

Bank Information Production Over the Business Cycle*

Cooper Howes[†] Gregory Weitzner[‡]

July 2025

Abstract

The information banks produce drives their lending decisions and macroeconomic outcomes, but this information is inherently difficult to analyze because it is private. We construct a novel measure of bank information quality from confidential regulatory data that include banks' private risk assessments for US corporate loans. Information quality improves as local economic conditions deteriorate, particularly for new loans, large loans, and loans with higher expected losses. Information quality also declines during periods of rapid local house price appreciation. Our results provide empirical support for theories of countercyclical information production in credit markets.

*We thank Thomas Chaney (the editor) and five anonymous referees as well as Hassan Afrouzi, Cynthia Balloch, Javier Bianchi, Olivier Coibion, Mariela Dal Borgo, Adolfo De Motta, Miguel Faria-e-Castro, Daniel Greenwald, Stefan Jacewitz, Gustavo Joaquim, Artashes Karapetyan, Anya Kleymenova, Alexandre Kohlhas, Yueran Ma, Blake Marsh, Johannes Matschke, Ralf Meisenzahl, Karel Mertens, Atanas Mihov, Lars Norden, Guillermo Ordoñez, Pablo Ottonello, Matt Pritsker, Samuel Rosen, Kasper Roszbach, Guillaume Roussellet, Jane Ryngaert, Padma Sharma, Lee Smith, Wenting Song, Andrea Vedolin, Jeff Wooldridge, Yufeng Wu and Choongryul Yang as well as seminar and conference participants at Arizona State University (Economics), the Bank of Canada, Bentley University (Economics), BSE Summer Forum, EFA, Federal Reserve Bank of Kansas City, Federal Reserve Board, IBEFA Annual Meeting, McGill University (Finance), Norges Bank, UNC Greensboro (Economics), University of Notre Dame (Economics), University of Oklahoma (Economics), Oxford Saïd – Risk Center at ETH Zürich Macro-finance Conference, George Washington University (Finance), UVA Darden (Finance), CAFRAL, CICF, FDIC Bank Research Conference, the Finance Forum, Fixed Income and Financial Institutions Conference, MFA, NFA, SEA annual meetings, NAMES summer meetings and the Procyclicality Symposium for helpful comments and discussions. We also thank Alex Zhang for excellent research assistance. These views are those of the authors and do not reflect the views of the Federal Reserve Board of Governors or the Federal Reserve System.

[†]Federal Reserve Board of Governors. Email: cooper.a.howes@frb.gov.

[‡]McGill University. Email: gregory.weitzner@mcgill.ca.

1 Introduction

A fundamental role of banks is to produce information about prospective borrowers.¹ Banks use this information to determine the recipients and terms of financing; hence, their information production decisions can affect real economic activity and financial stability through the supply of credit to firms. If the returns to distinguishing between different types of borrowers change with economic conditions, banks’ incentives to produce information can affect and be affected by business cycles. Despite policymaker interest and an extensive theoretical literature emphasizing the importance of banks’ information, there is little evidence of its empirical properties.

The key empirical challenge to testing theories of bank information production is that banks’ information is intrinsically private and, therefore, unobservable to the econometrician. Because of this data limitation, the existing literature typically relies on indirect evidence; however, without access to banks’ private information, researchers are severely constrained in their ability to test these theories. In this paper, we address this challenge using confidential regulatory data that contain banks’ private risk assessments for corporate bank loans over one million dollars in the US. We use county-level variation in unemployment rates to show that banks’ risk assessments discriminate better across borrowers in local downturns, suggesting that banks’ information quality is countercyclical. We then provide evidence that this countercyclical information quality results from endogenous information production. Specifically, we find that the cyclical variation in banks’ information quality is concentrated in loans that theory predicts to be more information sensitive: new loans, larger loans, and loans with higher expected losses. Finally, consistent with higher collateral values reducing information production incentives, we also show that information quality is lower during periods of rapid local house price appreciation. Overall, our results provide empirical support for theories of countercyclical information production in credit markets.

Our analysis uses the Federal Reserve’s Y-14Q Schedule H.1 data, which include all corporate loans larger than one million dollars extended by large bank holding companies (BHC). In addition to detailed loan and borrower characteristics, qualified BHCs must report their internal estimate of the borrower’s probability of default (PD) for each loan. Because the data also reveal whether loans ultimately default, these PDs—which incorporate both “hard” and “soft” information—allow us to quantify bank information quality.

We first show, using linear regressions and random forest regressions that allow for nonlinearities and interactions, that banks’ PDs predict realized default even after controlling for a rich set of loan- and firm-level controls. These results suggest that banks’ risk assessments contain private information that is i) relevant for predicting default and ii) not captured by other observables.

¹E.g., [Leland and Pyle \(1977\)](#), [Diamond \(1984\)](#) and [Boyd and Prescott \(1986\)](#).

In many models of information production in credit markets, including the simple one developed in this paper, banks have stronger incentives to produce idiosyncratic information during downturns, enabling them to better distinguish borrower quality.² To test this prediction, we estimate the area under the receiver operating characteristics curve (AUC) based on banks' reported PDs and the subsequent default realizations of those loans. The AUC, described in much further detail below, measures the discriminatory ability of forecasts of binary outcomes and is the most commonly used approach by practitioners ([Engelmann and Rauhmeier \(2011\)](#)), bank regulators ([Basel Committee on Banking Supervision \(2005\)](#)), and academics ([Puri, Rocholl, and Steffen \(2017\)](#), [Berg and Koziol \(2017\)](#) and [Berg, Puri, and Rocholl \(2020\)](#)). A higher AUC implies that banks' PDs discriminate better across borrowers.

We then analyze how the AUC, our measure of information quality, evolves over the business cycle. Specifically, we split our sample into periods of high and low unemployment based on whether a county's unemployment rate was above or below its median across our sample period. We find that the AUC derived from banks' PDs for newly originated loans is higher in periods of high unemployment and that this difference is statistically significant based on a DeLong test ([DeLong, DeLong, and Clarke-Pearson \(1988\)](#)).

While our main result is consistent with banks producing more information when economic conditions weaken, bank information quality could also vary exogenously over the business cycle. For example, if more firms become delinquent in periods of high unemployment, banks may be able to better distinguish between borrower types from the information they receive exogenously from borrowers. Next, we conduct tests that provide support for the endogenous information production channel by analyzing how banks' information quality varies across loans based on their information sensitivity ([Dang, Gorton, and Holmström \(2013\)](#)), i.e., the value of information for a given loan.

First, banks' information production incentives should be more sensitive to the business cycle for new loans because they require risking additional capital, in contrast to existing loans, for which banks' capital has already been sunk. To test this hypothesis, we separately compare the AUC of high- and low-unemployment periods for new and existing loans. Consistent with this prediction, we find that banks' information quality increases more during periods of high unemployment for newly originated loans. When we expand our sample to include existing loans, we also show that economic conditions at origination have persistent effects on information quality, even more so than current economic conditions. This result is consistent with the theories that motivate our analysis, in which economic conditions drive banks' information production decisions when deciding whether to grant a loan.

Second, we test whether banks' information quality is higher for larger loans and loans with higher expected losses. According to several theories of information production in credit

²See [Ruckes \(2004\)](#), [Dell'Ariccia and Marquez \(2006\)](#), [Gorton and He \(2008\)](#), [Dang, Gorton, and Holmström \(2013\)](#), [Gorton and Ordonez \(2014\)](#), [Gorton and Ordonez \(2020\)](#), [Fishman, Parker, and Straub \(2020\)](#), [Petriconi \(2015\)](#), [Farboodi and Kondor \(2020\)](#) and [Asriyan, Laeven, and Martin \(2022\)](#).

markets, such as [Dang, Gorton, and Holmström \(2012\)](#) and [Gorton and Ordonez \(2014\)](#), banks should produce more information about these loans, as they will have higher returns to distinguishing between borrowers. Consistent with these predictions, when we split loans by their size or expected loss and re-estimate each AUC separately, we find that information quality increases with loan size and expected losses. Moreover, for both characteristics, we show that the difference in AUCs between high- and low-unemployment periods is larger for the top quartile than for the bottom quartile, suggesting that information quality is more cyclically sensitive for large loans and loans with higher expected losses.

A key determinant of expected losses is collateral values, which recent theories suggest play an important role in how banks' information production incentives evolve over the business cycle. For example, rapid increases in collateral values can lead to reduced incentives for banks to screen borrowers as expected losses decrease (e.g., [Gorton and Ordonez \(2020\)](#), [Asriyan, Laeven, and Martin \(2022\)](#)). Motivated by these theories, we first show that local housing prices—a commonly used proxy for collateral values in the literature—are associated with lower LGDs (loss given default) and expected losses. We then show that the AUC is lower in areas with high house price growth, suggesting that increasing collateral values dampen banks' information production incentives.

Literature review. While the theory literature has long recognized the importance of bank information production (e.g., [Leland and Pyle \(1977\)](#), [Diamond \(1984\)](#) and [Boyd and Prescott \(1986\)](#)), testing for it is notoriously difficult given the private nature of banks' information. For this reason, the existing empirical literature focuses on proxies or indirect evidence of information production (e.g., [James \(1987\)](#), [Cerqueiro, Ongena, and Roszbach \(2016\)](#), [Gustafson, Ivanov, and Meisenzahl \(2020\)](#), [Iyer et al. \(2016\)](#) and [Bedayo et al. \(2020\)](#)). However, there are many different interpretations of these proxies. For example, [Gustafson, Ivanov, and Meisenzahl \(2020\)](#) create a measure of monitoring based on the number of visits banks take to firms. However, it is unclear whether banks visit firms to collect private information, in response to receiving information (public or private), or both. In contrast, our data allow us to observe banks' private information directly and test how well this information predicts subsequent defaults.

The closest to our empirical approach is [Becker, Bos, and Roszbach \(2020\)](#). They also find that banks' internal credit ratings better predict default in bad times, but they attribute this to exogenous variation in information over the business cycle. There are three key differences in our analysis and the interpretation of our results. First, their data are at the firm level rather than the loan level. This difference allows us to explore the relationship between loan characteristics and information production and how this relationship changes over the business cycle. Second, we provide evidence that the countercyclicality of information quality is driven by endogenous bank information production by showing that the effects are stronger for new loans, large loans, and loans with higher expected losses, which is difficult to rationalize solely through exogenous variation in information quality over the business cycle. Third, their approach uses time series variation in country-wide aggregate economic conditions for a single Swedish bank, while we

exploit rich cross-sectional variation in economic conditions across US counties.

