970s may have affected the behavior of our explanatory variables and their impact on the probability of future financial crises. More generally, the underlying data-generating process may be time-varying, suggesting a trade-off between sample length and sample homogeneity.

Table C.5 presents results when restricting our sample to exclude all observations prior to 1980Q1. For the logit model, we find that out-of-sample prediction performance based on this smaller set of information is lower than when using the full dataset. Thus, the trade-off between sample length and sample homogeneity appears to be tilted in favor of sample length. The relative ordering of machine learning methods compared with the logit method is unchanged. Thus, we conclude that, while sample length appears to be important, our main finding regarding relative prediction performance between methods is strikingly robust to the change in the sample. We also note that this robustness appears

---

[16]Optimal hyperparameters change only slightly for some models, such that the effect of the preference parameter is primarily driven by the change in the loss function.

to be driven by the robustness of the logit method, while machine learning methods sometimes react quite strongly to this moderate change in the sample.

**Crisis database:**  In our fourth robustness check we replace the ESCB/ESRB crisis database by the well-known database of Laeven and Valencia (2012). Results are shown in table C.6. It turns out that changing the dependent variable of our models induces the biggest changes in the results. Looking at relative usefulness (and AUC), there are winners and losers across all datasets and methods. However, results for logit.4 and logit.3 remain strikingly robust to this change. As a consequence, they continue to outperform their machine learning competitors and logit.4 remains the best overall model according to relative usefulness. BPS performances are more mixed. While logit continues to have the lowest BPS across models on dataset 4 and 3, the overall lowest BPS occur for all methods on dataset 1. However, when also taking into account relative usefulness and AUC, the logit.4 model still seems clearly preferable to the models on dataset 1.

   As a general pattern, out-of-sample performance is weaker for the majority of models when using Laeven & Valencia crises. This effect is particularly pronounced for BPS which is markedly higher than for our benchmark crisis database. The reason for this lies in two features of the Laeven & Valencia database. First, the share of pre-crisis periods in the training sample (prior to 2005Q3) is comparably low. This leads to low predicted crisis probabilities at the beginning of the out-of-sample exercise. Second, the out-of-sample period ends in 2012 and is dominated by the great financial crisis, leading to a share of pre-crisis periods of 68%. As a consequence, the estimated unconditional probability of pre-crisis periods is (too) low until the great financial crisis simultaneously hits the majority of countries in the sample. This constellation leads to high BPS values and often weaker performance in terms of relative usefulness and AUC, with logit.4 and logit.3 being comparably robust.

   Overall, our robustness checks confirm the finding that a logit model using all variables offers the best predictive performance among all models considered. Moreover, we saw that changes to model specifications, as those considered in this section, can induce substantial changes to some models' performance. Given this, we see the robustness of the logit.4 (and logit.3) model across specifications, as an additional feature of these models. By contrast, performance of machine learning methods is substantially less robust.

## 4.4   Interpretation and Discussion

A comparison of the performance of in-sample and recursive out-of-sample estimations gives an indication as to why machine learning methods do not outperform the logit approach in this application. Figure 1 displays the relationship between in-sample and recursive out-of-sample performance across models based on our benchmark results shown in Tables 3 and 4. A striking result is that many of the machine learning models (including all random forest models) achieve a near-perfect in-sample fit (relative usefulness close to its theoretical maximum of 1), but, at the same time, show much lower out-of-sample performance. This suggests that overfit may be a major issue for at least some of the models. Moreover, even for those machine learning models where in-sample fit is not perfect, their out-of-sample performance is in most cases markedly below their in-sample

performance. By contrast, among logit models this is only the case for logit.2, which we saw is a special case of negative relative usefulness for all methods on this dataset. Figures C.1 through C.4 in the appendix show that this pattern holds true more generally across all robustness checks.

Figure 1: Relative usefulness of in- and out-of-sample estimation by model.



In addition to the empirical evidence in figure 1, a theoretical argument pointing to overfit (relative to the true model) can be made. As a thought experiment, suppose we knew the true data generating process (DGP) and we had a model (and data) at hand that would give us for each observation the true conditional probability of a crisis. Even in this case, the prediction error as measured by $U_r$, AUC or BPS would still be positive (except in the degenerate case where conditional crisis probabilities could only be either 1 or zero). For that reason, even with a perfect model, we would not expect relative usefulness to approach its maximum value of 1, but rather to converge (both in-sample and out-of-sample) to a DGP specific maximum between 0 and 1 as the sample size increases.[17] This can formally be seen in Table 5 of Boissay, Collard, and Smets (2016). They present a DSGE model generating credit boom driven crises, and run an early warning exercise on simulated data from their model. It turns out that crisis prediction using the true model-implied conditional probability still leads to considerable error rates (around 1/3 missed crises). Their results also show that a logit model estimated on binary crisis realizations is able to converge to a performance similar to that of the true model. The remaining error rates under the true model reflect the fact that crises can only be predicted in probability and not with certainty. In other words, crises are driven by

---

[17]As a simple illustration, suppose for example that in 50% of the cases the true crisis probability was 80%, while in 50% of the cases, it was 20%, and that our signaling threshold was 50%. Then, even when knowing the true model, we would still have a false positive rate of 20% and a false negative rate of 20%, leading to a relative usefulness of only 60% (assuming $\mu = 0.5$).

a predictable component, captured by the true model, and a substantial unpredictable component (given the observables), which cannot be forecasted by any model.

To the extent that a logit model is already able to closely approximate the true model (as in Boissay et al. (2016)), it will not be possible to substantially outperform this model. This implies that when machine learning methods fit the data beyond (or below) the predictable component, this comes at the cost of worse performance in the recursive out-of-sample estimation. This may be an issue even for those machine learning models where in-sample fit is not perfect. Theoretically, the hyperparameters of the machine learning methods should provide some safeguard against overfit. However, despite devoting considerable effort to the calibration of these hyperparameters via a sophisticated cross-validation procedure (see Appendix A.4), the overfit still persists for many of the considered models. In sum, it appears that logit models naturally limit the amount of overfit, while being sufficiently flexible in their approximation of the data generating process.

The conclusion from our out-of-sample forecast comparison is different than that of Alessi and Detken (2018) and Tanaka et al. (2016), who argue that a random forest has a better prediction performance than a logit model in an early warning setting. However, Alessi and Detken (2018) do not run an out-of-sample comparison of the two methods. Their argument is rather based on results from k-fold-cross-validation, where they find some differences between the AUC of one random forest specification (AUC = 0.94) and two logit specifications (AUC = 0.84, and 0.93 respectively). Setting aside the question of whether this difference is statistically significant, the high levels of AUC (close to the maximum of 1) suggest that the cross-validation procedure may provide an inflated estimate of the performances of these methods. In fact, cross-validation estimates often appear to be closer to in-sample performance than to out-of-sample performance. The tendency of cross-validation to provide inflated estimates of performance, particularly in the presence of cross-sectional and serial correlation, has also been recognized by Holopainen and Sarlin (2017). As a consequence, these estimates are likely to be biased towards (more complex) machine learning methods, given their above-mentioned tendency to overfit in-sample data. Another important point has been noted by Neunhoeffer and Sternberg (2018) based on an example of civil war prediction from the political science literature. They show that performance of machine learning methods has been seriously over-estimated, in studies using cross-validation for both hyperparameter selection and model evaluation. Tanaka et al. (2016) find somewhat more pronounced differences between logit and random forest performance for a bank-level early warning model. However, they also focus on cross-validation estimates of performance.

The importance of conducting model comparisons via out-of-sample experiments has also been advocated by Holopainen and Sarlin (2017). However, in their out-of-sample forecasting exercise, logit models are outperformed by machine learning methods (except for trees). We conjecture that differences in the employed training sample are a key driver of the contrasting results. In their application, more than 95% of the pre-crisis periods are located in the recursive out-of-sample period and are thus unavailable for the first recursive estimations. This means that their model estimations are driven by a few influential pre-crisis observations in their training sample, such that the selection of these observations is a critical determinant of their out-of-sample results. By contrast, our broader data basis enables us to use nearly 60% of all pre-crisis periods in the first

recursive estimation, mimicking more closely actual out-of-sample prediction tasks. The different sample appears to be the most important explanation for our differing results.

As a final word regarding interpretation, we want to make clear that while we think it is important to establish a robust and valid out-of-sample prediction exercise, we do not want to over-interpret its results. After all, our results hold for the given finite dataset at hand. In particular, our out-of-sample window is naturally dominated by the great financial crisis of 2007-2008. Moreover, we cannot fully exclude the possibility that some other (potentially more sophisticated) modeling approach is able to outperform the logit model, or that machine learning methods may be preferable on other (possibly larger) datasets. However, we think that we have established that the logit model is surprisingly hard to beat, in line with findings in the forecasting literature more generally, that simple forecasting models often outperform more complex models. We have provided theoretical and empirical arguments, as well as a discussion of the literature, suggesting systematic issues related to overfit driving this result. Further research is needed, to gain a more complete understanding of the conditions under which machine learning methods can be successfully applied, in general, and in particular to early warning models of financial crises.

## 4.5   Economic Interpretation of the Logit Model

Besides the statistical results presented up to here, we also want to provide an economic interpretation of our early warning model. We do this for the example of the logit.4 baseline model. This model had the best performance in our horse race, and encompasses all other models in terms of variables used. Table 5 contains coefficients of the in-sample estimation. The most important variables in our model are the current account balance, credit-to-GDP, residential real-estate prices and gross fixed capital formation. This is consistent with existing theoretical and empirical evidence, which documents the vulnerability of the banking sector to many different channels, as described in Section 3. We find that the vast majority of coefficients have the expected sign. Credit growth above long-run trend increases the probability of being in an early warning window, as do high residential real estate prices and high equity prices. Thus, debt-financed asset price booms are found to be major drivers of crises (Kindleberger and Aliber, 2005; Jordà et al., 2015). The positive coefficient on the gap of gross fixed capital formation can be interpreted in a similar way: high levels of investment in fixed capital may be driven by overly optimistic expectations, leading to problems when future returns are lower than expected. Economic downturns, indicated by lower growth and lower inflation rates, also increase the crisis probability (albeit insignificantly). Last but not least, current account deficits and overvaluation of the real effective exchange rate are signs of an uncompetitive and (external) debt-financed economy.

Table 5 also allows us to look at the average marginal effects, which are of particular interest for the significant variables. Their average marginal effects, approximating the effect of a one standard deviation change in the respective indicator on the predicted crisis probability, are 1.4 percentage points for equity price gap, 3.0 percentage points for residential real estate price gap, 3.6 percentage points for gross-fixed capital formation-to-GDP gap, 3.8 percentage points for credit-to-GDP gap, and -4.5 percentage points for the current account balance. Compared to the unconditional probability of being in an early

Table 5: Logit coefficients (full sample)

|  | Coefficient | Std. Error | Marg. Effects | Stdev. recursive |
|---|---|---|---|---|
| Constant | -2.771*** | 0.113 | -0.199 | 0.151 |
| Total credit-to-GDP gap | 0.531*** | 0.102 | 0.038 | 0.033 |
| Real residential real estate price gap | 0.425*** | 0.096 | 0.030 | 0.073 |
| Real equity price gap | 0.200* | 0.108 | 0.014 | 0.029 |
| Real GDP gap | -0.141 | 0.123 | -0.010 | 0.049 |
| Inflation rate | -0.104 | 0.119 | -0.007 | 0.102 |
| Gross fixed capital formation-to-GDP gap | 0.499*** | 0.106 | 0.036 | 0.018 |
| Real 3-month money market rate | 0.007 | 0.100 | 0.001 | 0.091 |
| Current account-to-GDP ratio | -0.63*** | 0.093 | -0.045 | 0.077 |
| Real effective exchange rate gap | 0.015 | 0.101 | 0.001 | 0.031 |
| Real oil price gap | 0.002 | 0.090 | 0.000 | 0.031 |

*Note:* This table reports coefficients, standard errors and marginal effects for the logit model using all variables, estimated on all available observations. The model is estimated with standardized data to make coefficients comparable. The last column shows the standard deviation of coefficient estimates across recursive estimations.

warning window, which is just above 9.5%, these effects are substantial. In comparison to that, the marginal effects of the insignificant variables are mostly negligible. Even though the effects of the significant variables are sizeable, it has to be noted that the model never implies a probability above 90% of being in an early warning window. For such a probability, all important variables need to be around two standard deviations away from their mean at the same time, which is an extremely rare event. This is in line with the view that, while our observables may signal the buildup of vulnerabilities in probability, there remains a substantial unpredictable component driving the ultimate realization or non-realization of crises.