Our paper also relates to the theoretical work analyzing the cyclicity of information production in credit markets. This includes many theories in which information production is countercyclical (e.g., [Ruckes \(2004\)](#), [Dell’Ariccia and Marquez \(2006\)](#), [Gorton and He \(2008\)](#), [Dang, Gorton, and Holmström \(2013\)](#), [Gorton and Ordonez \(2020\)](#), [Fishman, Parker, and Straub \(2020\)](#), [Petriconi \(2015\)](#), [Farboodi and Kondor \(2020\)](#) and [Asriyan, Laeven, and Martin \(2022\)](#)). A common feature in these models is that lending standards tighten in downturns as banks produce more information. While countercyclical lending standards are widely acknowledged empirically (e.g., [Asea and Blomberg \(1998\)](#), [Lown and Morgan \(2006\)](#), [Maddaloni and Peydró \(2011\)](#), [Dell’Ariccia, Igan, and Laeven \(2012\)](#), [Bassett et al. \(2014\)](#) and [Rodano, Serrano-Velarde, and Tarantino \(2018\)](#)), the mechanism behind them remains unclear given the unobservability of banks’ information. For example, countercyclical lending standards could arise purely from banks imposing stricter thresholds for lending. Conversely and consistent with our evidence, banks could lend more selectively exactly because they produce more information. Hence, our results speak to the mechanisms behind changes in lending standards over the cycle.

2 Data

Our main data source is Schedule H.1 of the Federal Reserve’s Y-14Q filings. The Federal Reserve began collecting these data to support the Dodd-Frank mandated stress tests and the Comprehensive Capital and Analysis Review (CCAR). The sample includes commercial and industrial (C&I) loans from bank holding companies (BHCs) with \$50bn or more in total assets³, accounting for 85.9% of all assets in the banking sector ([Frame, McLemore, and Mihov \(2020\)](#)). Qualified institutions are required to report detailed quarterly loan-level data on corporate loans of at least \$1mm. The universe of loans we analyze is large: [Bidder, Krainer, and Shapiro \(2020\)](#) show that the Y-14Q data cover 70% of all commercial and industrial loan volume extended by BHCs that file an FR Y-9C report.

The data include detailed loan characteristics (interest rates, maturity, amount, collateral, and purpose) and performance measures (defaults, past-due payments, non-accruals, and charge-offs). They also include income statement, balance sheet, and geographic information about borrowers. Crucially, banks must also report their internal estimates of the borrower’s probability of default (PD) and loss given default (LGD) for each loan. According to the Basel Committee on Banking Supervision, internal estimates of PD and LGD “must incorporate all relevant, material and available data, information and methods. A bank may utilize internal data and data from external sources (including pooled data).”⁴

³In 2019, this threshold was increased to \$100bn. The most recent list of participating institutions can be found in Table 3 of the [2024 Federal Reserve Stress Test Results](#).

⁴The most recent instructions are available at [Calculation of RWA for credit risk](#).

Our primary analysis includes only newly originated loans to study banks' information production incentives at the time financing is committed; however, we also consider several extensions including existing loans. We exclude demand loans, which can be recalled by the lender at any time, as well as loans with government guarantees⁵, tax-exempt loans, loans to foreign borrowers, and loans to firms in the finance, insurance, and real estate (FIRE) sectors. We drop loans with negative interest rates and those with missing company identifiers, PD, or loan amount at origination. We follow [Brown, Gustafson, and Ivanov \(2021\)](#) and exclude loans to companies with under \$100k in reported assets at origination; given that the minimum reporting threshold for loans is \$1 million, these observations are likely reporting errors. Finally, we drop all publicly traded firms and private firms with assets above the 99th percentile, as these firms are likely to be more geographically diverse and, thus, less sensitive to changes in local economic conditions.

We define the following firm-level financial variables: profitability (EBITDA/assets), size (log assets), tangibility (tangible assets/assets), and leverage (debt/assets), which we winsorize at the 1% and 99% level. Our primary measure of loan performance is default, a dummy variable that equals one if the borrower defaults within two years after origination. Focusing on a two-year default window strikes a balance between our data's limited time series and the fact that the median loan maturity is close to five years. Our sample starts in 2014Q4 when the PD variable first becomes well populated. To allow for consistent measurement of default rates, we include loans on banks' balance sheets until 2021Q4 and track whether they ultimately default through 2023Q4.

Table 1 includes firm, loan, and county summary statistics. Panel A shows summary statistics at the loan level for newly originated loans, where the average and median loan size is approximately \$13.3mm and \$3.6mm, respectively. We average each reported measure at the firm-quarter level across all outstanding loans to calculate the firm-level statistics in Panel B. The median firm has \$21.5mm in assets and a leverage ratio of 0.26. These loan and firm sizes are small relative to other sources of loan data, such as DealScan, because our sample contains many small, private firms. Over our sample period, 1.13% of loans default within two years after origination, compared to an average ex-ante expected PD of 1.63%. Panel C shows characteristics aggregated at the county level. The median number of loans outstanding for each county-quarter is 5, while the median new loan volume is about \$40mm. Finally, Table 2 includes additional loan-level summary statistics, including splits by high and low-unemployment periods.

The first three panels of Figure 1 show the distributions of PD and log(PD). If PD contains valuable information for predicting default, then there should be a positive correlation between PD and future realized default. In the bottom-right panel, we place loans into PD quintiles where the number below each column indicates the average PD in that quintile of loans, while

⁵We conduct a test in Online Appendix Section C.4 using loans with government guarantees, which provides further additional support for the endogenous information production mechanism.

the vertical axis indicates the average realized default rate. The figure shows a clear positive relationship between PDs and realized defaults, suggesting that PD contains valuable information regarding the borrower’s default risk. We test this relationship more formally in the next section.

3 PDs Predict Realized Default

In this section, we validate PD as a measure of banks’ private information by showing that it is a statistically and economically significant predictor of realized default, even after controlling for various loan and firm characteristics. We start by estimating the following regression:

$$Default_i = \beta PD_i + \Omega X_i + \delta_{b,t} + \gamma_{j,t} + \sigma_{b,c} + \epsilon_i, \quad (1)$$

where i , b , t , j and c , index loan, bank, quarter, industry, and county, respectively. $Default_i$ is a dummy variable that equals one if loan i defaults within eight quarters following origination. PD_i is the bank’s estimated probability of default described in Section 2. X_i is a vector of firm and loan characteristics which include firm size (log of total assets), leverage ratio (total debt to total assets), profitability ratio (EBITDA to total assets), tangibility ratio (tangible assets to total assets), log loan size, the log of the original loan maturity in months, and the bank’s estimate of loss given default per dollar of exposure (LGD), as well as loan type fixed effects. We include bank-quarter fixed effects ($\delta_{b,t}$) to absorb any differences in banks’ risk assessment models and cost of capital, industry-quarter fixed effects ($\gamma_{j,t}$) to absorb variation in average loan performance across industries, and bank-county fixed effects ($\sigma_{b,c}$) to absorb persistent differences in risk assessment models or credit analysts across counties. In all regressions, we cluster standard errors by county.

The results are shown in Table 3. The primary coefficient of interest is β , which represents the expected increase in realized default (measured in percentage points) from a one percentage point increase in a loan’s PD. In Column (1), the coefficient estimate is 0.407, which means that an increase in PD of 1pp increases the probability of realized default by about 41bps.⁶ In Column (2), we display the results with firm and loan characteristics and find a similar coefficient of 0.444. Under strict rational expectations, regressing a realized outcome on its forecast should yield a coefficient of one (Muth (1961)); however, empirically, this is often not the case (e.g., Mincer and Zarnowitz (1969)). We discuss potential reasons why the coefficient estimates may deviate from one in Online Appendix Section C.5.

For comparison, Columns (3) and (4) repeat the same exercise using interest rates as an alternative measure of borrower risk.⁷ As expected, loans with higher interest rates default more

⁶Online Appendix Table OA.4 shows that these results are robust to alternative measures of loan performance.

⁷The number of observations drops when we include interest rates in the regression because we drop undrawn credit lines, for which banks are instructed to record an interest rate of zero.

frequently. However, when we include both interest rates and PD in the same specification in Column (5), the coefficient for interest rate attenuates substantially, while the coefficient for PD remains basically unchanged from Column (2).

One shortcoming of the regression approach in (1) is that it imposes a linear relationship between default and firm/loan characteristics. If firms’ actual PDs reflect nonlinearities or interactions between these characteristics, the linear specification in (1) could understate the explanatory power of observable characteristics. This would lead us to incorrectly attribute the predictive power of PD to private information when, in fact, banks’ PDs would simply be capturing public information in a more sophisticated way.⁸ To address this concern, we produce predicted default estimates using the random forest regression algorithm developed in Breiman (2001). This approach, described in detail in Appendix A, generates predictions by sequentially partitioning observations based on their observable characteristics and calculating the average default rate within each partition. The ability to flexibly accommodate complex nonlinearities and high-dimensional data has made random forests a popular tool across a range of fields, including econometrics (Wager and Athey (2018)), asset pricing (Gu, Kelly, and Xiu (2020)), and macroeconomics (Goulet Coulombe et al. (2022)).

We first estimate predicted default using several different specifications that combine the same firm and loan characteristics in (1) with combinations of dummy variables for industry, bank, and time. In each case, we first estimate the random forest using half of our new loans. We then estimate the following regression using the predictions PD^{RF} generated for the remaining half of our sample:

$$Default_i = \beta PD_i + \theta PD_i^{RF} + \epsilon_i.$$

The results are reported in Table 4. Across all specifications, both β and θ remain statistically and economically significant. In addition to providing out-of-sample validation that the random forest algorithm produces reasonable default estimates, this result shows that banks’ PDs contain information useful for predicting default that is not fully reflected in observables and validate PD as a meaningful measure of banks’ private information regarding firms’ default risk.⁹

While we refer to all loan and firm characteristics as “observables”, in reality, banks only observe a subset of these (i.e., from the firms they specifically lend to). Hence, these empirical tests, which control for all these characteristics across the entire sample, likely underestimate

⁸In principle, collinearity between PD and other observables could also lead to potential difficulties in accurately assessing the private information content of banks’ PDs in this framework. However, we show in Online Appendix Table OA.2 that PD is only weakly correlated with other firm and loan characteristics, and in Online Appendix Table OA.3 that statistical tests soundly reject collinearity.

⁹In Online Appendix Tables OA.5 and OA.6, we perform similar exercises using interest rates to show that PDs contain valuable information for pricing loans, which is consistent with Beyhaghi, Fracassi, and Weitzner (2025). In Appendix A, we show that PDs provide incremental *discriminatory power* over observables using the area under the receiver operating characteristic curve approach, which we describe in the next section.

banks’ private information. That we still find banks’ actual PDs predict future defaults, even with these stringent controls, strengthens our conclusion that PDs capture meaningful private information.