In addition to their economic and statistical significance, coefficients are quite stable over time. In our recursive out-of-sample forecasting exercise, we can observe how (re-)estimated coefficients change across time, as more and more information becomes available. To summarize this, the last column of Table 5 (Stdev. recursive) reports the standard deviation of coefficients across recursive estimations. As we can see, the magnitude of changes in coefficients during the out-of-sample window is relatively small for the significant coefficients, which is even more remarkable given the occurrence of the great financial crisis during this time period. This suggests a degree of robustness of the estimated model with respect to the addition of new information, which is promising regarding the potential use of such models for future (true) out-of-sample predictions.

## 4.6   Probability Predictions with the Logit Model

Given its economic content and strong performance relative to machine learning models, it is worthwhile to take a closer look at the performance of the logit model in different scenarios. Table 6 shows the performance of the logit.4 model (using all variables) for different estimation setups. In the first part of Table 6, we report performance measures

for four different estimation strategies. In the first row, we estimate and evaluate the model on the full sample (in-sample estimation *is.full*). The second row contains results for the recursive out-of-sample estimation on which the previous model comparisons were based (*oos.full*). As an alternative to recursive estimation, we also explore one-off splits, where an estimation on observations until 2005Q2 is used to predict probabilities after that (*forecast.full*) or the other way around (*backcast.full*) in rows three and four.[18]

In general, we find that the performance of the model is quite stable independent of the specification. For the in-sample estimation, recursive forecasting and forecasting in a one-off split, all three performance measures (relative usefulness, AUC and Brier probability score) are quite similar. This is another piece of evidence that the logit model provides a stable approximation of the true data generating process. The one-off split for the backcasting exercise, where we use only information from the great financial crisis in an early warning model on all data prior to 2005Q2 seems to go somewhat against that result. However, we should be clear that we are predicting more than 83% of the dataset based on the remaining 17% of observations. This is a very hard task, especially since the pre-crisis probability in the training sample is considerably larger than in the rest of the sample, as also indicated by the strong increase of the threshold. Yet, even in this extreme scenario relative usefulness remains positive.

Table 6: Performance of logit model in different setups

|  | Threshold | TP | FP | TN | FN | FPrate | FNrate | $U_r$ | AUC | BPS |
|---|---|---|---|---|---|---|---|---|---|---|
| is.full | 0.109 | 119 | 297 | 1312 | 52 | 0.185 | 0.304 | 0.511 | 0.810 | 0.073 |
| oos.full | 0.091 | 45 | 35 | 206 | 15 | 0.145 | 0.250 | 0.605 | 0.852 | 0.125 |
| forecast.full | 0.081 | 47 | 22 | 219 | 13 | 0.091 | 0.217 | 0.692 | 0.881 | 0.132 |
| backcast.full | 0.392 | 72 | 514 | 854 | 39 | 0.376 | 0.351 | 0.273 | 0.681 | 0.237 |
| is.sig | 0.109 | 72 | 118 | 827 | 26 | 0.125 | 0.265 | 0.610 | 0.842 | 0.064 |
| oos.sig | 0.090 | 33 | 5 | 95 | 6 | 0.050 | 0.154 | 0.796 | 0.893 | 0.142 |
| forecast.sig | 0.081 | 29 | 3 | 111 | 6 | 0.026 | 0.171 | 0.802 | 0.925 | 0.129 |
| backcast.sig | 0.392 | 21 | 40 | 134 | 10 | 0.230 | 0.323 | 0.448 | 0.752 | 0.231 |

*Note*: The table shows the performance of the logit model using all variables when using different strategies for estimating and evaluating the model. See text for detailed explanation.

Looking at the absolute performance of the model both in-sample and out-of-sample, where it correctly classifies almost 70% and 75% of pre-crisis periods respectively (1 - FNrate) while keeping the rate of false alarms below 20%, we think that this model may add value to a quantitative assessment of financial stability risks. That being said, the estimation of the logit models is still accompanied by considerable uncertainty. This also leads at times to large confidence bands around point estimates of crisis probabilities. It is plausible that policymakers might differentiate between a signal which is based on a probability that is significantly different from the signaling threshold, and a signal based on a probability that is insignificantly above or below the threshold. In this case, one could

---

[18]That is, *backcast.full* assumes knowledge of observations from 2005Q3 until 2016Q2 for estimating the model, and predicts all previous observations based on these estimates. The performance measures for that model may serve as another check of model stability.

divide policy recommendations from early warning models into three categories: (a) a clear recommendation to act based on a crisis probability significantly above the signaling threshold; (b) a clear recommendation to abstain from acting based on a crisis probability significantly below the signaling threshold; and (c) a recommendation to start further analyses and investigations if the crisis probability is insignificantly different from the threshold. In the second part of the table, we incorporate this line of thought and report performance measures for the four estimation strategies based on only those observations where the estimated probability is significantly different from the threshold.[19]

In the backcasting case, only 14% of signals are significant, while 46%-50% are significant in the two forecasting cases and nearly 60% in the in-sample estimation. Zooming into the details, we can see that the share of significant periods differs across the different classes in the confusion matrix. False negatives and, especially, false positives are much less often statistically significant than true positives and negatives. This is reassuring because it implies that a focus on significant signals increases relative usefulness. In fact, focusing on significant signals would lead to only five false positives and six false negatives in the recursive out-of-sample forecasting exercise (oos.sig), while still correctly signaling a substantial share of pre-crisis and tranquil periods. Thus, while such an approach would not be able to issue signals before every crisis, it would allow signaling before a substantial number of crisis events with a high degree of confidence. In this sense, taking uncertainty explicitly into account when making predictions may help to increase the accuracy of the early warning model.

# 5    Conclusion

This paper has presented an analysis of early warning models for systemic banking crises, based on a dataset covering 15 advanced countries over the period 1970-2016. It is one of the first papers to use the latest version of the ESCB/ESRB crisis database (Lo Duca et al., 2017). This database is extended using the Laeven & Valencia database in order to cover not only European countries, but also the U.S. and Japan. Regarding potential predictor variables, we build on the existing empirical and theoretical literature and include indicators representing several important channels for the buildup of systemic banking crises.

Our analysis confirms that indicators representing credit and asset prices, but also those related to external imbalances carry information which can be used to predict the likelihood of systemic banking crises. A logit model taking into account all of these channels would have been able to issue relatively accurate warnings before the financial crisis of 2007/2008 for many countries. At the same time, we also emphasize that there always remains an unpredictable component of systemic banking crises, which leads to classification errors in ex post measures of prediction performance.

We assess how different methods – a benchmark logit approach and several machine learning methods – perform in a quasi real-time (pseudo) out-of-sample forecasting experiment. It turns out that the logit approach is surprisingly hard to beat, generally leading

---

[19]We use 90% bootstrap confidence bands around estimated probabilities. The threshold is always based on all observations, without taking estimation uncertainty into account. This is of course only one option for incorporating estimation uncertainty, serving as an illustration. We leave the many potential variations and refinements for future research.

to lower out-of-sample prediction errors than the machine learning methods. This result holds under different performance measures and different selections of variables, and is robust to alternative choices of crisis variable, variable transformation, sample length or loss function preference parameter.

Our interpretation of this result is that a strong in-sample fit of machine learning methods should not necessarily be taken as an indication of strong out-of-sample prediction performance, since it could alternatively be a sign of overfitting. Moreover, the stability of these methods' performances to different variations of the setup seems to be less pronounced than that of the logit model. This suggests that performance of machine learning methods in real-world out-of-sample prediction situations cannot be taken for granted. Instead, the circumstances under which these methods offer clear advantages as well as potential modifications for improving their stability and performance in early warning applications need further investigation.

# References

Alessi, L. and C. Detken (2011). Quasi Real Time Early Warning Indicators for Costly Asset Price Boom/Bust Cycles: A Role for Global Liquidity. *European Journal of Political Economy 27*(3), 520–533.

Alessi, L. and C. Detken (2018). Identifying Excessive Credit Growth and Leverage. *Journal of Financial Stability 35*, 215–225.

Allen, F. and D. Gale (2007). *Understanding Financial Crises.* Oxford University Press, Oxford.

Arlot, S., A. Celisse, et al. (2010). A Survey of Cross-Validation Procedures for Model Selection. *Statistics Surveys 4*, 40–79.

Babeckỳ, J., T. Havránek, J. Matějů, M. Rusnák, K. Šmídková, and B. Vašíček (2014). Banking, Debt, and Currency Crises in Developed Countries: Stylized Facts and Early Warning Indicators. *Journal of Financial Stability 15*, 1–17.

Basel Committee on Banking Supervision (2010). Guidance for National Authorities Operating the Countercyclical Capital Buffer. Technical report.

Berg, A. and C. Pattillo (1999). What Caused the Asian Crises: An Early Warning System Approach. *Economic Notes 28*(3), 285–334.

Boissay, F., F. Collard, and F. Smets (2016). Booms and Banking Crises. *Journal of Political Economy 124*(2), 489–538.

Breiman, L. (1996). Bagging Predictors. *Machine Learning 24*(2), 123–140.

Breiman, L. (2001). Random Forests. *Machine Learning 45*(1), 5–32.

Breiman, L., J. Friedman, R. Olshen, and C. Stone (1984). *Classification and Regression Trees.* Wadsworth.

Brier, G. W. (1950). Verification of Forecasts Expressed in Terms of Probability. *Monthly weather review 78*(1), 1–3.

Brunnermeier, M. K. (2009). Deciphering the Liquidity and Credit Crunch 2007-2008. *Journal of Economic Perspectives 23*(1), 77–100.

Brunnermeier, M. K. and M. Oehmke (2013). Bubbles, Financial Crises, and Systemic Risk. In *Handbook of the Economics of Finance*, Chapter 18, pp. 1221–1288.

Bussière, M. and M. Fratzscher (2006). Towards a New Early Warning System of Financial Crises. *Journal of International Money and Finance 25(6)*, 953–973.

Calvo, G. A. (1998). Capital Flows and Capital-Market Crises: The Simple Economics of Sudden Stops. *Journal of Applied Economics 1*(1), 35–54.

Cerutti, E., S. Claessens, and L. Laeven (2017). The Use and Effectiveness of Macroprudential Policies: New Evidence. *Journal of Financial Stability 28*, 203–224.

Claessens, S. (2015). An Overview of Macroprudential Policy Tools. *Annual Review of Financial Economics 7*, 397–422.

Cover, T. and P. Hart (1967). Nearest Neighbor Pattern Classification. *IEEE Transactions on Information Theory 13*(1), 21–27.

DeLong, E. R., D. M. DeLong, and D. L. Clarke-Pearson (1988). Comparing the Areas under Two or More Correlated Receiver Operating Characteristic Curves: a Nonparametric Approach. *Biometrics*, 837–845.

Detken, C., O. Weeken, L. Alessi, D. Bonfim, M. M. Boucinha, S. Frontczak, G. Giordana, J. Giese, N. Jahn, J. Kakes, B. Klaus, J. H. Lang, N. Puzanova, and P. Welz (2014). Operationalising the Countercyclical Capital Buffer: Indicator Selection, Threshold Identification and Calibration Options. ESRB Occasional Paper 5.

Diebold, F. X. and G. D. Rudebusch (1989). Scoring the Leading Indicators. *Journal of Business 62*(3), 369–391.

Drehmann, M., C. Borio, and K. Tsatsaronis (2011). Anchoring Countercyclical Capital Buffers: the Role of Credit Aggregates. *International Journal of Central Banking 7*(4), 189–240.

Drehmann, M., C. Borio, and K. Tsatsaronis (2012). Characterising the Financial Cycle: Don't Lose Sight of the Medium Term! BIS Working Papers 380.

Drehmann, M. and M. Juselius (2014). Evaluating Early Warning Indicators of Banking Crises: Satisfying Policy Requirements. *International Journal of Forecasting 30*(3), 759–780.

Drehmann, M., M. Juselius, and A. Korinek (2017, June). Accounting for Debt Service: The Painful Legacy of Credit Booms. Technical Report 645, BIS Working Papers.

Duprey, T., B. Klaus, and T. Peltonen (2015). Dating Systemic Financial Stress Episodes in the EU Countries. ECB Working Paper 1873.

El-Shagi, M., T. Knedlik, and G. von Schweinitz (2013). Predicting Financial Crises: The (Statistical) Significance of the Signals Approach. *Journal of International Money and Finance 35*, 76–103.

El-Shagi, M., A. Lindner, and G. von Schweinitz (2016). Real Effective Exchange Rate Misalignment in the Euro Area: A Counterfactual Analysis. *Review of International Economics 24*(1), 37–66.

European Central Bank (2010). Financial Stability Review, June 2010.

European Central Bank (2017). Financial Stability Review, May 2017.

Frankel, J. A. and A. K. Rose (1996). Currency Crashes in Emerging Markets: An Empirical Treatment. *Journal of International Economics 41*(3), 351–366.

Gourinchas, P.-O. and M. Obstfeld (2012). Stories of the Twentieth Century for the Twenty-First. *American Economic Journal: Macroeconomics 4*(1), 226–265.

Holopainen, M. and P. Sarlin (2017). Toward Robust Early-Warning Models: A Horse Race, Ensembles and Model Uncertainty. *Quantitative Finance 17*(12), 1–31.

Inoue, A. and L. Kilian (2005). In-Sample or Out-of-Sample Tests of Predictability: Which One Should We Use? *Econometric Reviews 23*(4), 371–402.

James, G., D. Witten, T. Hastie, and R. Tibshirani (2013). *An Introduction to Statistical Learning with Applications in R.* Springer.

Janes, H., G. Longton, and M. Pepe (2009). Accommodating Covariates in ROC Analysis. *The Stata Journal 9*(1), 17.

Jordà, Ò., M. Schularick, and A. M. Taylor (2011). Financial Crises, Credit Booms, and External Imbalances: 140 Years of Lessons. *IMF Economic Review 59*(2), 340–378.

Jordà, O., M. Schularick, and A. M. Taylor (2015). Leveraged Bubbles. *Journal of Monetary Economics 76*, S1–S20.

Kaminsky, G. L. and C. M. Reinhart (1999). The Twin Crises: the Causes of Banking and Balance-of-Payments Problems. *American Economic Review 89*(3), 473–500.

Kaminsky, G. L. and C. M. Reinhart (2000). On Crises, Contagion, and Confusion. *Journal of International Economics 51*, 145–168.

Kilian, L. (1998). Small-Sample Confidence Intervals for Impulse Response Functions. *Review of Economics and Statistics 80*(2), 218–230.

Kindleberger, C. P. and R. Z. Aliber (2005). *Manias, Panics and Crashes – A History of Financial Crises.* Palgrave Macmillan, Hampshire and New York.

Knedlik, T. and G. von Schweinitz (2012). Macroeconomic Imbalances as Indicators for Debt Crises in Europe. *JCMS: Journal of Common Market Studies 50*(5), 726–745.

Laeven, L. and F. Valencia (2012). Systemic Banking Crises Database: An Update. IMF Working Paper 12/163.

LeCun, Y., Y. Bengio, and G. Hinton (2015). Deep learning. *Nature 521*(7553), 436.

Lim, C. H., A. Costa, F. Columba, P. Kongsamut, A. Otani, M. Saiyid, T. Wezel, and X. Wu (2011). Macroprudential Policy: What Instruments and How to Use Them? Lessons from Country Experiences. IMF Working Paper 11/238.

Lo Duca, M., A. Koban, M. Basten, E. Bengtsson, B. Klaus, P. Kusmierczyk, J. H. Lang, C. Detken, and T. Peltonen (2017). A New Database for Financial Crises in European Countries - ECB/ESRB EU Crises Database. ECB Occasional Paper No 194.

Lo Duca, M. and T. A. Peltonen (2013). Assessing Systemic Risks and Predicting Systemic Events. *Journal of Banking & Finance 37*(7), 2183–2195.

Maddaloni, A. and J.-L. Peydró (2011). Bank Risk-taking, Securitization, Supervision, and Low Interest Rates: Evidence from the Euro-area and the U.S. Lending Standards. *The Review of Financial Studies 24*(6), 2121–2165.

McFadden, D. L. (1984). Econometric Analysis of Qualitative Response Models. *Handbook of econometrics 2*, 1395–1457.

Minsky, H. P. (1982). *Can "It" happen again? – Essays on Instability and Finance.* M.E. Sharpe Inc., Armonk, N.Y.

Murphy, K. (2012). *Machine Learning: A Probabilistic Perspective.* MIT Press.

Neunhoeffer, M. and S. Sternberg (2018). How cross-validation can go wrong and what to do about it. *Political Analysis. Forthcoming.*

Rajan, R. G. (2005). Has Financial Development Made the World Riskier? NBER Working paper 11728.

Reinhart, C. M. and K. S. Rogoff (2008). Is the 2007 US Sub-Prime Financial Crisis so Different? An International Historical Comparison. *The American Economic Review 98*(2), 339–344.

Reinhart, C. M. and K. S. Rogoff (2009). *This Time is Different: Eight Centuries of Financial Folly.* Princeton University Press, Princeton and Woodstock.

Rose, A. K. and M. M. Spiegel (2012). Cross-Country Causes and Consequences of the 2008 Crisis: Early Warning. *Japan and the World Economy 24*(1), 1–16.

Sarlin, P. (2013). On Policymakers' Loss Functions and the Evaluation of Early Warning Systems. *Economics Letters 119*(1), 1–7.

Sarlin, P. and G. von Schweinitz (2017). Optimizing Policymakers' Loss Functions in Crisis Prediction: Before, Within or After? ECB Working Paper 2025.

Schüler, Y. S., P. Hiebert, and T. A. Peltonen (2015). Characterising the Financial Cycle: a Multivariate and Time-Varying Approach. ECB Working Paper Series 1846, European Central Bank.

Shalev-Shwartz, S. and S. Ben-David (2014). *Understanding Machine Learning: from Theory to Algorithms.* Cambridge University Press.

Tanaka, K., T. Kinkyo, and S. Hamori (2016). Random Forests-Based Early Warning System for Bank Failures. *Economics Letters 148*, 118–121.

# 6    Appendix A: Methodology

## A.1    Definition of Dependent Variable

Early warning models (typically) perform window forecasts of crisis probabilities and use thresholds to derive binary signals from these probabilities. The rationale behind this approach has two aspects. First, window forecasts are used since it is hard to predict the exact quarterly start date of a crisis, which may be driven to a large extent by unforecastable shocks. However, recurring patterns before crises may still be informative about their likelihood of occurrence during a given time *interval*. Thus, window forecasts of the probability of a systemic banking crisis can be used to reflect potential buildups of vulnerabilities, which might require for instance the activation of macroprudential policy measures. Second, converting probabilities into clear signals (taking into account the policymaker's preferences) helps to inform policymakers' ultimate decision on whether and when to take action. Moreover, it allows for a straightforward evaluation of predictions in terms of correct or incorrect signals.

To implement these ideas, we follow the literature and define the dependent variable for our estimations as follows. Starting from a crisis database, where $C_{t,n}$ is 1 if a crisis was ongoing in country $n$ at time $t$ and zero otherwise, we define another binary variable $\bar{C}_{t,n}$ as our dependent variable. This dependent variable $\bar{C}_{t,n}$ is set to one during early warning windows between $h_1$ and $h_2$ periods before a crisis (pre-crisis periods) and zero for observations that are not followed by a crisis within the next $h_2$ quarters (tranquil periods).

The resulting gap of length $h_1-1$ between early warning windows and crises is excluded from the estimation, as these periods can neither be classified as being in an early warning window, nor as being tranquil periods. Moreover, it is standard to exclude periods where a country is already in a crisis (crisis periods). The reason for excluding crisis periods is that the extreme imbalances during these periods are typically due to being in a crisis (which is assumed to be known), instead of reflecting the buildup of imbalances prior to a crisis.[20]

The timing of the early warning window is chosen to fulfill two criteria. First, the gap $h_1 \geq 0$ between the window and the start of the crisis is chosen to allow for policy action. Second, the window needs to be sufficiently close to the predicted crisis for economic variables to show informative developments. Following the literature, we set the limits of early warning windows to $h_1 = 5$ and $h_2 = 12$ quarters (see Drehmann and Juselius, 2014; Alessi and Detken, 2018; Holopainen and Sarlin, 2017). This allows at least one year for policy measures to become effective and to issue warnings up to three years before a crisis.

In sum, this leads to the following definition of the dependent variable:

$$
\bar{C}_{t,n} = \begin{cases} 0 & \text{, if } C_{t+h,n} = 0, \text{ for all } h \in \{0, \dots, h_2\} \\ 1 & \text{, if } C_{t+h,n} = 1, \text{ for some } h \in \{h_1, \dots, h_2\} \text{ and} \\ & \quad C_{t+h,n} = 0, \text{ for all } h \in \{0, \dots, h_1 - 1\} \\ NA & \text{, otherwise.} \end{cases} \tag{2}
$$

---

[20]We do not exclude additional periods after a crisis, as crises in our database are defined such that they already account for the post-crisis bias discussed in Bussière and Fratzscher (2006).

We thus estimate the probability of a crisis *starting* between the next $h_1 = 5$ to $h_2 = 12$ quarters, conditional on not already being in an acute crisis period. To do this, the binary dependent variable $\bar{C}$ is linked to a set of early warning indicators $X$ using different modeling choices. Each model is estimated and then used to predict the probability of being in an early warning window at time $t$ in country $n$ conditional on the observables $X$, $P(\bar{C}_{t,n}|X_{t,n})$. For the sake of brevity, we simply refer to this as 'crisis probability'.

## A.2 Description of Estimation Methods

This section provides a brief overview of each method and highlights some key aspects for each method. Table A.1 summarizes the discussion in this section by highlighting some key benefits and drawbacks of the employed methods. Of course, this is only a snapshot of the more complete description of these methods in the mentioned references.[21]

Table A.1: Comparison of employed methods: benefits and drawbacks

|  | **Benefits** | **Drawbacks** |
|---|---|---|
| **logit** | explicit probabilistic foundations high interpretability | pre-specified functional form |
| **knn** | simple approach | strong curse of dimensionality |
| **trees** | automatic variable selection intuitive approach | instability across time / different samples |
| **rf** | more stable than trees improves on tree accuracy | risk of overfitting identifying drivers of predictions complex |
| **svm** | flexible nonlinear fitting computationally efficient | risk of overfitting ad hoc in probabilistic setups difficult to communicate |

**Logistic regression (logit):** Logit models are the workhorse models in the early warning literature (Frankel and Rose, 1996; Bussière and Fratzscher, 2006; Lo Duca and Peltonen, 2013). They are based on two assumptions. First, the dependent binary variable is assumed to be driven by a latent process $y^*$, which is in turn linearly related to the employed explanatory variables: $y^* = X\beta + \varepsilon$. Second, the latent process is assumed to be linked to the binary variable by a logistic transformation (or, equivalently, estimation errors $\varepsilon$ follow a logistic distribution). Hence, a key advantage of logit models is that they are based on a clear and straightforward statistical model, which explicitly takes uncertainty into account. Compared to machine learning methods, they are easy to interpret (for instance, in terms of coefficients), but, at the same time, restricted to the specific functional form just described. A key issue in their estimation is to make sure that a sufficient number of observations in each category is available (McFadden, 1984). In the

---

[21]More detailed introductions to these methods may be found, for instance, in Murphy (2012), James et al. (2013), or Shalev-Shwartz and Ben-David (2014).

context of early warning models, it is crucial to have a sufficient number of pre-crisis periods (which are much less frequent than tranquil periods) available for estimation. When the number of crisis events contained in the sample is reduced, estimation uncertainty increases, and, in the extreme case, perfect discrimination can prevent a proper estimation of the model's parameters. To put the logit method on equal footing with the machine learning methods, we estimate a non-dynamic logit model, pooling observations both in the cross-section and the time dimension.