4 Bank Information Quality Over the Business Cycle

In the previous section, we show that PDs contain information regarding future default that is not reflected in observable characteristics, suggesting that PDs reflect banks’ private information. In this section, we analyze how the quality of banks’ information evolves over the business cycle. We present our approach to assessing information quality in Section 4.1 and then implement this approach in Section 4.2.

4.1 Approach to Measuring Information Quality

In the theories that motivate our analysis, banks have increased incentives to produce *idiosyncratic* information about borrowers during bad times (e.g., [Ruckes \(2004\)](#) [Dang, Gorton, and Holmström \(2012\)](#), [Gorton and Ordonez \(2014\)](#), [Asriyan, Laeven, and Martin \(2022\)](#), [Fishman, Parker, and Straub \(2020\)](#), [Farboodi and Kondor \(2020\)](#)).¹⁰ Thus, we would like to have a statistical measure of how well PDs *discriminate across* borrowers. Henceforth, we refer to information quality and discrimination interchangeably.

The main way that banks, practitioners, and regulators measure the discriminatory power of PDs is by estimating the area under the receiver operating characteristics (ROC) curve (AUC).¹¹ At a high level, the AUC measures how well banks’ internal risk assessments rank borrowers based on their relative default risk. Estimating the AUC first requires constructing the ROC curve. The ROC curve considers every possible value of PD in the data as a threshold for classifying loans: any loan with a PD above this threshold is classified as a predicted default, while any loan with a PD below this threshold is classified as a predicted non-default. For each threshold PD, the true positive rate is defined as the ratio of correctly predicted defaults over the total number of defaults, whereas the false positive rate is the ratio of incorrectly predicted defaults over the total number of non-defaults. Using each observed value of PD as a classification threshold, the ROC curve plots the false positive rate on the x-axis and the true positive rate on the y-axis. The AUC is the area under the ROC curve.

¹⁰This can be due to several reasons. For example, in recessions, the difference in payoff between lending to a “good” and “bad” borrower may increase (e.g., [Dang, Gorton, and Holmström \(2012\)](#)), or there may be more “bad” borrowers (e.g., [Farboodi and Kondor \(2020\)](#)).

¹¹For example, in referring to the receiver operating characteristics curve and the cumulative accuracy curve (whose test statistic is a simple transformation of the AUC), [Engelmann and Rauhmeier \(2011\)](#) say they are “...the most important and the most widely applied in practice.” The Basel Committee on Banking Supervision stated, “The Group has found that the Accuracy Ratio (AR) and the ROC measure appear to be more meaningful than the other above-mentioned indices because of their statistical properties” ([Basel Committee on Banking Supervision \(2005\)](#)). The AUC is also highly prevalent in other contexts measuring the performance of binary classifier models, particularly in medicine (e.g., [Pepe \(2003\)](#)) and machine learning (e.g., [Bradley \(1997\)](#)).

The AUC also has a simple probabilistic interpretation: given a randomly chosen defaulting loan and solvent loan, the AUC is the probability that the defaulting loan’s PD is higher than the solvent one’s. Hence, a higher AUC indicates that banks’ PDs have higher discriminatory power. A completely random prediction model would have an AUC of 0.5, while a perfect prediction model would have an AUC of 1. As a rule of thumb, an AUC of 0.6 is generally considered desirable in environments with less information, whereas AUCs of 0.7 or greater are desirable in information-rich environments (Iyer et al. (2016) and Berg, Puri, and Rocholl (2020)). We use the Stata function `roccomp` to construct ROC curves and numerically integrate them to estimate the AUC. We test for statistical significance of differences in AUCs using the DeLong test (DeLong, DeLong, and Clarke-Pearson (1988)), which is the standard approach for testing differences in AUCs.

The AUC is particularly well suited for measuring bank information quality for several reasons. First, because it depends only on the relative ordering of PDs, it is unaffected by any transformation of PDs that does not affect the relative ordering of loans.¹² Second, the AUC primarily measures the type of information quality we are focused on, i.e., PDs’ discrimination ability. This contrasts with other measures of forecast accuracy. For example, the sum of mean-squared forecast errors, also known as the Brier Score (Brier (1950)) when applied to binary outcomes, also reflects other factors unrelated to discrimination (Murphy (1973)). In Online Appendix Section C.8, we show that there is a strong mechanical relationship between the Brier Score and the underlying risk of borrowers, particularly when PDs are close to zero. Moreover, a well-known problem with forecast errors is that they do not distinguish forecast accuracy well for rare events. We discuss these issues and provide numerical examples in Online Appendix Section C.8; however, based on our understanding, these are the main reasons forecast errors are not as commonly used in practice for credit scoring, particularly when measuring the discrimination ability of PD models. In contrast, the AUC is far less affected by the underlying level of borrowers and is better able to distinguish rare events.¹³

Figure 2 displays the estimated AUC of 0.703 over our sample of new loans. While the thresholds for interpreting AUCs will vary depending on the context, the average AUC in our sample is slightly larger than recent studies analyzing consumer loans in a large German bank (Puri, Rocholl, and Steffen (2017), Berg and Koziol (2017), and Berg, Puri, and Rocholl (2020)).¹⁴ As further validation of the AUC as a measure of information quality, in Appendix B we develop a simple model in which banks produce more information in bad times, resulting in a higher AUC.

¹²For example, banks may have incorrect priors, incentives to misreport PDs (e.g., Behn, Haselmann, and Vig (2016) and Plosser and Santos (2018)), or behavioral biases that cause them to underreact to information. See Section C.5 for a further discussion of these issues.

¹³See Appendix Section B.4 and Online Appendix Section C.8 for a formal analysis of these issues. Another common approach is the accuracy ratio (AR) obtained from the cumulative accuracy profile (CAP) curve (Engelmann and Rauhmeier (2011)). However, the AR is a simple linear transformation of the AUC: $AR = 2AUC - 1$. We use the AUC because it has a simpler probabilistic interpretation.

¹⁴See also Lessmanna et al. (2013) and Hayashi (2022) for surveys of AUCs in different credit scoring contexts.

4.2 Testing the Cyclicalities of Bank Information Quality

In this section, we test the cyclicalities of bank information quality by analyzing how the AUC varies over local economic conditions. Our measure of county-level economic conditions is the unemployment rate from the BLS.¹⁵ Figure 3 highlights the substantial cross-sectional variation in the changes in county-level unemployment rates over this period, with roughly one-quarter of counties experiencing an increase. The top-right panel displays a histogram of defaults across county unemployment rates, suggesting that the variation in defaults for new loans is not coming solely from high-unemployment areas.

In the top-left panel of Figure 4, we split our primary sample of new loans based on whether the county's unemployment rate is above or below its county-specific median over the sample. The AUC in periods of high unemployment is 0.724 versus 0.681 in periods of low unemployment. This difference in AUCs implies that, given a randomly chosen defaulting loan and non-defaulting loan, the probability that the defaulting loan's PD is higher than the solvent one is 4.3pp higher in periods of high unemployment. Hence, this result suggests that PDs better discriminate across borrowers in bad times. In the other three panels, we find qualitatively similar results if we instead define a high-unemployment period as one in which 1) the county's unemployment rate is higher than the median county-level unemployment rate across counties with at least five new loans in a given quarter (top right), 2) the county's unemployment rate increased from the previous quarter (bottom left), or 3) the quarterly change in a county's unemployment rate was greater than the quarterly change in the aggregate US unemployment rate (bottom right).¹⁶

Local economic conditions should be more relevant for loans whose cash flows are more sensitive to local economic conditions. We test this hypothesis by comparing firms in tradeable and nontradeable industries.¹⁷ Because firms in nontradeable industries will be more dependent on local markets, the same change in local economic conditions should have a more significant effect on their underlying businesses. As a result, to the extent that the business cycle drives banks' information quality, we would expect this effect to be larger for firms in nontradeable industries. In Figure 5, we test this prediction by comparing the AUC across high and low unemployment periods separately for each group of industries. Consistent with our hypothesis, the difference in AUCs across high and low unemployment periods is larger and only statistically significant for nontradeable firms. This result suggests that local economic conditions are an important driver of the cyclicalities of banks' information quality.

One concern could be that differences in the distribution of underlying borrowers evolve

¹⁵The Y-14Q data use ZIP codes as geographical identifiers, so we first use the ZIP-to-county crosswalk from the Department of Housing and Urban Development to assign a county to each zip code before merging it with the unemployment rate data.

¹⁶Online Appendix Figure OA.3 also shows that our main results hold when we exclude counties with few loans, and Online Appendix Table OA.8 shows that they hold when we estimate AUCs separately by bank.

¹⁷The list of nontradeable industries includes utilities, construction, wholesale trade, retail trade, transportation, accommodation, food services, information and communication, and professional and administrative services.

over the business cycle, which could mechanically explain these results. First, as shown below in Section 5.3, loan and firm characteristics do not vary significantly over the business cycle. Second, the bottom-left panel of Figure 1 shows that the distribution of PDs appears fairly similar across high and low unemployment periods, while Table 2 shows only small differences in the mean and standard deviation of PD for newly originated loans (1.66pp versus 1.61pp and 2.88pp versus 2.58pp, respectively). In Appendix Section B.4, we show that, in the absence of differences in banks’ information over the cycle, the observed difference in the distribution in PDs across high and low unemployment regimes has a quantitatively small effect on the AUC and is hence unlikely to explain our results.¹⁸

Another concern is that banks may have incentives to adjust their PDs because they are used as inputs to calculate capital requirements. Because the ROC curve is purely ordinal, any adjustments to PDs would not affect the AUC so long as they do not change their relative ordering. For example, if a bank reduced all of its PDs by 0.1pp, the AUC would be completely unaffected. Nonetheless, in Online Appendix Figure OA.9, we split our sample into loans granted by banks with high amounts of capital (as defined by having a total risk-based capital ratio above the median of all banks in our data in each quarter) versus those with low amounts of capital. We find that the AUC is actually *higher* for the low-capital banks despite the potentially stronger incentives to manipulate PDs.

Overall, the results in this section suggest that periods of elevated unemployment are associated with improvements in bank information quality. Hence, we conclude that bank information quality is countercyclical.

5 Mechanisms

In Section 4, we show that banks have better information about borrowers during downturns; however, these results are entirely agnostic to the underlying mechanism. In this section, we provide evidence that banks’ improved information during downturns is the result of endogenous information production.