**K nearest neighbors (knn):** The idea of knn[22] (Cover and Hart, 1967) is to predict the probability of an event for a given observation $(t, n)$ by the share of such an event among its $K$ closest (nearest) neighbors. Closeness of two observations $\mathbf{x}$ and $\mathbf{x}'$ is measured by the Euclidean distance, i.e. $||\mathbf{x} - \mathbf{x}'|| = \sqrt{\sum_{i=1}^{d}(x_i - x_i')^2}$. That is, two observations are close if the realizations of the explanatory variables associated with these observations are similar. Formally, we define a neighborhood $N_K(X_{t,n})$ around every observation $X_{t,n}$, containing the $K$ closest observations to $X_{t,n}$ in the training sample. The probability of an event is the average occurrence of the event in the neighborhood, i.e. $P(\bar{C}_{t,n} = 1) = \frac{1}{K}\sum_{k \in N_K(X_{t,n})} \bar{C}_k$. The hyperparameter K is chosen by cross-validation. Moreover, we use a knn algorithm that refines the method by weighting each of the neighboring points by their distance to the given point $X_{t,n}$. A key problem of knn is that it is subject to a strong "curse of dimensionality". Shalev-Shwartz and Ben-David (2014) show that the sample size required to achieve a given error grows exponentially with the number of explanatory variables in the dataset. In our empirical application, we use sets of explanatory variables of different dimension in order to gain insights on the tradeoff between additional information and additional complexity.

**Decision trees (trees):** Binary decision trees[23] (Breiman, Friedman, Olshen, and Stone, 1984) consist of a root, interior nodes (branches) and final nodes (leafs). The root and every branch consist of a decision rule based on a single explanatory variable $x_i$ and a threshold $\tau_i$. The decision rules assign observations to the left subtree if $x_i > \tau_i$ and to the right subtree otherwise. Starting at the root, observations are thus passed down the tree until they end up in a final node. For every node, the (predicted) probability of an event is equal to the average occurrence of said event among observations from the training sample assigned to the same final node.

The estimation of the tree entails choosing simultaneously the variables $x$ and thresholds $\tau$ to split on. Efficient algorithms have been developed for approximating the optimal solution to this non-trivial task. These proceed by starting at the root and recursively constructing the tree from there, based on a measure of gain from each considered split and several potential stopping criteria for limiting the complexity of the tree. In our case, the number of branches is determined by a "pruning" parameter which balances increasing complexity against the homogeneity in final leaves.[24]

---

[22]We implement knn using the R-package 'kknn'.

[23]We implement decision trees using the R package 'rpart'.

[24]This is, of course, only one way to limit tree complexity. Other approaches, for example, set a minimum number of observations per final node (used in our implementation of random forest), or, alternatively, a maximum number of final nodes.

The selection of the pruning parameter (the hyperparameter of this method) thus decides on the complexity of the tree. Lower complexity costs imply additional splits which decrease classification errors on the training sample and thus increase the sharpness of estimated probabilities (pushing them closer to either zero or one). At the same time, the larger number of final nodes implies fewer training observations per final node, which increases estimation uncertainty and the potential for overfit. As the sensitivity of estimated trees to small changes in the underlying dataset can be high, the method of random forests has been developed to mitigate this undesirable feature.

**Random forest (rf):** Random forests[25] ([Breiman, 1996, 2001](#)) generalize decision trees by averaging over the predictions of a large number of different decision trees. This can reduce the variance of estimates and, hence, prediction errors. Random forests generate heterogeneity among its trees by (a) estimating trees on randomly chosen subsets of observations (also called bootstrap aggregating or bagging), and (b) considering only a randomly chosen subset of early warning indicators at each split (also called random subspace method, or attribute bagging). Both components are needed in order to de-correlate individual trees sufficiently, so as to achieve the desired variance reduction, while maintaining a high degree of prediction accuracy.

To put the random forest method into practice, we have to select three different hyperparameters. First, we set the number of trees used in each random forest to 1'000 such that the average prediction of the trees in the forest converges. Cross-validation is used to set the further two hyperparameters of this method. Heterogeneity between trees is driven largely by the number of randomly drawn variables to be considered at each split (the second hyperparameter). Third, complexity of the trees in the forest is limited by setting a minimum number of observations per terminal node, which is the third hyperparameter of this method .

Random forests have so far been the most frequently employed machine learning method in the early warning literature ([Alessi and Detken, 2018](#); [Holopainen and Sarlin, 2017](#); [Tanaka et al., 2016](#)). However, their success in reducing variance and improving out-of-sample performance depends on achieving a sufficiently low correlation between the randomly generated trees ([Breiman, 2001](#)). We conjecture that achieving such a low degree of correlation could be especially challenging in the presence of serial and cross-sectional correlation of the underlying training data.

**Support vector machine (svm):** svm[26] constructs a hyperplane in order to separate observations into distinct groups, pre-crisis and tranquil periods in our case. When the data is linearly separable, the main question is which hyperplane to choose from an infinite space of possible separating hyperplanes. svm uniquely determines the hyperplane by maximizing the distance of the two closest observations to the separating hyperplane (this distance is called margin).

For illustration, let us consider a one-dimensional example. Suppose a dataset is univariate, with observations given as points on the real line, namely $x = \{-3, -2, -1, 1, 2, 3\}$ and suppose that observations are linearly separable, namely $y(x = \{-3, -2\}) = 1$ and

---

[25] We implement random forests using the R package 'randomForest'.

[26] We implement support vector machines using the R package 'e1071'.

$y(x = \{-1, 1, 2, 3\}) = 0$. Then, obviously, any rule which assigns $y(x < -2) = 1$ and $y(x > -1) = 0$ perfectly separates the observations. The svm method would choose the point $-1.5$ for a separating rule that maximizes the margin. Obviously, we achieve separation by a point in this one-dimensional example, by a line in 2-d, and by a hyperplane in 3-d or higher.

However, in typical applications observations are not linearly separable. Consider a modification of the above example where $y = 1$ for all observations where $\mid x \mid > 2$ and zero otherwise (this example is inspired by Shalev-Shwartz and Ben-David (2014), p. 179). This is not linearly separable in the original space, but can be made separable by mapping to a two-dimensional space, for instance by using $\phi : \mathbb{R} \to \mathbb{R}^2$ where $\phi(x) = (x, x^2)$. Then, any rule that assigns, $y = 1$ whenever $x^2 > 4$ separates the observations. This simple example illustrates the more general idea that is typically used in combination with svm: By mapping non-linear transformations of the original data into a higher dimensional space, linear separability of the dataset can be achieved or, at least, enhanced.

Mapping data into higher dimensional feature spaces enhances the expressiveness of methods (enlarging the space of functions considered for describing the data), but, obviously, higher dimensionality comes at the cost of increased complexity (the number of parameters can rise exponentially in the multivariate case, for example when polynomials using cross-products of variables are considered). To deal with this issue, the machine learning literature has developed what is known as 'kernel trick'. This allows an efficient computation of svm classifiers when such non-linear mappings into high-dimensional spaces are used. Broadly speaking, kernel functions describe similarities between observations and have special properties that allow the svm calculation to be based on these kernel functions without explicitly handling the high-dimensional representation of the data. A formal description of the kernel trick may be found, for instance, in Shalev-Shwartz and Ben-David (2014), pp. 181. The complexity of the high-dimensional function space is controlled by a hyperparameter gamma inside the kernel function, with higher values of gamma leading to more complexity (we use the standard choice of a radial basis kernel).

While it is possible to linearly separate observations after mapping them into an arbitrarily complex space (n observations can always be perfectly fitted using a n-1 degree polynomial), this is generally not desirable. Instead, a penalty term for misclassified observations is added to the svm classifier loss function. Allowing for misclassification makes it possible to use the parameter svmg (see table A.2) to separately control the complexity of the classifier. Moreover, it induces a tradeoff between large margins and misclassification, which is controlled by a second hyperparameter svmc (cost of soft margin constraint violation). The more tolerant we are towards misclassification on the training sample, the larger the margin can be (ceteris paribus). A larger margin then makes the classification more robust towards perturbations of the original data, for example when predicting the label of new data points. Thus, both hyperparameters seek to strike a balance between perfectly fitting the training data (potentially overfitting relative to the true model), and correctly classifying new data.

Support vector machines are among the most frequently used machine learning algorithms. Their ability to flexibly fit complex functions to the data at the same time entails the risk of overfitting. Moreover, the probabilistic foundations for the svm method are rather ad-hoc (Murphy, 2012). In the early warning context, the presence of a substantial unpredictable component as well as of cross-sectional and serial correlation may

dampen the method's out-of-sample performance (see results section and discussion). As can be seen from the relatively long explanation in this paragraph, understanding and communicating this method may present additional challenges.

## A.3    Panel-block-bootstrap

The aim of our panel-block-bootstrap is to draw random datasets with similar autocorrelation and cross-sectional dependence patterns as in the original dataset. To achieve this, we construct bootstrap datasets from blocks of observations that are jointly sampled from the original dataset. Drawing blocks of consecutive observations retains the autocorrelation structure of the data, while the panel structure of blocks captures cross-sectional correlation.

For every estimation, we sample $R = 1'000$ different bootstrap datasets from the respective training sample, which covers the time-country specific observations $\{(t,n)|t \in \{1,\ldots,T\},\ n \in \{1,\ldots,N\}\}$. A block $B_t$ (with blocklength $b = 8$) starts at time $t$ and contains the following observations of both $X$ and $\bar{C}$:

$$B_t = \begin{bmatrix} (t,1) & \cdots & (t,N) \\ \vdots & \ddots & \vdots \\ (t+b-1,1) & \cdots & (t+b-1,N) \end{bmatrix} \tag{3}$$

Bootstrap samples $r \in \{1,\ldots,R\}$ are drawn randomly from the original data such that every observation has an equal probability of entering the random sample. Thus, we proceed as follows:

1. Initialize with an empty bootstrap sample $r = \emptyset$.

2. Draw a random starting period $t^* \in \{2-b,\ldots,T\}$. If we would not allow for early or late starting periods (that effectively generate blocks with missing observations), observations at the beginning or end of the original sample would have a lower probability of entering the bootstrap sample.

3. Obtain $B_{t^*}$, corresponding to $t^*$, from the original training dataset. Some observations may be missing due to (a) shorter sample length for an individual country, (b) an early starting period $t^* < 1$ or (c) a late starting period $t^* > T - b$. In this case, only include nonempty observations in $B_{t^*}$.

4. Concatenate the bootstrap sample $r$ and $B_{t^*}$.

5. If the bootstrap sample $r$ has fewer observations than the original in-sample dataset, return to step 2. Otherwise, return the bootstrap sample $r$.

We estimate every model $m$ on every bootstrap sample $r \in \{1,\ldots,R\}$. Results are used to predict probability estimates $p_{t,n}^{m,r}$ for every observation $(t,n)$ in the test sample. From this, we calculate the different performance measures (relative usefulness, AUC and BPS) for every bootstrap sample $r$ and model $m$. The bootstrap distribution of

performance measures across r yields estimates of confidence bands for each model m.[27] Moreover, it allows us to test whether model $m_1$'s performance is statistically significantly better than model $m_2$'s performance. For example, the probability that the relative usefulness of model $m_1$ is higher than that of model $m_2$ is given by $\frac{1}{R}\sum_{r=1}^{R} 1_{U_r^{m_1,r} > U_r^{m_2,r}}$.

## A.4 Cross-validation

We use cross-validation to select optimal hyperparameters from a predefined grid. The idea of cross-validation is to obtain an estimate of how well a model is able to make predictions on previously unseen data. To this end, the sample is cut repeatedly into an estimation sample and a test sample. In this sense, cross-validation is similar to out-of-sample prediction, but without paying as much attention to the time dimension of the dataset. In particular, we use panel block leave-p-out cross-validation (Arlot, Celisse, et al., 2010). In this variant, a block of twelve consecutive quarters (corresponding to the horizon of the early warning window) across all countries is used as test sample, while all other observations are included in the training sample. This is done repeatedly for all possible blocks until the sample observations are exhausted.

Using whole blocks of observations in the test sample instead of randomly selected observations has two advantages. First, it captures the serial and cross-sectional correlation of the data in a similar way as recursive out-of-sample estimation. Second, the number of possible splits of the dataset is limited, making an exhaustive cross-validation over all possible combinations possible. Using all possible splits into blocks of twelve quarters causes each observation (time $t$, country $n$) to be contained in twelve different panel blocks. That is, for each observation $(t, n)$, the panel block leave-p-out cross-validation results in twelve different predictions for every model (with each model being defined by a combination of method, hyperparameters, and explanatory variables). In order to calculate the performance measure for a given model, we average performance over all cross-validation predictions from that model. We can then perform a grid search to select for each method the hyperparameters maximizing its cross-validation performance. Table A.2 displays the resulting optimal hyperparameters.

---

[27]A well-known problem of bootstrapping is its potential for bias in small samples (Kilian, 1998). We need to correct for this in the bootstrap distribution of our performance measures. In case of relative usefulness and the Brier probability score, we mean-adjust the bootstrap distribution underlying confidence bands and p-values. In the case of the AUC, confidence bands based on this bias correction method appear implausible, such that we recourse to a non-parametric approach developed specifically for this measure (DeLong, DeLong, and Clarke-Pearson, 1988).

Table A.2: Hyperparameters for machine learning methods (for baseline results)

| Method | Hyperparameter name | Opt. value | Hyperparameter name | Opt. value |
|---|---|---|---|---|
| trees.1 | cp (tree complexity | 0.0212 | | |
| trees.2 | parameter controlling | 0.0273 | | |
| trees.3 | cost of adding another | 0.0121 | | |
| trees.4 | split to the tree) | 0 | | |
| knn.1 | k (number of nearest | 49 | | |
| knn.2 | neighbours to use for | 7 | | |
| knn.3 | each prediction) | 29 | | |
| knn.4 | | 15 | | |
| rf.1 | nodesize (minimum | 16 | rfmtry (number of | 2 |
| rf.2 | number of observations | 16 | variables randomly | 2 |
| rf.3 | per terminal node of | 2 | sampled as | 2 |
| rf.4 | each tree in the forest) | 16 | candidates at each | 9 |
| | | | split) | |
| svm.1 | | 0.5 | svmc (cost of soft | 0.0156 |
| svm.2 | | 0.5 | margin constraint | 0.0211 |
| svm.3 | svmg (parameter in | 0.0005 | violation) | 0.0267 |
| svm.4 | radial basis function) | 0.1250 | | 8 |

35

# 7  Appendix B: Data

Table B.1: Country coverage and crisis dates

| Country | Data availability | | No. of quarters | Crisis dates | | | | | |
| | Start | End | | Start | End | Start | End | Start | End |
|---|---|---|---|---|---|---|---|---|---|
| BE | 1975 Q1 | 2016 Q2 | 166 | no crisis | | | | | |
| DE | 1971 Q1 | 2016 Q2 | 182 | 1974 Q2 | 1974 Q4 | 2001 Q1 | 2003 Q4 | | |
| DK | 1975 Q1 | 2016 Q1 | 165 | 1987 Q1 | 1995 Q1 | 2008 Q1 | 2013 Q4 | | |
| ES | 1975 Q1 | 2016 Q1 | 165 | 1978 Q1 | 1985 Q3 | 2009 Q1 | 2013 Q4 | | |
| FI | 1981 Q4 | 2016 Q2 | 139 | 1991 Q3 | 1996 Q4 | | | | |
| FR | 1973 Q1 | 2016 Q1 | 173 | 1991 Q2 | 1995 Q1 | 2008 Q2 | 2009 Q4 | | |
| GB | 1972 Q1 | 2016 Q1 | 177 | 1973 Q4 | 1975 Q4 | 1991 Q3 | 1994 Q2 | 2007 Q3 | 2010 Q1 |
| IE | 1990 Q4 | 2016 Q1 | 102 | 2008 Q3 | 2013 Q4 | | | | |
| IT | 1971 Q1 | 2016 Q1 | 181 | 1991 Q3 | 1997 Q4 | 2011 Q3 | 2013 Q4 | | |
| JP | 1971 Q1 | 2016 Q2 | 182 | 1997 Q4 | 2001 Q4 | | | | |
| NL | 1971 Q1 | 2016 Q1 | 181 | 2008 Q1 | 2013 Q1 | | | | |
| NO | 1975 Q1 | 2016 Q2 | 166 | 1988 Q3 | 1992 Q4 | | | | |
| PT | 1988 Q1 | 2016 Q2 | 114 | 2008 Q4 | ongoing | | | | |
| SE | 1975 Q1 | 2016 Q1 | 165 | 1991 Q1 | 1997 Q2 | 2007 Q4 | 2010 Q4 | | |
| US | 1971 Q1 | 2016 Q2 | 182 | 1988 Q1 | 1995 Q4 | | | | |

Table B.2: Data transformation and sources - variables used

| Category | Variable name | Source | Transformation / Description (for non-transformed variables) | Inputs | Combinations |
|---|---|---|---|---|---|
| Asset prices | Real equity price gap | OECD, Eurostat and own calculations | Relative gap using HP filter with lambda of 400,000 | Equity prices, consumer price index | |
| Asset prices | Real residential real estate price gap | OECD, Eurostat and own calculations | Relative gap using HP filter with lambda of 400,000 of real residential real estate price | Residential real estate prices, consumer price index | |
| Credit | Total credit-to-GDP gap | BIS, OECD, Eurostat and own calculations | Absolute gap using HP filter with lambda of 400,000 of total credit-to-GDP ratio | Total credit, nominal GDP | |
| External | Current account-to-GDP ratio | OECD, Eurostat and own calculations | OECD data: Current account to GDP without transformations, Eurostat data: Current account to GDP calculated as ratio of current account balance to GDP (both summed up over four quarters) | Current account balance, Nominal GDP | Take longest time series available of OECD data or Eurostat data |
| External | Real oil price gap | OECD and own calculations | Relative gap using HP filter with lambda of 1,600 of real oil price | Oil price, consumer price index (US) | |
| External | Real effective exchange rate gap | OECD, IMF and own calculations | Relative gap using HP filter with lambda of 400,000 of real effective exchange rate | Real effective exchange rate | |
| Macro | 3-month real money market rate | OECD, ECB, Eurostat and own calculations | Real interbank lending rate | Nominal 3-month money market rate, consumer price index | |
| Macro | Inflation rate | Eurostat, OECD | Annual rate of inflation (y-o-y growth rate of quarterly data) | Consumer price index | |
| Macro | Real GDP gap | OECD, Eurostat and own calculations | Relative gap using HP filter with lambda of 1,600 of real GDP | Nominal GDP, consumer price index | |
| Macro | Gross fixed capital formation-to-GDP gap | OECD, Eurostat, Bundesbank and own calculations | Absolute gap using HP filter with lambda of 1,600 of gross fixed capital formation-to-GDP ratio | Gross fixed capital formation, nominal GDP | |

Table B.3: Data transformation and sources - input data

| Variable name | Source | Transformation / Description (for non-transformed variables) | Combinations |
|---|---|---|---|
| Consumer price index | Eurostat, OECD | Consumer price index, end of quarter values, re-based to 2015=100 | Take longest time series available of OECD data or Eurostat data |
| Total credit | BIS | Total credit to the private non-financial sector, domestic currency, billions | |
| Oil price | OECD | Brent crude oil price, USD per barrel | |
| Real effective exchange rate | OECD, IMF | Real effective exchange rate, CPI based index, base year: 2010 | Take longest time series available of OECD data or IMF data |
| Current account balance | OECD, Eurostat | OECD: Current acount balance as percentage of GDP Eurostat: Current account balance (own calculations: sum of last four quarters, as percentage of GDP) | Take longest time series available of OECD data or Eurostat data |
| Nominal GDP (national currency) | OECD, Eurostat | Gross domestic product at market prices, seasonally adjusted, domestic currency, billions, sum of last four quarters | Take longest time series available of OECD data or Eurostat data |
| Nominal GDP (in EUR, for current account-to-GDP calculation) | Eurostat | Gross domestic product at market prices, seasonally adjusted, euro, millions, sum of last four quarters | |
| Gross fixed capital formation | OECD, Eurostat, Bundesbank | Gross fixed capital formation, domestic currency, millions. For DE: Bundesbank data (including calculations) for long time series of GFCF | Take longest time series available of OECD data or Eurostat data |
| 3-month nominal money market rate | OECD, ECB, Datastream | Interbank interest rate, average through quarter | Take longest time series available of OECD, ECB and Datastream data |
| Equity prices | OECD, Bloomberg, Datastream | Equity price index, 2010=100, average through quarter | Take longest time series available of OECD, Bloomberg and Datastream data |
| Residential real estate prices | OECD | Index of residential real estate price, based in 2010, seasonally adjusted. | |

Table B.4: Summary statistics: Growth rate dataset

| Variable name | Pre-crisis periods | | | | | Tranquil periods | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | mean | sd | min | max | obs | mean | sd | min | max | obs |
| 4-year growth rate of credit-to-GDP ratio | 0.62 | 1.04 | -1.19 | 3.46 | 176 | -0.07 | 0.97 | -2.34 | 3.46 | 1525 |
| 4-year growth rate of real residential real estate prices | 0.52 | 1.09 | -1.27 | 2.63 | 176 | -0.06 | 0.97 | -2.47 | 2.63 | 1525 |
| 4-year growth rate of real equity prices | 0.55 | 1.01 | -1.12 | 3.67 | 176 | -0.06 | 0.98 | -1.35 | 3.67 | 1525 |
| 1-year growth rate of real GDP | 0.09 | 0.88 | -3.38 | 1.78 | 176 | -0.01 | 1.01 | -3.98 | 4.65 | 1525 |
| Inflation rate | -0.07 | 0.90 | -1.22 | 3.87 | 176 | 0.01 | 1.01 | -1.37 | 5.11 | 1525 |
| 1-year growth rate of gross fixed capital formation-to-GDP ratio | 0.27 | 1.09 | -3.96 | 3.75 | 176 | -0.03 | 0.98 | -3.96 | 3.75 | 1525 |
| 3-month real money market rate | 0.07 | 1.10 | -3.33 | 2.38 | 176 | -0.01 | 0.99 | -3.33 | 2.67 | 1525 |
| Current account-to-GDP ratio | -0.75 | 1.14 | -3.22 | 2.06 | 176 | 0.09 | 0.95 | -3.22 | 2.60 | 1525 |
| 4-year growth rate of real effective exchange rate | 0.24 | 0.82 | -1.45 | 3.48 | 176 | -0.03 | 1.02 | -2.21 | 3.48 | 1525 |
| 1-year growth rate of real oil price | 0.02 | 0.88 | -1.82 | 3.10 | 176 | 0.00 | 1.01 | -1.82 | 3.79 | 1525 |

*Note:* Data have been standardized using the unconditional mean and standard deviation across all periods. Since data is winsorized at the 1% and 99% level, minimum (maximum) values may be the same in pre-crisis and tranquil periods.

Table B.5: Variables used in specifications (1) - (4)

| Credit and Asset Prices (1) | Macro (2) | External (3) | All (4) |
|---|---|---|---|
| Total credit-to-GDP gap | Real GDP gap | Current account-to-GDP ratio | Total credit-to-GDP gap |
| Real residential real estate price gap | Inflation rate | Real effective exchange rate gap | Real residential real estate price gap |
| Real equity price gap | 3-month real money market rate | Real oil price gap | Real equity price gap |
| | Gross fixed capital formation-to-GDP gap | | Real GDP gap |
| | | | Inflation rate |
| | | | 3-month real money market rate |
| | | | Gross fixed capital formation-to-GDP gap |
| | | | Current account-to-GDP ratio |
| | | | Real effective exchange rate gap |
| | | | Real oil price gap |