As previously discussed, many models of endogenous information production, including the one developed in Appendix B, predict that banks will produce more information during downturns. However, because we cannot directly observe banks’ information production decisions—only the ability of their PDs to predict default—it is possible that our empirical results could instead arise purely through exogenous variation in information quality over the business cycle. For example, if more firms become delinquent in periods of high unemployment, this may exogenously provide banks with more private information, which allows them to distinguish better across borrower types. Alternatively, downturns can *publicly* reveal default-relevant information. For example, Warren Buffett famously said, “Only when the tide goes

¹⁸We also show, that for similar reasons, differences in the distribution of aggregate shocks over the business cycle are unlikely to explain our results.

out do you learn who has been swimming naked,” suggesting that downturns can publicly reveal which firms are performing well and which are performing poorly.¹⁹ Finally, if the set of borrowers who approach banks for loans changes, banks may exogenously have better information regarding these specific borrowers. While exogenous variation in information quality is not necessarily mutually exclusive with banks endogenously producing more information as economic conditions deteriorate, this section develops several additional tests to rule out the possibility that it is the primary driver of our results.

In Section 5.1, we first show that the cyclicalities of information quality is higher for new loans and that, when we expand the sample to include existing loans, the unemployment rate at origination has persistent effects on information quality. In Section 5.2, we show that the effects are concentrated in larger loans and loans with higher expected losses. Finally, in Section 5.3, we analyze how lending volume and borrower risk evolve over the business cycle and argue that both are consistent with models of endogenous information production.

5.1 Information Sensitivity of New Loans

We first compare the cyclicalities of bank information quality for newly issued loans to those that were issued in prior quarters. Intuitively, the marginal value of information about a borrower’s creditworthiness should be highest prior to the capital being sunk, and thus, new loans should be more information sensitive. If banks’ incentives are driving them to produce more information in bad times, we would thus expect these effects to be stronger for new loans than for loans made in the past.

To test this hypothesis, we extend our sample to include observations for every quarter that each loan was on banks’ balance sheets rather than focusing exclusively on the quarter of origination.²⁰ In the left panel of Figure 6, we reestimate the ROC curves in high and low unemployment periods based on this larger sample. The AUC is larger during periods of high unemployment, confirming our earlier findings that PDs discriminate better in bad times, even for previously issued loans. However, while this difference in AUC is statistically significant, the magnitude of the difference (0.830 versus 0.806) is smaller in both absolute and relative terms than our baseline sample of newly issued loans shown in Figure 4 (0.724 versus 0.681).

The previous test suggests that banks’ information quality is most sensitive to the business cycle when loans are originated. If banks are indeed producing more information at origination, we would expect the economic conditions at origination to have a persistent effect on information quality. To test this hypothesis, we estimate AUCs using current PDs for separate groups based on each loan’s unemployment rate at *origination*. We report the results of this test in the right panel of Figure 6, which shows the estimated ROC curves for high and low unemployment periods based on the origination date of the loan. The difference in AUC

¹⁹See [Berkshire Hathaway 2001 Annual Report](#).

²⁰The summary statistics for this extended sample, which includes recent observations of loans originated prior to the start of our sample in 2014Q4, are shown in Online Appendix Table OA.1.

across periods of high and low unemployment at origination (0.806 versus 0.768) is larger than the differences based on the current unemployment rate shown in the left panel (0.830 versus 0.806).²¹ This result suggests that economic conditions at origination have large and persistent effects on information quality and provide support for endogenous information production as a driving force behind the cyclicalities of information quality that we observe in the data.

5.2 Information Sensitivity, Loan Size, and Expected Losses

In the previous section, we document that information quality is more cyclical for new loans than existing loans. In this section, we analyze how information quality varies across different types of new loans, which theory predicts to be more information sensitive. Specifically, several theories of endogenous information production, such as [Manove, Padilla, and Pagano \(2001\)](#), [Dang, Gorton, and Holmström \(2012\)](#) and [Asriyan, Laeven, and Martin \(2022\)](#) predict that lenders have stronger incentives to produce information about loans that are larger and have higher expected losses. Intuitively, because these are loans for which banks face more severe consequences for lending to low quality borrowers, they should produce more information, and hence, their information quality should be both better and more cyclically sensitive.

We first analyze the effects of loan size on information quality by grouping loans based on their exposure at default (EAD), which is a measure of expected loan size at the time of default.²² The top panel of Figure 7 shows ROC curves for loans split into quartiles of EAD within each bank/quarter. As the Figure shows, the AUC increases from 0.666 for the smallest quartile to 0.716 for the largest, suggesting that the discriminatory power of PDs grows with the size of the loan.

We next examine how the sensitivity of information quality to EAD evolves over the business cycle. In addition to suggesting that information quality should improve with loan size, theories such as [Dang, Gorton, and Holmström \(2012\)](#) and [Biswas \(2022\)](#) also predict that information quality should be more cyclically sensitive for these loans. We test this prediction in the bottom panels of Figure 7 by comparing differences in information quality for new loans made during periods of high and low unemployment for the largest and smallest EAD quartiles. During periods of low unemployment (lower-left panel), the difference in AUCs between the largest and smallest loans is small (0.010) and statistically insignificant. However, when the unemployment rate is high (lower-right panel), the difference in AUCs is much larger (0.109) and statistically significant. These results suggest that information quality is both higher and more cyclically sensitive for larger loans, which provides additional support for the endogenous information production channel. They are also, to our knowledge, the first direct

²¹The loan-level correlations between the current and origination unemployment indicators are 0.19 for the entire sample and 0.14 after excluding new loans (where the two will be the same).

²²Specifically, the EAD incorporates the expected utilization and prepayment of principal by the borrower. Using this as a measure of loan size is particularly advantageous for credit lines in which the utilization rate changes over time. For this reason, EAD is a key parameter for calculating regulatory capital under Basel (e.g., [Calculation of RWA for credit risk](#)).

empirical evidence that a loan’s size affects its information sensitivity as well as the cyclical nature of its information sensitivity.

All else equal, larger loans expose banks to higher losses. However, expected losses are also affected by the loss given default (LGD) and the likelihood of default itself (PD). Thus, we next analyze how information quality varies across a bank’s *expected losses* (EL) for each loan, which are defined as the product of EAD, LGD, and PD. Beyond providing a more direct measure of banks’ exposure to loans, expected losses are a key measure of credit risk used to calculate regulatory capital under Basel ([Basel Committee on Banking Supervision \(2015\)](#)). If banks endogenously produce information in response to their incentives, we expect information quality to be both higher and more cyclically sensitive for loans with higher expected losses.

To test this hypothesis, we repeat our previous analysis but instead use quartiles of expected losses in Figure 8. As with EAD, the top panel shows that information quality is lower for loans with the smallest ELs and monotonically improves as EL becomes larger. The bottom panels compare the cyclical sensitivity of information quality for the top and bottom quartiles across periods of low and high local unemployment. The bottom-left panel shows that differences in information quality across loans with the largest and smallest EL remain significant even during periods of low unemployment (0.655 vs 0.552). However, similar to the results regarding EAD in Figure 7, the bottom-right panel shows that this difference is larger during periods of high local unemployment (0.761 vs 0.592).

One potential concern is that the difference in the underlying distribution of PDs in high and low unemployment periods may differ across the various loan characteristics we analyze in Sections 5.1 and 5.2. As discussed in Section 4.2, differences in the underlying distribution of PDs can mechanically affect the AUC. However, in Online Appendix Section C.4, we show that, in the absence of differences in information production, mechanical differences in the distribution of underlying PDs are unlikely to explain our cross-sectional results.

The results in Sections 5.1 and 5.2 help rule out the possibility that our results in Section 4.2 are driven purely by exogenous variation in information over the cycle. Specifically, if downturns exogenously reveal which firms are “naked,” we would expect this to be true for all loans rather than just new loans, large loans, or loans with higher expected losses. Similarly, if banks choose to lend only to the borrowers they happen to have more information about in bad times, we do not see how this alternative mechanism can explain why banks have exogenously better information for large loans and those with higher expected losses.²³

5.3 Lending Outcomes Over the Business Cycle

The evidence in the previous two sections is consistent with banks endogenously producing more information in downturns because they have stronger incentives to do so. In this section,

²³As we discuss below in Section 5.3, we do believe that riskier borrowers are approaching banks for loans; however, this should not explain our main result so long as banks do not have exogenously better information about these firms, which our results in this and the previous section help rule out.

we show that while the characteristics of new loans in a county do not meaningfully change as the local unemployment rate rises, the number and volume of new loans decline sharply. We argue that this is consistent with banks lending more selectively in downturns due to their increased information production.

We first estimate the following regression across different outcome variables y_i :

$$y_i = \beta U R_{c,t} + \Omega X_i + \delta_{b,t} + \gamma_{j,t} + \sigma_{b,c} + \epsilon_i.$$

This regression includes the same firm-level characteristics and fixed effects that we use in (1); however, we exclude loan characteristics as controls and instead include them as dependent variables. The coefficient β reflects how each of these characteristics changes when the local unemployment rate is above its median. We cluster standard errors by county. Table 5 displays the results.

Loan amounts and loan maturities do not vary over the business cycle in a statistically significant way. Interest rates and PDs are only marginally higher in bad times—about 3bps and 6bps higher when the local unemployment rate is above its median, respectively—and neither difference is statistically significant. Similarly, Table 6, which reports the results of a similar regression using firm outcomes as the dependent variable while controlling for loan characteristics, shows that firm characteristics do not appear to move meaningfully over the cycle. Hence, while the pool of potential borrowers is likely to be riskier in downturns, the pool of loans actually granted does not seem substantially riskier.

Despite minimal changes in loan and firm characteristics over the business cycle, we do observe large changes in lending volume over the business cycle. We aggregate the number and total volume of loans to the county level, take logs, and then regress these measures on the high-unemployment indicator and county fixed effects. The results are reported in Table 7, which shows a decline in both the number and total volume of loans in a county as its economic conditions worsen. Specifically, periods of above-median local unemployment are associated with a 6.7% decrease in the number of loans and an 8.9% decrease in total loan volume. Together, these results suggest that local downturns primarily affect the number and volume of loans issued by banks rather than the composition of loan types or borrowers.