# 8    Appendix C: Results

Table C.1: *In-sample* performance using different performance measures

| | Threshold | TP | FP | TN | FN | FPrate | FNrate | U_r | | AUC | | BPS | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| logit.1 | 0.109 | 102 | 402 | 1207 | 69 | 0.25 | 0.40 | 0.347*** | [0.088, 0.457] | 0.736*** | [0.703, 0.769] | 0.077*** | [0.069, 0.091] |
| trees.1 | 0.097 | 84 | 96 | 1513 | 87 | 0.06 | 0.51 | 0.432*** | [0.172, 0.576] | 0.776*** | [0.742, 0.810] | 0.063*** | [0.046, 0.089] |
| knn.1 | 0.126 | 152 | 315 | 1294 | 19 | 0.20 | 0.11 | 0.693*** | [0.419, 0.825] | 0.928*** | [0.915, 0.942] | 0.057*** | [0.045, 0.074] |
| rf.1 | 0.198 | 170 | 56 | 1553 | 1 | 0.03 | 0.01 | **0.959** | [0.609, 1.000] | **0.994** | [0.992, 0.996] | **0.029** | [0.011, 0.053] |
| svm.1 | 0.087 | 61 | 48 | 1561 | 110 | 0.03 | 0.64 | 0.327*** | [0.080, 0.460] | 0.845*** | [0.817, 0.874] | 0.073*** | [0.050, 0.085] |
| logit.2 | 0.098 | 107 | 682 | 927 | 64 | 0.42 | 0.37 | 0.202*** | [0.114, 0.306] | 0.613*** | [0.576, 0.651] | 0.086*** | [0.079, 0.099] |
| trees.2 | 0.096 | 171 | 1609 | 0 | 0 | 1.00 | 0.00 | 0.000*** | [-0.099, 0.135] | 0.500*** | [0.500, 0.500] | 0.087*** | [0.074, 0.112] |
| knn.2 | 0.315 | 171 | 57 | 1552 | 0 | 0.04 | 0.00 | **0.965** | [0.693, 1.000] | 0.988*** | [0.984, 0.991] | **0.036** | [0.014, 0.068] |
| rf.2 | 0.167 | 168 | 42 | 1567 | 3 | 0.03 | 0.02 | 0.956 | [0.683, 1.000] | **0.993** | [0.991, 0.996] | 0.041 | [0.022, 0.063] |
| svm.2 | 0.093 | 122 | 324 | 1285 | 49 | 0.20 | 0.29 | 0.512*** | [0.340, 0.676] | 0.852*** | [0.827, 0.877] | 0.073** | [0.050, 0.097] |
| logit.3 | 0.097 | 118 | 482 | 1127 | 53 | 0.30 | 0.31 | 0.390*** | [0.271, 0.461] | 0.745*** | [0.712, 0.778] | 0.081*** | [0.077, 0.089] |
| trees.3 | 0.096 | 171 | 1609 | 0 | 0 | 1.00 | 0.00 | 0.000*** | [-0.207, 0.147] | 0.500*** | [0.500, 0.500] | 0.087*** | [0.074, 0.107] |
| knn.3 | 0.150 | 148 | 291 | 1318 | 23 | 0.18 | 0.13 | 0.685*** | [0.470, 0.832] | 0.919*** | [0.905, 0.932] | 0.062*** | [0.051, 0.075] |
| rf.3 | 0.352 | 171 | 0 | 1609 | 0 | 0.00 | 0.00 | **1.000** | [0.702, 1.000] | **1.000** | [1.000, 1.000] | **0.012** | [-0.009, 0.035] |
| svm.3 | 0.136 | 64 | 232 | 1377 | 107 | 0.14 | 0.63 | 0.230*** | [0.015, 0.428] | 0.610*** | [0.569, 0.650] | 0.085*** | [0.062, 0.082] |
| logit.4 | 0.109 | 119 | 297 | 1312 | 52 | 0.18 | 0.30 | 0.511*** | [0.294, 0.655] | 0.810*** | [0.779, 0.840] | 0.073** | [0.060, 0.094] |
| trees.4 | 0.100 | 138 | 214 | 1395 | 33 | 0.13 | 0.19 | 0.674*** | [0.406, 0.883] | 0.901*** | [0.879, 0.923] | 0.041** | [0.015, 0.075] |
| knn.4 | 0.339 | 166 | 26 | 1583 | 5 | 0.02 | 0.03 | 0.955 | [0.660, 1.000] | 0.997** | [0.996, 0.999] | 0.021** | [-0.004, 0.046] |
| rf.4 | 0.248 | 171 | 16 | 1593 | 0 | 0.01 | 0.00 | **0.990** | [0.672, 1.000] | **0.999** | [0.999, 1.000] | 0.018* | [-0.007, 0.046] |
| svm.4 | 0.090 | 162 | 28 | 1581 | 9 | 0.02 | 0.05 | 0.930 | [0.566, 1.000] | 0.989*** | [0.983, 0.995] | **0.013** | [-0.028, 0.025] |

*Note*: Best performance on given dataset in bold. Stars indicate if the respective method's performance is significantly worse than the best model's performance on the same dataset (\*\*\*/\*\*/\* for the 1%/5%/10% significance level). For $U_r$ and AUC higher values indicate better prediction performance, while for BPS lower values are preferable. Numbers in brackets indicate 90% confidence bands. Model names reported in format "method.dataset", where datasets consist of the following sets of variables (1): credit and asset prices, (2) macro variables, (3) external imbalances, and (4) all variables (see also table B.5).

Table C.2: *Out-of-sample* performance using different performance measures

| | Threshold | TP | FP | TN | FN | FPrate | FNrate | U_r | | AUC | | BPS | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| logit.1 | 0.096 | 46 | 96 | 145 | 14 | 0.40 | 0.23 | **0.368** | [0.237,0.487] | **0.737** | [0.682,0.792] | **0.139** | [0.119,0.159] |
| trees.1 | 0.065 | 24 | 48 | 193 | 36 | 0.20 | 0.60 | 0.201** | [0.084,0.329] | 0.591*** | [0.51,0.672] | 0.148 | [0.125,0.171] |
| knm.1 | 0.124 | 30 | 61 | 180 | 30 | 0.25 | 0.50 | 0.247* | [0.123,0.374] | 0.704 | [0.644,0.765] | 0.153* | [0.134,0.172] |
| rf.1 | 0.193 | 22 | 29 | 212 | 38 | 0.12 | 0.63 | 0.246* | [0.134,0.358] | 0.73 | [0.669,0.79] | 0.15 | [0.131,0.17] |
| svm.1 | 0.079 | 27 | 50 | 191 | 33 | 0.21 | 0.55 | 0.243* | [0.118,0.355] | 0.492*** | [0.405,0.579] | 0.159 | [0.125,0.198] |
| logit.2 | 0.100 | 7 | 85 | 156 | 53 | 0.35 | 0.88 | -0.236* | [-0.394,-0.074] | 0.739 | [0.682,0.797] | 0.182 | [0.171,0.197] |
| trees.2 | 0.076 | 8 | 178 | 63 | 52 | 0.74 | 0.87 | -0.605*** | [-0.703,-0.499] | **0.824** | [0.773,0.874] | 0.204 | [0.18,0.232] |
| knm.2 | 0.309 | 4 | 37 | 204 | 56 | 0.15 | 0.93 | -0.087 | [-0.147,-0.023] | 0.385*** | [0.339,0.431] | 0.231** | [0.21,0.253] |
| rf.2 | 0.183 | 0 | 66 | 175 | 60 | 0.27 | 1.00 | -0.274*** | [-0.343,-0.194] | 0.771 | [0.726,0.815] | 0.22* | [0.203,0.239] |
| svm.2 | 0.087 | 15 | 80 | 161 | 45 | 0.33 | 0.75 | **-0.082** | [-0.198,0.039] | 0.687 | [0.619,0.756] | **0.182** | [0.148,0.222] |
| logit.3 | 0.087 | 40 | 55 | 186 | 20 | 0.23 | 0.33 | **0.438** | [0.34,0.532] | **0.762** | [0.699,0.825] | **0.153** | [0.142,0.164] |
| trees.3 | 0.077 | 56 | 131 | 110 | 4 | 0.54 | 0.07 | 0.39*** | [0.27,0.516] | 0.618*** | [0.567,0.67] | 0.181*** | [0.163,0.203] |
| knm.3 | 0.151 | 30 | 50 | 191 | 30 | 0.21 | 0.50 | 0.293** | [0.182,0.399] | 0.663*** | [0.594,0.732] | 0.159** | [0.144,0.175] |
| rf.3 | 0.431 | 9 | 8 | 233 | 51 | 0.03 | 0.85 | 0.117*** | [0.036,0.203] | 0.602*** | [0.527,0.677] | 0.164** | [0.148,0.18] |
| svm.3 | 0.089 | 31 | 110 | 131 | 29 | 0.46 | 0.48 | 0.06*** | [-0.11,0.215] | 0.696* | [0.639,0.753] | 0.176*** | [0.148,0.216] |
| logit.4 | 0.091 | 45 | 35 | 206 | 15 | 0.15 | 0.25 | **0.605** | [0.481,0.727] | **0.852** | [0.797,0.906] | **0.125** | [0.101,0.152] |
| trees.4 | 0.065 | 18 | 42 | 199 | 42 | 0.17 | 0.70 | 0.126*** | [0.007,0.248] | 0.544*** | [0.47,0.617] | 0.255*** | [0.223,0.284] |
| knm.4 | 0.287 | 4 | 31 | 210 | 56 | 0.13 | 0.93 | -0.062*** | [-0.142,0.024] | 0.366*** | [0.312,0.419] | 0.217*** | [0.197,0.237] |
| rf.4 | 0.267 | 10 | 41 | 200 | 50 | 0.17 | 0.83 | -0.003*** | [-0.126,0.124] | 0.521*** | [0.454,0.588] | 0.199*** | [0.176,0.221] |
| svm.4 | 0.068 | 9 | 81 | 160 | 51 | 0.34 | 0.85 | -0.186*** | [-0.305,-0.063] | 0.629*** | [0.566,0.691] | 0.24*** | [0.21,0.279] |

*Note*: Best performance on given dataset in bold. Stars indicate if the respective method's performance is significantly worse than the best model's performance on the same dataset (***/**/* for the 1%/5%/10% significance level). For $U_r$ and AUC higher values indicate better prediction performance, while for BPS lower values are preferable. Numbers in brackets indicate 90% confidence bands. Model names reported in format "method.dataset", where datasets consist of the following sets of variables (1): credit and asset prices, (2) macro variables, (3) external imbalances, and (4) all variables (see also table B.5).

Table C.3: Robustness: Out-of-sample performance for $\mu = 0.6$ (affects only signals / $U_r$)

| | Threshold | TP | FP | TN | FN | FPrate | FNrate | $U_r$ | |
|---|---|---|---|---|---|---|---|---|---|
| logit.1 | 0.060 | 53 | 132 | 109 | 7 | 0.55 | 0.12 | **0.277** | [0.080, 0.451] |
| trees.1 | 0.045 | 32 | 106 | 135 | 28 | 0.44 | 0.47 | -0.140*** | [-0.329, 0.063] |
| knn.1 | 0.094 | 39 | 72 | 169 | 21 | 0.30 | 0.35 | 0.176 | [-0.017, 0.367] |
| rf.1 | 0.185 | 22 | 32 | 209 | 38 | 0.13 | 0.63 | -0.083*** | [-0.254, 0.096] |
| svm.1 | 0.056 | 39 | 149 | 92 | 21 | 0.62 | 0.35 | -0.143*** | [-0.351, 0.058] |
| logit.2 | 0.069 | 37 | 175 | 66 | 23 | 0.73 | 0.38 | -0.301 | [-0.569, -0.304] |
| trees.2 | 0.076 | 8 | 178 | 63 | 52 | 0.74 | 0.87 | -1.039*** | [-1.215, -0.810] |
| knn.2 | 0.145 | 7 | 74 | 167 | 53 | 0.31 | 0.88 | -0.632** | [-0.732, -0.528] |
| rf.2 | 0.178 | 0 | 68 | 173 | 60 | 0.28 | 1.00 | -0.782*** | [-0.879, -0.679] |
| svm.2 | 0.083 | 43 | 200 | 41 | 17 | 0.83 | 0.28 | **-0.255** | [-0.513, 0.021] |
| logit.3 | 0.070 | 47 | 75 | 166 | 13 | 0.31 | 0.22 | **0.364** | [0.225, 0.480] |
| trees.3 | 0.071 | 56 | 148 | 93 | 4 | 0.61 | 0.07 | 0.286 | [0.085, 0.494] |
| knn.3 | 0.102 | 37 | 60 | 181 | 23 | 0.25 | 0.38 | 0.176** | [0.005, 0.326] |
| rf.3 | 0.320 | 13 | 23 | 218 | 47 | 0.10 | 0.78 | -0.270*** | [-0.402, -0.131] |
| svm.3 | 0.073 | 55 | 178 | 63 | 5 | 0.74 | 0.08 | 0.136** | [-0.060, 0.311] |
| logit.4 | 0.070 | 51 | 63 | 178 | 9 | 0.26 | 0.15 | **0.514** | [0.331, 0.685] |
| trees.4 | 0.049 | 20 | 61 | 180 | 40 | 0.25 | 0.67 | -0.253*** | [-0.427, -0.064] |
| knn.4 | 0.274 | 4 | 30 | 211 | 56 | 0.12 | 0.93 | -0.524*** | [-0.649, -0.394] |
| rf.4 | 0.267 | 10 | 41 | 200 | 50 | 0.17 | 0.83 | -0.420*** | [-0.600, -0.233] |
| svm.4 | 0.108 | 16 | 74 | 167 | 44 | 0.31 | 0.73 | -0.407*** | [-0.580, -0.222] |

*Note:* For $U_r$ and AUC higher values indicate better prediction performance, while for BPS lower values are preferable. Numbers in brackets indicate 90% confidence bands. For machine learning methods, we report significance levels (\*\*\*/\*\*/\* for the 1%/5%/10% level) if the respective performance measure is significantly below (for $U_r$ and AUC) or above (for BPS) the logit model on the same dataset. Model names reported in format "method.dataset", where datasets consist of the following sets of variables (1): credit and asset prices, (2) macro variables, (3) external imbalances, and (4) all variables (see also table B.5).

Table C.4: Robustness: Out-of-sample performance when using growth rates (instead of gaps) as data transformation

| | Threshold | TP | FP | TN | FN | FPrate | FNrate | U$_r$ | | AUC | | BPS | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| logit.1 | 0.092 | 39 | 74 | 179 | 34 | 0.29 | 0.47 | 0.242 | [0.133, 0.344] | **0.746** | [0.696, 0.796] | 0.169 | [0.156, 0.184] |
| trees.1 | 0.076 | 56 | 206 | 47 | 17 | 0.81 | 0.23 | -0.047*** | [-0.187, 0.090] | 0.536*** | [0.464, 0.608] | **0.163** | [0.144, 0.180] |
| knn.1 | 0.111 | 37 | 54 | 199 | 36 | 0.21 | 0.49 | **0.293** | [0.194, 0.393] | 0.638*** | [0.570, 0.707] | 0.181** | [0.168, 0.193] |
| rf.1 | 0.216 | 25 | 28 | 225 | 48 | 0.11 | 0.66 | 0.232* | [0.115, 0.296] | 0.703* | [0.639, 0.758] | 0.167 | [0.149, 0.182] |
| svm.1 | 0.099 | 30 | 85 | 168 | 43 | 0.34 | 0.59 | 0.075*** | [-0.049, 0.200] | 0.537*** | [0.496, 0.621] | 0.195* | [0.161, 0.246] |
| logit.2 | 0.095 | 10 | 130 | 123 | 63 | 0.51 | 0.86 | -0.377*** | [-0.502, -0.252] | 0.764 | [0.713, 0.815] | **0.209** | [0.198, 0.223] |
| trees.2 | 0.088 | 13 | 181 | 72 | 60 | 0.72 | 0.82 | -0.537*** | [-0.647, -0.432] | 0.807 | [0.768, 0.846] | 0.225* | [0.206, 0.247] |
| knn.2 | 0.159 | 1 | 55 | 198 | 72 | 0.22 | 0.99 | -0.204 | [-0.264, -0.140] | 0.784 | [0.742, 0.826] | 0.230*** | [0.222, 0.238] |
| rf.2 | 0.230 | 0 | 37 | 216 | 73 | 0.15 | 1.00 | **-0.146** | [-0.205, -0.097] | 0.792 | [0.749, 0.836] | 0.249*** | [0.234, 0.265] |
| svm.2 | 0.096 | 3 | 64 | 189 | 70 | 0.25 | 0.96 | -0.212 | [-0.305, -0.125] | **0.835** | [0.798, 0.874] | 0.220 | [0.187, 0.265] |
| logit.3 | 0.079 | 47 | 81 | 172 | 26 | 0.32 | 0.36 | **0.324** | [0.222, 0.419] | **0.693** | [0.628, 0.758] | **0.172** | [0.160, 0.184] |
| trees.3 | 0.070 | 49 | 194 | 59 | 24 | 0.77 | 0.33 | -0.096*** | [-0.218, 0.031] | 0.583*** | [0.508, 0.658] | 0.196** | [0.179, 0.215] |
| knn.3 | 0.096 | 48 | 96 | 157 | 25 | 0.38 | 0.34 | 0.278 | [0.169, 0.390] | 0.679 | [0.625, 0.734] | 0.180 | [0.166, 0.195] |
| rf.3 | 0.273 | 10 | 25 | 228 | 63 | 0.10 | 0.86 | 0.038*** | [-0.018, 0.115] | 0.674 | [0.610, 0.717] | 0.193*** | [0.180, 0.208] |
| svm.3 | 0.089 | 27 | 107 | 146 | 46 | 0.42 | 0.63 | -0.053*** | [-0.155, 0.083] | 0.594** | [0.543, 0.653] | 0.202* | [0.166, 0.246] |
| logit.4 | 0.092 | 51 | 54 | 199 | 22 | 0.21 | 0.30 | **0.485** | [0.373, 0.603] | **0.712** | [0.648, 0.776] | **0.163** | [0.146, 0.180] |
| trees.4 | 0.083 | 3 | 7 | 246 | 70 | 0.03 | 0.96 | 0.013*** | [-0.081, 0.119] | 0.661 | [0.609, 0.712] | 0.224*** | [0.199, 0.251] |
| knn.4 | 0.309 | 2 | 20 | 233 | 71 | 0.08 | 0.97 | -0.052*** | [-0.116, 0.023] | 0.377*** | [0.329, 0.426] | 0.230*** | [0.214, 0.246] |
| rf.4 | 0.264 | 1 | 15 | 238 | 72 | 0.06 | 0.99 | -0.046*** | [-0.150, 0.055] | 0.503*** | [0.438, 0.564] | 0.205*** | [0.188, 0.221] |
| svm.4 | 0.066 | 27 | 96 | 157 | 46 | 0.38 | 0.63 | -0.010*** | [-0.128, 0.105] | 0.483*** | [0.425, 0.557] | 0.190 | [0.155, 0.238] |

*Note*: For U$_r$ and AUC higher values indicate better prediction performance, while for BPS lower values are preferable. Numbers in brackets indicate 90% confidence bands. For machine learning methods, we report significance levels (***/**/* for the 1%/5%/10% level) if the respective performance measure is significantly below (for U$_r$ and AUC) or above (for BPS) the logit model on the same dataset. Model names reported in format "method.dataset", where datasets consist of the following sets of variables (1): credit and asset prices, (2) macro variables, (3) external imbalances, and (4) all variables (see also table B.5).

Table C.5: Robustness: Out-of-sample performance when starting the dataset in 1980Q1 (instead of 1970Q1)

| | Threshold | TP | FP | TN | FN | FPrate | FNrate | $U_r$ | | AUC | | BPS | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| logit.1 | 0.104 | 37 | 81 | 160 | 23 | 0.34 | 0.38 | **0.281** | [0.133, 0.429] | **0.732** | [0.676,0.788] | **0.139** | [0.119, 0.161] |
| trees.1 | 0.077 | 26 | 52 | 189 | 34 | 0.22 | 0.57 | 0.218 | [0.097, 0.343] | 0.559*** | [0.473,0.645] | 0.156 | [0.132, 0.180] |
| knn.1 | 0.129 | 26 | 52 | 189 | 34 | 0.22 | 0.57 | 0.218 | [0.095, 0.336] | 0.622 | [0.546,0.699] | 0.158*** | [0.141, 0.174] |
| rf.1 | 0.225 | 19 | 29 | 212 | 41 | 0.12 | 0.68 | 0.196 | [0.090, 0.302] | 0.689 | [0.627,0.751] | 0.155* | [0.135, 0.174] |
| svm.1 | 0.121 | 21 | 55 | 186 | 39 | 0.23 | 0.65 | 0.122* | [0.013, 0.263] | 0.604*** | [0.532,0.677] | 0.152 | [0.120, 0.195] |
| logit.2 | 0.100 | 5 | 83 | 158 | 55 | 0.34 | 0.92 | -0.261** | [-0.383, -0.125] | 0.765 | [0.709,0.821] | 0.189 | [0.177, 0.204] |
| trees.2 | 0.089 | 0 | 34 | 207 | 60 | 0.14 | 1.00 | -0.141 | [-0.224, -0.054] | **0.827** | [0.785,0.869] | 0.234** | [0.203, 0.266] |
| knn.2 | 0.130 | 1 | 85 | 156 | 59 | 0.35 | 0.98 | -0.336*** | [-0.431, -0.232] | 0.779 | [0.73 ,0.829] | 0.205 | [0.195, 0.216] |
| rf.2 | 0.181 | 0 | 50 | 191 | 60 | 0.21 | 1.00 | -0.207* | [-0.274, -0.128] | 0.778 | [0.734,0.821] | 0.211* | [0.199, 0.226] |
| svm.2 | 0.110 | 16 | 79 | 162 | 44 | 0.33 | 0.73 | **-0.061** | [-0.202, -0.011] | 0.729 | [0.67 ,0.788] | **0.184** | [0.159, 0.214] |
| logit.3 | 0.096 | 39 | 59 | 182 | 21 | 0.24 | 0.35 | **0.405** | [0.304, 0.500] | **0.753** | [0.691,0.815] | **0.153** | [0.137, 0.168] |
| trees.3 | 0.067 | 41 | 107 | 134 | 19 | 0.44 | 0.32 | 0.239** | [0.097, 0.396] | 0.595*** | [0.526,0.664] | 0.175** | [0.153, 0.199] |
| knn.3 | 0.147 | 19 | 51 | 190 | 41 | 0.21 | 0.68 | 0.105*** | [0.001, 0.213] | 0.555*** | [0.486,0.625] | 0.175*** | [0.158, 0.191] |
| rf.3 | 0.377 | 8 | 12 | 229 | 52 | 0.05 | 0.87 | 0.084*** | [-0.012, 0.172] | 0.587*** | [0.515,0.658] | 0.171** | [0.156, 0.187] |
| svm.3 | 0.104 | 20 | 127 | 114 | 40 | 0.53 | 0.67 | -0.194*** | [-0.301, -0.072] | 0.734 | [0.672,0.797] | 0.185** | [0.156, 0.217] |
| logit.4 | 0.090 | 37 | 27 | 214 | 23 | 0.11 | 0.38 | **0.505** | [0.370, 0.619] | 0.809 | [0.754,0.863] | **0.143** | [0.119, 0.169] |
| trees.4 | 0.070 | 1 | 51 | 190 | 59 | 0.21 | 0.98 | -0.195*** | [-0.303, -0.066] | **0.828** | [0.781,0.875] | 0.298*** | [0.268, 0.330] |
| knn.4 | 0.267 | 2 | 36 | 205 | 58 | 0.15 | 0.97 | -0.116*** | [-0.204, -0.021] | 0.645 | [0.591,0.698] | 0.211*** | [0.192, 0.230] |
| rf.4 | 0.189 | 4 | 26 | 215 | 56 | 0.11 | 0.93 | -0.041*** | [-0.175, 0.095] | 0.555 | [0.491,0.619] | 0.179*** | [0.168, 0.189] |
| svm.4 | 0.103 | 13 | 66 | 175 | 47 | 0.27 | 0.78 | -0.057*** | [-0.159, 0.071] | 0.476*** | [0.411,0.541] | 0.210*** | [0.176, 0.247] |

*Note:* For $U_r$ and AUC higher values indicate better prediction performance, while for BPS lower values are preferable. Numbers in brackets indicate 90% confidence bands. For machine learning methods, we report significance levels (***/**/* for the 1%/5%/10% level) if the respective performance measure is significantly below (for $U_r$ and AUC) or above (for BPS) the logit model on the same dataset. Model names reported in format "method.dataset", where datasets consist of the following sets of variables (1): credit and asset prices, (2) macro variables, (3) external imbalances, and (4) all variables (see also table B.5).