Our model in Appendix B shows that these results, combined with our earlier results on the higher AUCs in bad times, can be rationalized if the average risk of *potential* borrowers is higher in downturns.²⁴ In our model, the value of producing information is to screen out lower quality borrowers. Hence, increased information production in downturns results in lower aggregate lending volume, a higher AUC due to more informative PDs, and the possibility that

²⁴See Proposition 3 for more details. At first glance, it may seem puzzling that more information is produced in downturns even though the pool of borrowers receiving credit is not riskier. What matters in both our model and other models in the literature is not the average risk of borrowers, but the difference in risk across unobservable borrower types, which is higher in downturns.

the average risk of borrowers receiving credit remains the same.²⁵

6 Information Quality and House Price Growth

The previous sections establish that bank information quality is countercyclical and that this effect is concentrated among loans that are more information sensitive, such as loans with higher expected losses. A key determinant of expected losses is collateral values, which recent theories suggest play an important role in how banks' information production incentives evolve over the business cycle. For example, [Asriyan, Laeven, and Martin \(2022\)](#) show that booms driven by high collateral values result in reduced information production as lenders face lower expected losses.²⁶ This section tests the implications of these theories by analyzing how information quality varies with changes in local house prices, which are commonly used in the literature as a proxy for collateral values.

We first provide evidence of the link between local house prices and loan-level measures of expected losses. Specifically, we estimate the following regression:

$$\Delta V_{i,t} = \beta R_t + \epsilon_{i,t},$$

where the dependent variable, $\Delta V_{i,t}$, is either the quarterly change in a loan's LGD or log expected loss. To the extent that local house price increases are associated with rising collateral values, higher house prices should lead to higher recovery values and lower expected losses.²⁷ The independent variable of interest is R_t , which is the county-level quarterly house price return (from $t - 1$ to t) obtained from Zillow. In some specifications, we also include county and/or loan fixed effects. We cluster our standard errors by county.

Table 8 displays the results. The heading at the top row reports the dependent variable in each regression. The first three columns under each heading report results with and without loan or county fixed effects. The second three columns in each heading repeat the same exercises while also controlling for whether the county-level unemployment rate was above or below its median. Across all specifications, stronger appreciation in house prices leads to a statistically significant decrease in LGD and expected losses, which is consistent with the positive

²⁵Of course, there is no requirement that the average risk of borrowers remains the same; the net effect could be positive or negative. We only argue that it is theoretically possible for this mechanism to explain what we observe in the data. The mechanism in our model in which lending volume and information production are inextricably linked is very similar to those in the literature, such as [Dang, Gorton, and Holmström \(2013\)](#), [Gorton and Ordonez \(2014\)](#), and [Asriyan, Laeven, and Martin \(2022\)](#), but increased information production need not result in lower lending volume. For instance, banks may produce more information but maintain the same lending standards.

²⁶See also [Gorton and Ordonez \(2020\)](#) for a theory in which “good booms” are driven by increases in productivity, while “bad booms” are driven by increases in collateral values.

²⁷While the Y-14Q data do include a field for the market value of collateral, this variable is either missing or zero for the vast majority of our sample. However, the value of collateral should be reflected in LGDs and, in turn, expected losses (see [Frye et al. \(2000\)](#)). While we could compare collateralized and non-collateralized loans, roughly 90% of all new loans in our sample report having some collateral. Moreover, in most models, collateral typically matters because it affects expected losses and we analyze expected losses directly in Figure 8.

relationship between collateral values and real estate prices documented in [Chaney, Sraer, and Thesmar \(2012\)](#). Specifically, a 1% increase in housing prices leads to just over a 1% decrease in expected losses. The fact that this relationship holds even after controlling for local economic conditions suggests that collateral values—and, in turn, expected losses—do not simply move uniformly with the business cycle.

Having verified the empirical link between local house prices and measures of expected losses, we next test whether changes in house prices drive differences in information quality. The left panel of Figure 9 shows ROC curves estimated from our baseline sample of new loans split by whether county-level house price growth at origination was above or below its median for that county across our sample period. The AUC is much higher in periods of low house price growth (0.744) than in periods of high house price growth (0.655), and this difference is statistically significant and is consistent with banks having weaker incentives to produce information when house prices increase. The correlation between the loan-level indicator for “high unemployment” and “low house price growth” is 0.08 in our baseline sample of new loans, suggesting that these two measures reflect independent sources of variation in local conditions.

We next show that house price growth affects banks’ information quality even when we restrict the sample to periods of low unemployment in the right panel of Figure 9. This exercise is motivated by [Asriyan, Laeven, and Martin \(2022\)](#), who show that booms driven by high collateral values result in information depletion due to diminished information production incentives. For this analysis, we restrict the sample to periods of low unemployment (i.e., booms) and test for differences in information quality based on house price returns. The difference in information quality remains statistically and economically significant in this sample, which provides direct empirical support for theories of endogenous information production, such as [Asriyan, Laeven, and Martin \(2022\)](#), by showing that variation in collateral values affects information quality.

Taken together, these exercises provide additional support for the countercyclical information quality we document being driven by banks’ endogenous information production decisions. However, they also highlight that not all economic expansions and downturns are alike in terms of their effect on information quality. Booms (or busts) accompanied by rapid growth in collateral values are more prone to reduced information production.

7 Conclusion

Information plays a crucial role in banks’ lending decisions and, in turn, macroeconomic outcomes, but it is difficult to analyze empirically. In this paper, we analyze bank information quality from confidential regulatory data containing banks’ private risk assessments of their borrowers. Using county-level variation in unemployment rates, we find that information quality improves as local economic conditions worsen. We provide evidence that these results are consistent with theories of endogenous information production by showing that information

quality is higher during downturns for newly originated loans and loans with higher expected losses, and that information quality is lower during booms accompanied by high house price growth. To our knowledge, our findings are the first in the empirical banking literature providing evidence of countercyclical information *production*.

While the focus of our analysis is on how banks' information production evolves over the business cycle, our paper demonstrates how banks' private risk assessments in Y-14Q data can be used to analyze how banks' information production decisions vary across other dimensions. This opens up promising avenues for future research examining how various factors—such as market structure, competition, technological change, organizational design, and capital requirements—shape banks' information production decisions.

References

- Asea, Patrick K and Brock Blomberg, 1998, Lending cycles, *Journal of Econometrics* 83, 89–128.
- Asriyan, Vladimir, Luc Laeven, and Alberto Martin, 2022, Collateral booms and information depletion, *The Review of Economic Studies* 89, 517–555.
- Basel Committee on Banking Supervision, 2005, Studies on the validation of internal rating systems, Working Paper No. 14, Bank for International Settlements.
- Basel Committee on Banking Supervision, 2015, Guidance on credit risk and accounting for expected credit losses.
- Bassett, William F, Mary Beth Chosak, John C Driscoll, and Egon Zakrajšek, 2014, Changes in bank lending standards and the macroeconomy, *Journal of Monetary Economics* 62, 23–40.
- Becker, Bo, Marieke Bos, and Kasper Roszbach, 2020, Bad times, good credit, *Journal of Money, Credit and Banking* 52, 107–142.
- Bedayo, Mikel, Gabriel Jiménez, José-Luis Peydró, and Raquel Vegas Sánchez, 2020, Screening and loan origination time: lending standards, loan defaults and bank failures .
- Behn, Markus, Rainer FH Haselmann, and Vikrant Vig, 2016, The limits of model-based regulation .
- Benedetti, Riccardo, 2010, Scoring rules for forecast verification, *Monthly Weather Review* 138, 203–211.
- Berg, Tobias and Philipp Koziol, 2017, An analysis of the consistency of banks' internal ratings, *Journal of Banking & Finance* 78, 27–41.
- Berg, Tobias, Manju Puri, and Jörg Rocholl, 2020, Loan officer incentives, internal rating models, and default rates, *Review of Finance* 24, 529–578.
- Beyhaghi, Mehdi, Cesare Fracassi, and Gregory Weitzner, 2025, Adverse selection in corporate loan markets, *Journal of Finance* forthcoming.
- Bidder, Rhys M, Nicolas Crouzet, Margaret Jacobson, and Michael Siemer, 2023, Debt flexibility .
- Bidder, Rhys M, John R Krainer, and Adam Hale Shapiro, 2020, De-leveraging or de-risking? how banks cope with loss, *Review of Economic Dynamics* .
- Biswas, Sonny, 2022, Collateral and bank screening as complements: A spillover effect, *Working Paper* .
- Boyd, John H and Edward C Prescott, 1986, Financial intermediary-coalitions, *Journal of Economic theory* 38, 211–232.

- Bradley, Andrew P, 1997, The use of the area under the roc curve in the evaluation of machine learning algorithms, *Pattern recognition* 30, 1145–1159.
- Breiman, Leo, 2001, Random forests, *Machine learning* 45, 5–32.
- Brier, Glenn W, 1950, Verification of forecasts expressed in terms of probability, *Monthly weather review* 78, 1–3.
- Brown, James R, Matthew T Gustafson, and Ivan T Ivanov, 2021, Weathering cash flow shocks, *The Journal of Finance* 76, 1731–1772.
- Cerqueiro, Geraldo, Steven Ongena, and Kasper Roszbach, 2016, Collateralization, bank loan rates, and monitoring, *The Journal of Finance* 71, 1295–1322.
- Chaney, Thomas, David Sraer, and David Thesmar, 2012, The collateral channel: How real estate shocks affect corporate investment, *American Economic Review* 102, 2381–2409.
- Dang, Tri Vi, Gary Gorton, and Bengt Holmström, 2012, Ignorance, debt and financial crises, *Unpublished, Yale SOM*.
- Dang, Tri Vi, Gary Gorton, and Bengt Holmström, 2013, The information sensitivity of a security, *Unpublished working paper, Yale University* 39–65.
- Dell’Ariccia, Giovanni and Robert Marquez, 2006, Lending booms and lending standards, *The Journal of Finance* 61, 2511–2546.
- Dell’Ariccia, Giovanni, Deniz Igan, and Luc UC Laeven, 2012, Credit booms and lending standards: Evidence from the subprime mortgage market, *Journal of Money, Credit and Banking* 44, 367–384.
- DeLong, Elizabeth R, David M DeLong, and Daniel L Clarke-Pearson, 1988, Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach, *Biometrics* 837–845.
- Diamond, Douglas W, 1984, Financial intermediation and delegated monitoring, *The review of economic studies* 51, 393–414.
- Engelmann, Bernd and Robert Rauhmeier, *The basel II risk parameters: estimation, validation, stress testing-with applications to loan risk management* (Springer Science & Business Media 2011).
- Farboodi, Maryam and Peter Kondor, 2020, Rational sentiments and economic cycles, Working paper, National Bureau of Economic Research.
- Fishman, Michael J, Jonathan A Parker, and Ludwig Straub, 2020, A dynamic theory of lending standards, Working paper, National Bureau of Economic Research.
- Frame, W Scott, Ping McLemore, and Atanas Mihov, 2020, Haste makes waste: Banking organization growth and operational risk.
- Frye, Jon, Lisa Ashley, Robert Bliss, Richard Cahill, Paul Calem, Matthew Foss, Michael Gordy, David Jones, Catherine Lemieux, Michael Lesiak et al., 2000, Collateral damage: A source of systematic credit risk, *Risk* 13, 91–94.
- Gorton, Gary and Guillermo Ordonez, 2014, Collateral crises, *American Economic Review* 104, 343–78.
- Gorton, Gary and Guillermo Ordonez, 2020, Good booms, bad booms, *Journal of the European Economic Association* 18, 618–665.
- Gorton, Gary B and Ping He, 2008, Bank credit cycles, *The Review of Economic Studies* 75, 1181–1214.
- Goulet Coulombe, Philippe, Maxime Leroux, Dalibor Stevanovic, and Stéphane Surprenant, 2022, How is machine learning useful for macroeconomic forecasting?, *Journal of Applied Econometrics* 37, 920–964.
- Granger, Clive WJ, 1969, Prediction with a generalized cost of error function, *Journal of the Operational Research Society* 20, 199–207.

- Grether, David M, 1980, Bayes rule as a descriptive model: The representativeness heuristic, *The Quarterly journal of economics* 95, 537–557.
- Gu, Shihao, Bryan Kelly, and Dacheng Xiu, 2020, Empirical asset pricing via machine learning, *The Review of Financial Studies* 33, 2223–2273.
- Gustafson, Matthew, Ivan Ivanov, and Ralf R Meisenzahl, 2020, Bank monitoring: Evidence from syndicated loans, *Available at SSRN 2831455*.
- Hayashi, Yoichi, 2022, Emerging trends in deep learning for credit scoring: A review, *Electronics* 11, 3181.
- Iyer, Rajkamal, Asim Ijaz Khwaja, Erzo FP Luttmer, and Kelly Shue, 2016, Screening peers softly: Inferring the quality of small borrowers, *Management Science* 62, 1554–1577.
- James, Christopher, 1987, Some evidence on the uniqueness of bank loans, *Journal of financial economics* 19, 217–235.
- Khwaja, Asim Ijaz and Atif Mian, 2008, Tracing the impact of bank liquidity shocks: Evidence from an emerging market, *American Economic Review* 98, 1413–1442.
- Leland, Hayne E and David H Pyle, 1977, Informational asymmetries, financial structure, and financial intermediation, *The journal of Finance* 32, 371–387.
- Lessmanna, Stefan, H Seowb, Bart Baesenscd, and Lyn C Thomasd, Benchmarking state-of-the-art classification algorithms for credit scoring: A ten-year update, *Credit Research Centre, Conference Archive* (2013).
- Lown, Cara and Donald P Morgan, 2006, The credit cycle and the business cycle: new findings using the loan officer opinion survey, *Journal of Money, Credit and Banking* 1575–1597.
- Maddaloni, Angela and José-Luis Peydró, 2011, Bank risk-taking, securitization, supervision, and low interest rates: Evidence from the euro-area and the us lending standards, *the review of financial studies* 24, 2121–2165.
- Manove, Michael, A Jorge Padilla, and Marco Pagano, 2001, Collateral versus project screening: A model of lazy banks, *Rand journal of economics* 726–744.
- Mincer, Jacob A and Victor Zarnowitz, The evaluation of economic forecasts, *Economic forecasts and expectations: Analysis of forecasting behavior and performance*, 3–46 (NBER 1969).
- Murphy, Allan H, 1973, A new vector partition of the probability score, *Journal of Applied Meteorology and Climatology* 12, 595–600.
- Muth, John F, 1961, Rational expectations and the theory of price movements, *Econometrica: journal of the Econometric Society* 315–335.
- O’Brien, Robert M, 2007, A caution regarding rules of thumb for variance inflation factors, *Quality & quantity* 41, 673–690.
- Pepe, Margaret Sullivan, *The statistical evaluation of medical tests for classification and prediction* (Oxford university press 2003).
- Petriconi, Silvio, 2015, Bank competition, information choice and inefficient lending booms, *Information Choice and Inefficient Lending Booms (December 9, 2015)*.
- Plosser, Matthew C and Joao AC Santos, 2018, Banks’ incentives and inconsistent risk models, *The Review of Financial Studies* 31, 2080–2112.
- Puri, Manju, Jörg Rocholl, and Sascha Steffen, 2017, What do a million observations have to say about loan defaults? opening the black box of relationships, *Journal of Financial Intermediation* 31, 1–15.
- Rodano, Giacomo, Nicolas Serrano-Velarde, and Emanuele Tarantino, 2018, Lending standards over the credit cycle, *The Review of Financial Studies* 31, 2943–2982.

Ruckes, Martin, 2004, Bank competition and credit standards, *Review of Financial Studies* 17, 1073–1102.

Schonlau, Matthias and Rosie Yuyan Zou, 2020, The random forest algorithm for statistical learning, *The Stata Journal* 20, 3–29.

Tversky, Amos and Daniel Kahneman, 1974, Judgment under uncertainty: Heuristics and biases: Biases in judgments reveal some heuristics of thinking under uncertainty., *science* 185, 1124–1131.

Wager, Stefan and Susan Athey, 2018, Estimation and inference of heterogeneous treatment effects using random forests, *Journal of the American Statistical Association* 113, 1228–1242.

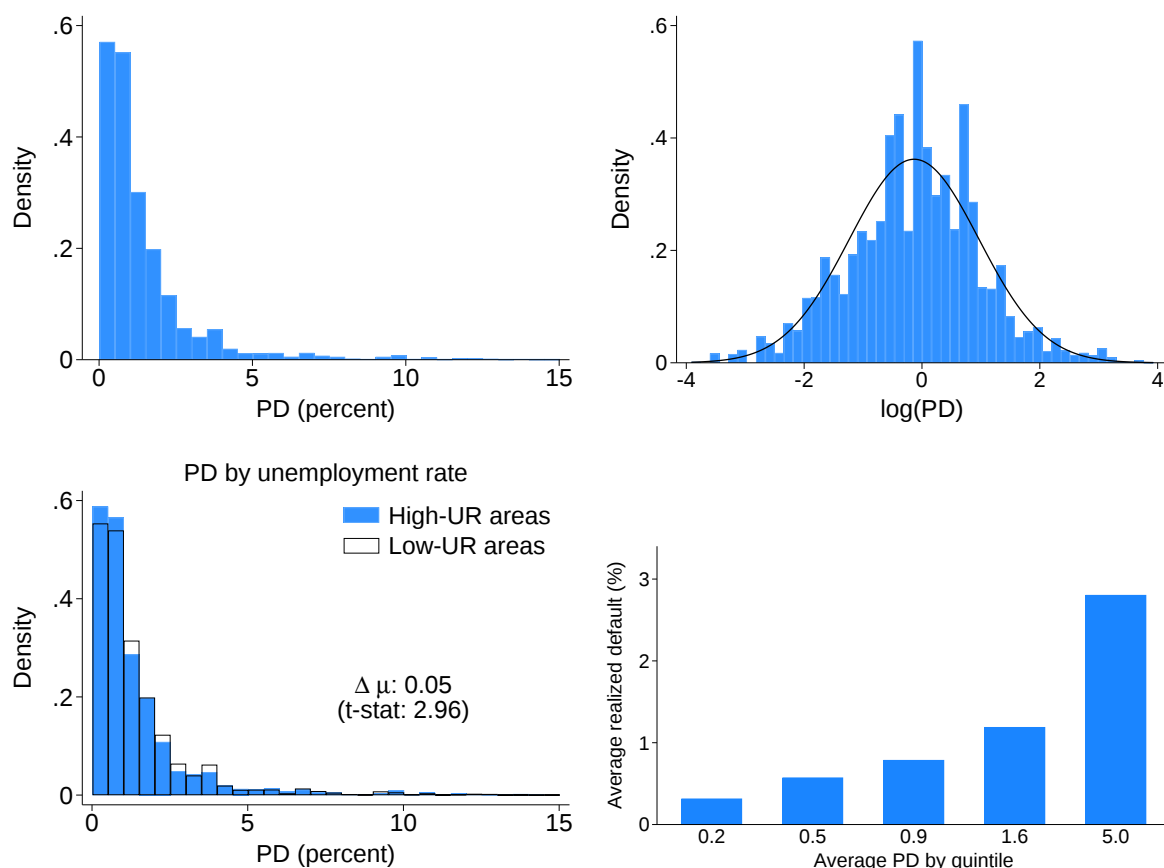


Figure 1: Distributions of PD

The top-left panel plots the distribution of PD. The top-right panel plots the distribution of $\log(\text{PD})$ with an overlaid normal distribution. The bottom-left panel plots overlaid distributions of PDs originated when a county's unemployment rate is above (solid bars) or below (hollow bars) its median from 2014Q4-2021Q4. $\Delta\mu$ reports the difference between the average PD at origination in high-UR areas minus the average PD in low-UR areas, with the corresponding t-statistic shown below in parentheses. The bottom-right panel reports average default rates within two years of origination on the y-axis by quintiles of PD at origination. The numbers beneath each bar correspond to the average PD in each quintile (rounded to the nearest 0.1pp).

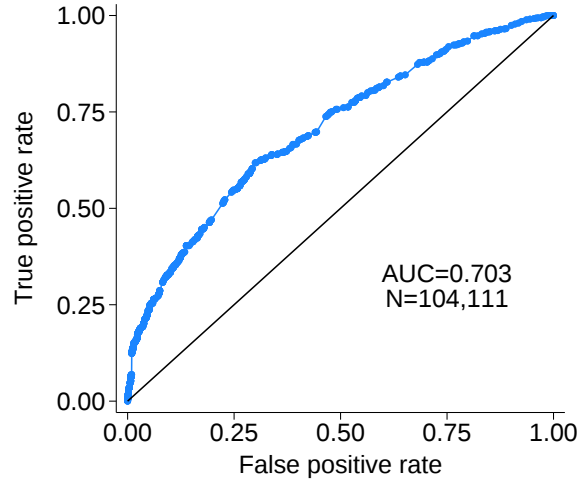


Figure 2: ROC for new loans

This figure shows the ROC curve for all new loans in our sample along with the area under the ROC curve (AUC).