Table C.6: Robustness: Out-of-sample performance when using Laeven & Valencia crisis database

| | Threshold | TP | FP | TN | FN | FPrate | FNrate | $U_r$ | | AUC | | BPS | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| logit.1 | 0.029 | 63 | 26 | 16 | 26 | 0.62 | 0.29 | 0.089*** | [-0.052, 0.216] | 0.616 | [0.527, 0.705] | 0.567** | [0.524, 0.611] |
| trees.1 | 0.027 | 17 | 15 | 27 | 72 | 0.36 | 0.81 | -0.166*** | [-0.340, 0.004] | 0.585 | [0.499, 0.672] | 0.556 | [0.512, 0.596] |
| knn.1 | 0.054 | 42 | 7 | 35 | 47 | 0.17 | 0.53 | **0.305** | [0.178, 0.437] | 0.613 | [0.532, 0.694] | 0.588*** | [0.546, 0.626] |
| rf.1 | 0.217 | 19 | 1 | 41 | 70 | 0.02 | 0.79 | 0.190** | [0.113, 0.294] | **0.617** | [0.537, 0.696] | 0.557** | [0.516, 0.596] |
| svm.1 | 0.053 | 29 | 5 | 37 | 60 | 0.12 | 0.67 | 0.207* | [0.092, 0.322] | 0.464** | [0.375, 0.554] | **0.518** | [0.471, 0.559] |
| logit.2 | 0.034 | 0 | 6 | 36 | 89 | 0.14 | 1.00 | -0.143 | [-0.286, 0.015] | 0.606 | [0.497, 0.714] | 0.649 | [0.625, 0.666] |
| trees.2 | 0.031 | 77 | 40 | 2 | 12 | 0.95 | 0.13 | -0.087 | [-0.250, 0.079] | 0.668 | [0.574, 0.761] | 0.647 | [0.630, 0.664] |
| knn.2 | 0.077 | 3 | 12 | 30 | 86 | 0.29 | 0.97 | -0.252 | [-0.345, -0.164] | 0.693 | [0.614, 0.773] | 0.667** | [0.654, 0.677] |
| rf.2 | 0.246 | 0 | 5 | 37 | 89 | 0.12 | 1.00 | -0.119 | [-0.201, -0.045] | 0.690 | [0.602, 0.778] | 0.678** | [0.669, 0.686] |
| svm.2 | 0.032 | 37 | 21 | 21 | 52 | 0.50 | 0.58 | **-0.084** | [-0.335, 0.043] | **0.723** | [0.641, 0.806] | **0.640** | [0.610, 0.666] |
| logit.3 | 0.034 | 37 | 0 | 42 | 52 | 0.00 | 0.58 | **0.416** | [0.292, 0.526] | **0.808** | [0.742, 0.873] | **0.628** | [0.610, 0.644] |
| trees.3 | 0.030 | 89 | 42 | 0 | 0 | 1.00 | 0.00 | 0.000*** | [-0.158, 0.140] | 0.630** | [0.531, 0.730] | 0.642 | [0.617, 0.665] |
| knn.3 | 0.046 | 28 | 11 | 31 | 61 | 0.26 | 0.69 | 0.053*** | [-0.065, 0.168] | 0.554*** | [0.471, 0.636] | 0.632 | [0.612, 0.650] |
| rf.3 | 0.129 | 7 | 0 | 42 | 82 | 0.00 | 0.92 | 0.079*** | [0.014, 0.178] | 0.520*** | [0.433, 0.606] | 0.643 | [0.619, 0.663] |
| svm.3 | 0.036 | 26 | 18 | 24 | 63 | 0.43 | 0.71 | -0.136*** | [-0.260, 0.112] | 0.558*** | [0.459, 0.657] | 0.639 | [0.617, 0.663] |
| logit.4 | 0.035 | 49 | 1 | 41 | 40 | 0.02 | 0.45 | **0.527** | [0.380, 0.663] | **0.770** | [0.703, 0.836] | 0.580 | [0.521, 0.632] |
| trees.4 | 0.030 | 6 | 18 | 24 | 83 | 0.43 | 0.93 | -0.361*** | [-0.523, -0.212] | 0.711 | [0.631, 0.791] | 0.640** | [0.593, 0.681] |
| knn.4 | 0.221 | 6 | 9 | 33 | 83 | 0.21 | 0.93 | -0.147*** | [-0.238, -0.052] | 0.430*** | [0.361, 0.499] | 0.656*** | [0.621, 0.690] |
| rf.4 | 0.133 | 14 | 4 | 38 | 75 | 0.10 | 0.84 | 0.062*** | [-0.077, 0.148] | 0.629** | [0.545, 0.712] | 0.624** | [0.601, 0.644] |
| svm.4 | 0.038 | 29 | 13 | 29 | 60 | 0.31 | 0.67 | 0.016*** | [-0.105, 0.133] | 0.526*** | [0.435, 0.618] | **0.578** | [0.536, 0.616] |

*Note*: For $U_r$ and AUC higher values indicate better prediction performance, while for BPS lower values are preferable. Numbers in brackets indicate 90% confidence bands. For machine learning methods, we report significance levels (***/**/* for the 1%/5%/10% level) if the respective performance measure is significantly below (for $U_r$ and AUC) or above (for BPS) the logit model on the same dataset. Model names reported in format "method.dataset", where datasets consist of the following sets of variables (1): credit and asset prices, (2) macro variables, (3) external imbalances, and (4) all variables (see also table B.5).

Figure C.1: Robustness (preference parameter): Relative usefulness of in- and out-of-sample estimation by model.
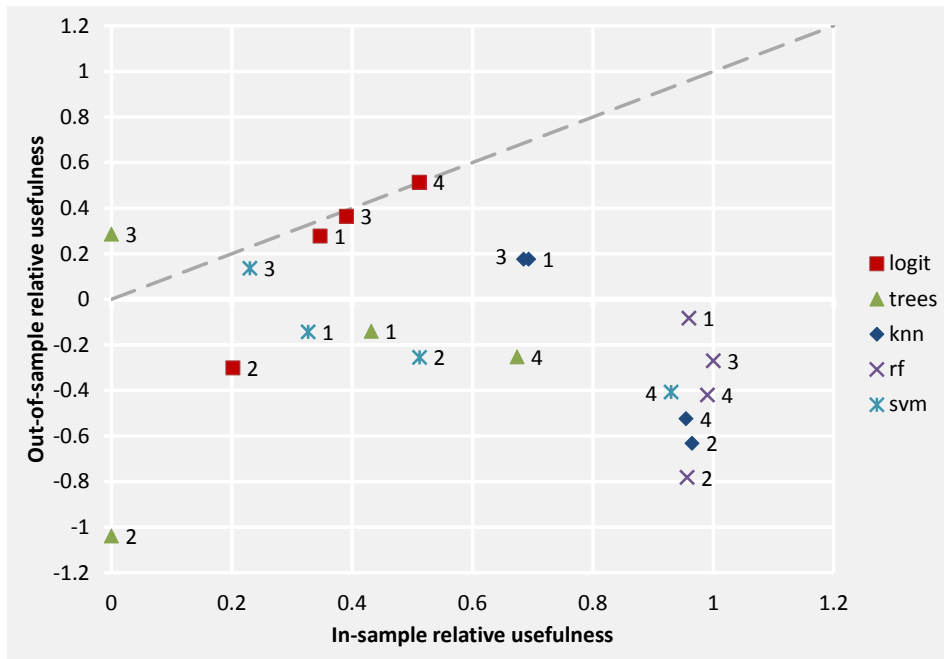


Figure C.2: Robustness (data transformation): Relative usefulness of in- and out-of-sample estimation by model.
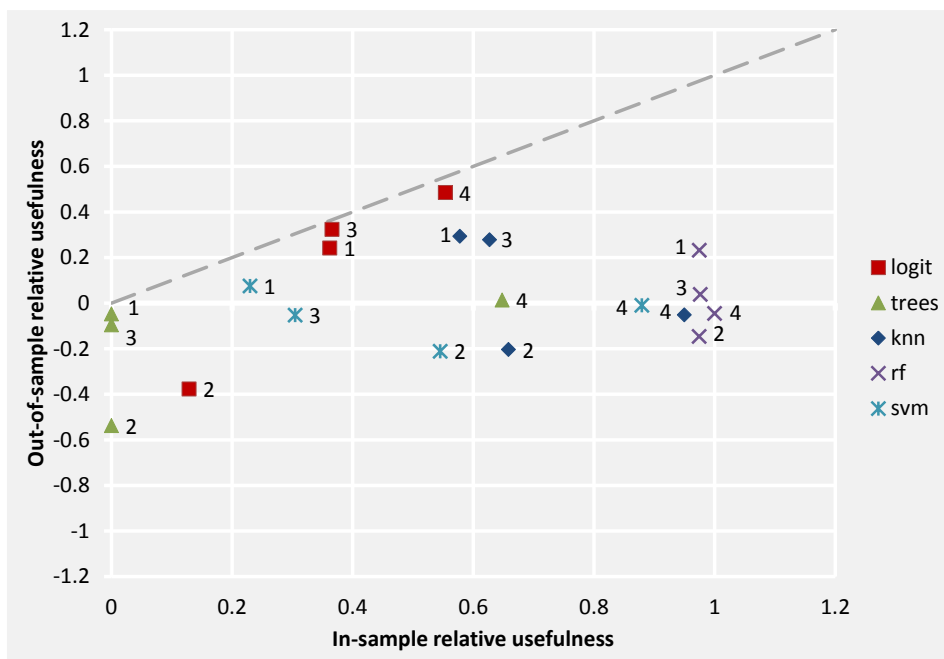
Figure C.3: Robustness (sample length): Relative usefulness of in- and out-of-sample estimation by model.
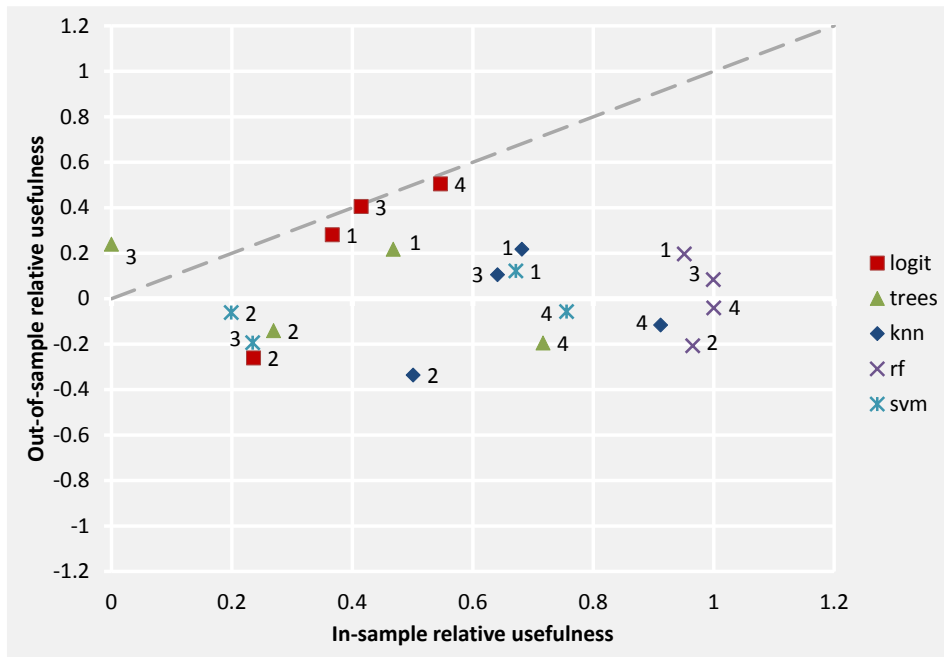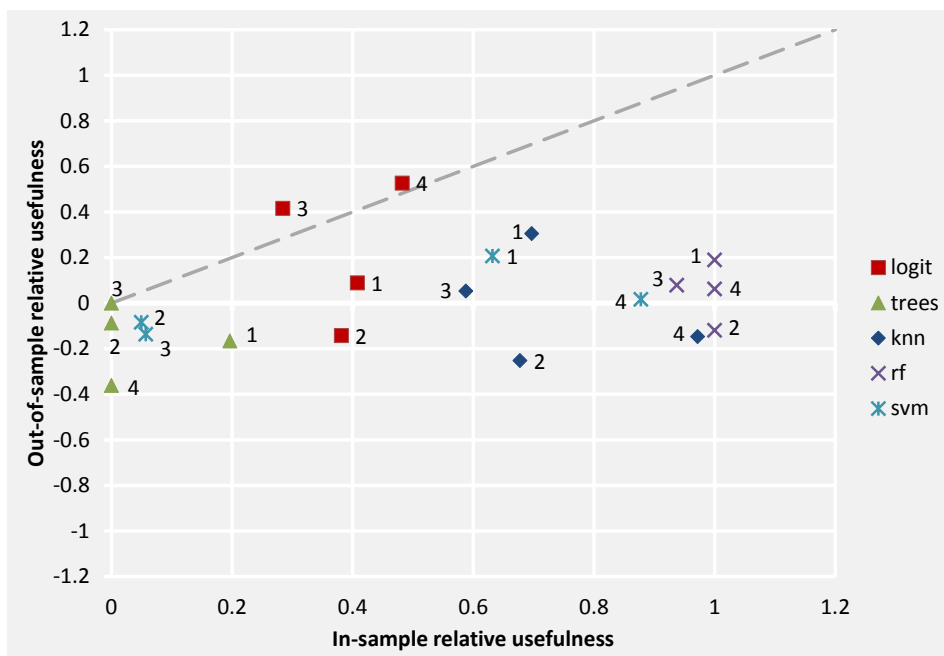


Figure C.4: Robustness (crisis database): Relative usefulness of in- and out-of-sample estimation by model.

Leibniz
Association