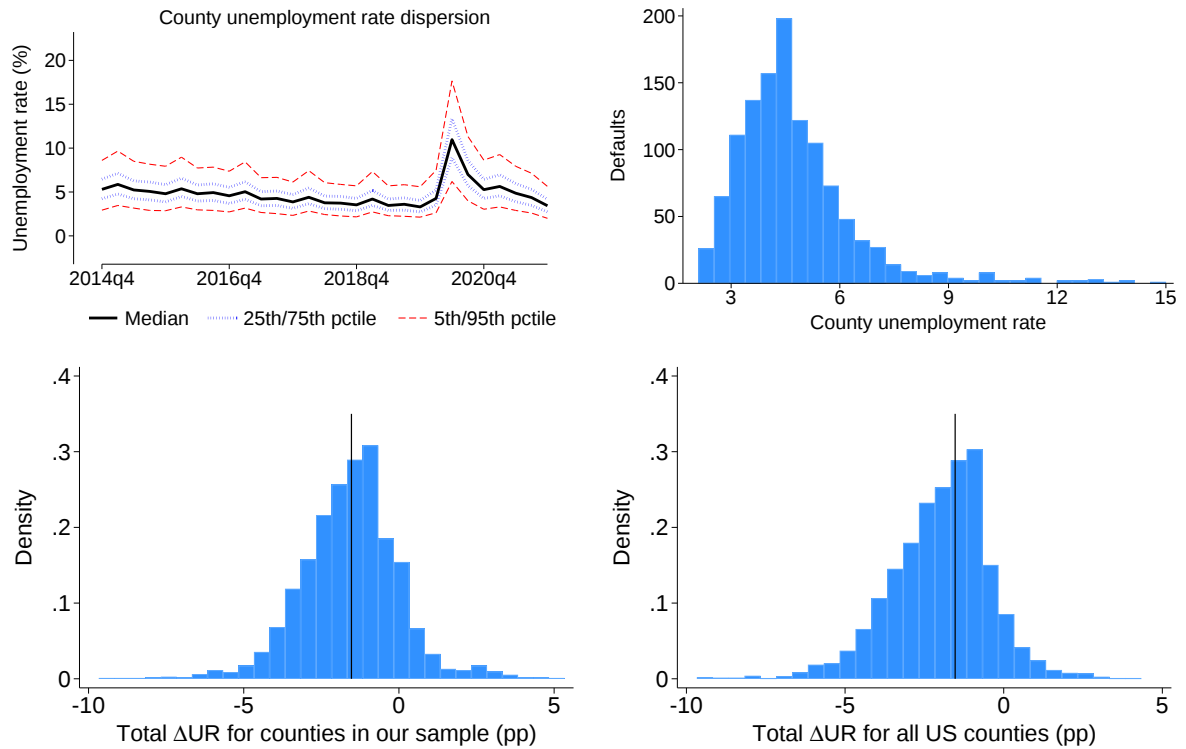


Figure 3: Variation in unemployment rates

The top-left panel shows the range of county-level unemployment rates for all county-quarter observations with at least one new loan in our sample. The top-right panel plots the distribution of the unemployment rate at origination for the 1,176 new loans in our sample that default within two years of origination. The bottom-left panel shows the distribution of county-level changes from the earliest observation of the unemployment rate to the latest for each county with at least two new loans in our sample. The bottom-right panel shows the distribution of county-level changes in the unemployment rate for all US counties from 2014Q4 through 2021Q4. The vertical black lines at -1.5pp show the change in the total US unemployment rate during this period. The bottom panels are truncated between -10pp and +5pp for readability. Sample sizes are shown in panel C of Table

1.

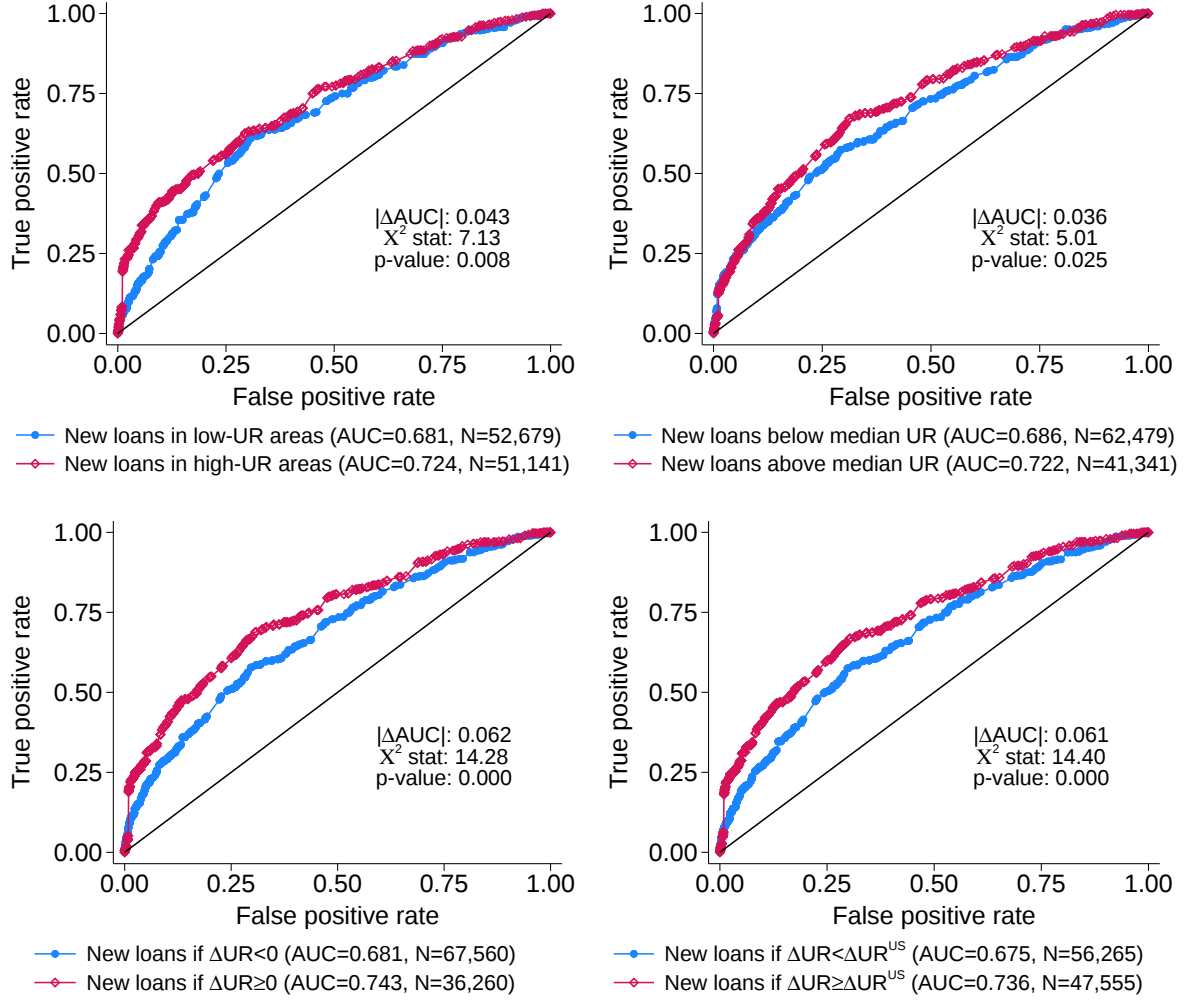


Figure 4: ROC for new loans over the business cycle

The top-left panel shows the ROC curve for all new loans in our sample split by whether the unemployment rate in each period was above or below its county-level median from 2014Q4-2021Q4. The top-right panel shows ROC curves split by whether the unemployment rate at origination was above or below the median unemployment rate of all counties with at least 5 new loans in that quarter. The bottom-left panel shows ROC curves split by whether the county unemployment rate increased from the prior quarter. The bottom-right panel shows ROC curves split whether the change in the county unemployment rate was greater than the change in the total US unemployment rate. The area under each ROC curve (AUC) is reported along with the number of observations in the legend. $|\Delta AUC|$ reports the difference between the two AUCs. Below $|\Delta AUC|$, the DeLong, DeLong, and Clarke-Pearson (1988) statistics are reported: the χ^2 test statistic and its corresponding p-value, which tests the null hypothesis that the difference between the two AUCs equals zero.

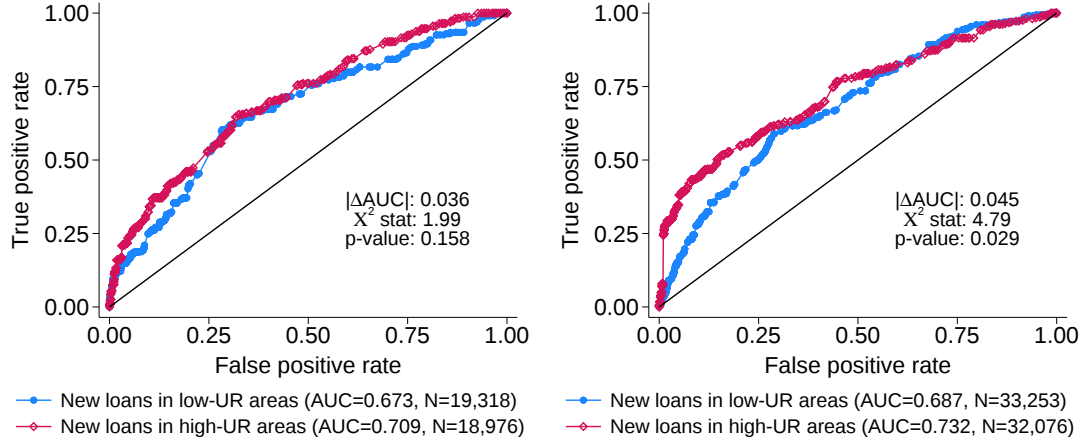


Figure 5: ROC for new loans to tradeable (left) vs nontradeable (right) industries

This figure shows ROC curves for new loans split by the local unemployment rate at origination for tradeable (left) and nontradeable (right) industries. Nontradeable industries include utilities, construction, wholesale trade, retail trade, transportation, accommodation, food services, information and communication, and professional and administrative services; all other loans in our sample with non-missing industry codes are considered tradeable. The area under each ROC curve (AUC) is reported along with the number of observations in the legend. $|\Delta AUC|$ reports the difference between the two AUCs. Below $|\Delta AUC|$, the DeLong, DeLong, and Clarke-Pearson (1988) statistics are reported: the χ^2 test statistic and its corresponding p-value, which tests the null hypothesis that the difference between the two AUCs equals zero.

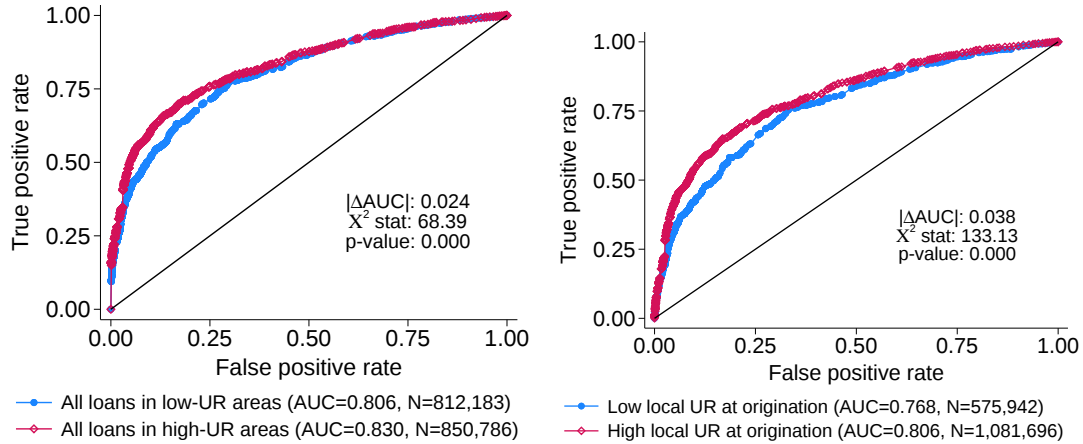


Figure 6: ROC for current (left) versus origination (right) unemployment rate among all loans

The left panel shows ROC curves for both new and existing loans split by whether the contemporaneous unemployment rate in each period was above or below its county-level median from 2014Q4-2021Q4. The right panel shows ROC curves split by whether the unemployment rate at origination was above or below its county-level median from 2014Q4-2021Q4. The area under each ROC curve (AUC) is reported along with the number of observations in the legend. $|\Delta AUC|$ reports the difference between the two AUCs. Below $|\Delta AUC|$, the DeLong, DeLong, and Clarke-Pearson (1988) statistics are reported: the χ^2 test statistic and its corresponding p-value, which tests the null hypothesis that the difference between the two AUCs equals zero.

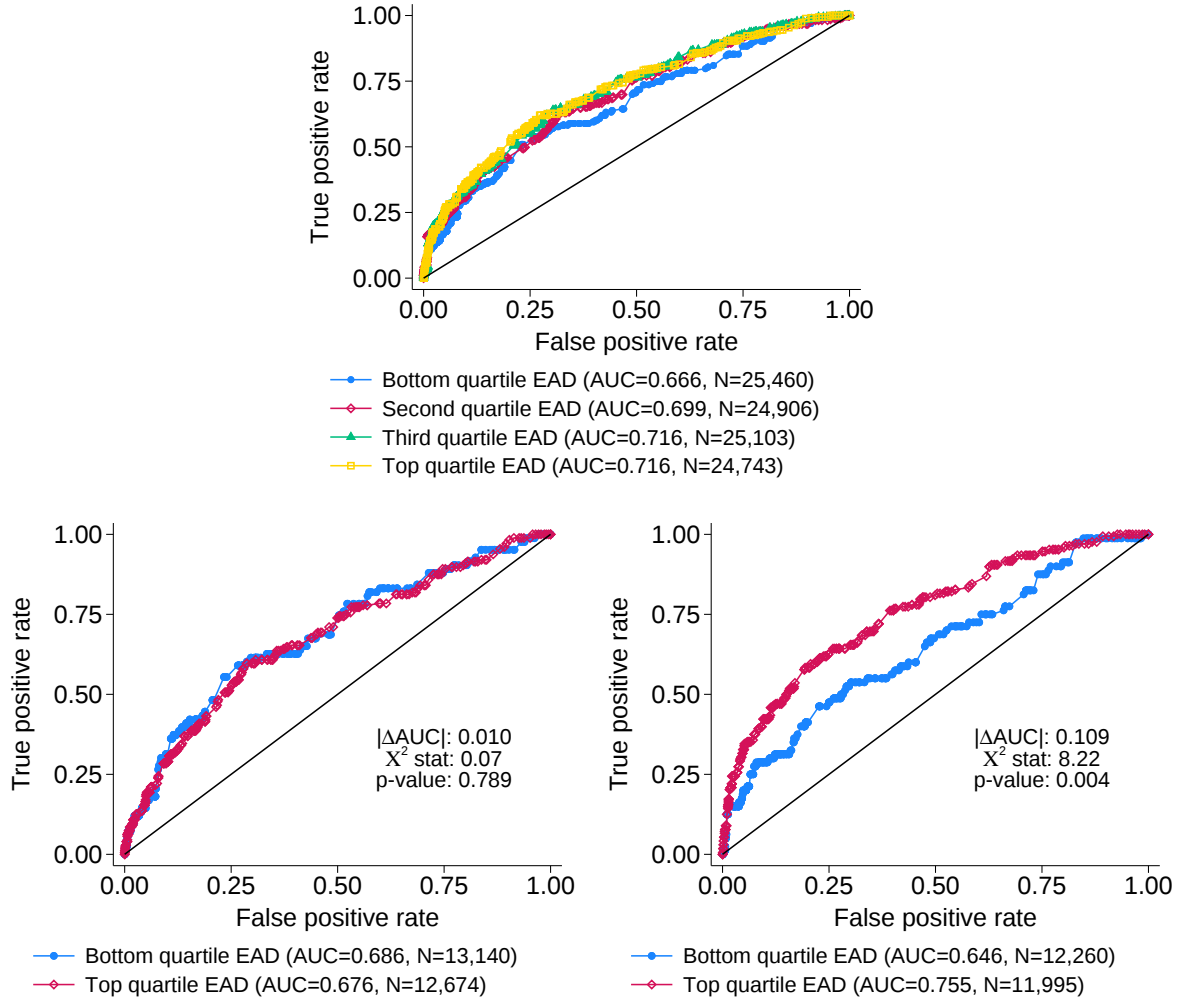


Figure 7: ROC by quartiles of exposure at default

The top panel shows ROC curves for new loans based on their quartile of exposure at default (EAD) calculated within each bank-quarter. The bottom-left panel shows the difference between the highest and lowest quartiles in low-UR areas. The bottom-right panel shows the difference between the highest and lowest quartiles in high-UR areas. The area under each ROC curve (AUC) is reported along with the number of observations in the legend. $|\Delta AUC|$ reports the difference between the two AUCs. Below $|\Delta AUC|$, the [DeLong, DeLong, and Clarke-Pearson \(1988\)](#) statistics are reported: the χ^2 test statistic and its corresponding p-value, which tests the null hypothesis that the difference between the two AUCs equals zero.

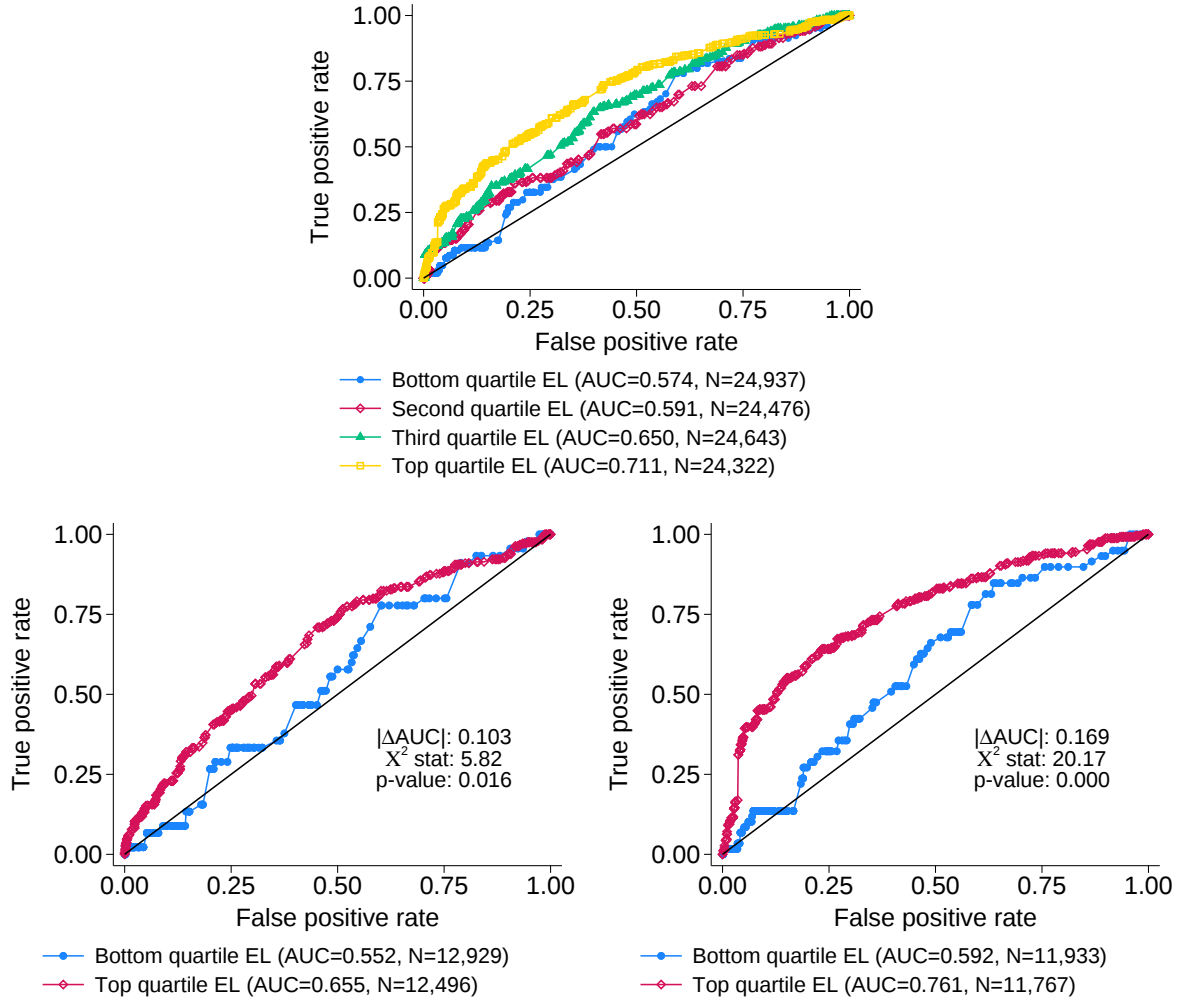


Figure 8: ROC by quartiles of expected loss

This figure shows ROC curves for new loans based on their quartile of expected loss ($EAD \times LGD \times PD$) calculated within each bank-quarter. The bottom-left panel shows the difference between the highest and lowest quartiles in low-UR areas. The bottom-right panel shows the difference between the highest and lowest quartiles in high-UR areas. The area under each ROC curve (AUC) is reported along with the number of observations in the legend. $|\Delta AUC|$ reports the difference between the two AUCs. Below $|\Delta AUC|$, the [DeLong, DeLong, and Clarke-Pearson \(1988\)](#) statistics are reported: the χ^2 test statistic and its corresponding p-value, which tests the null hypothesis that the difference between the two AUCs equals zero.

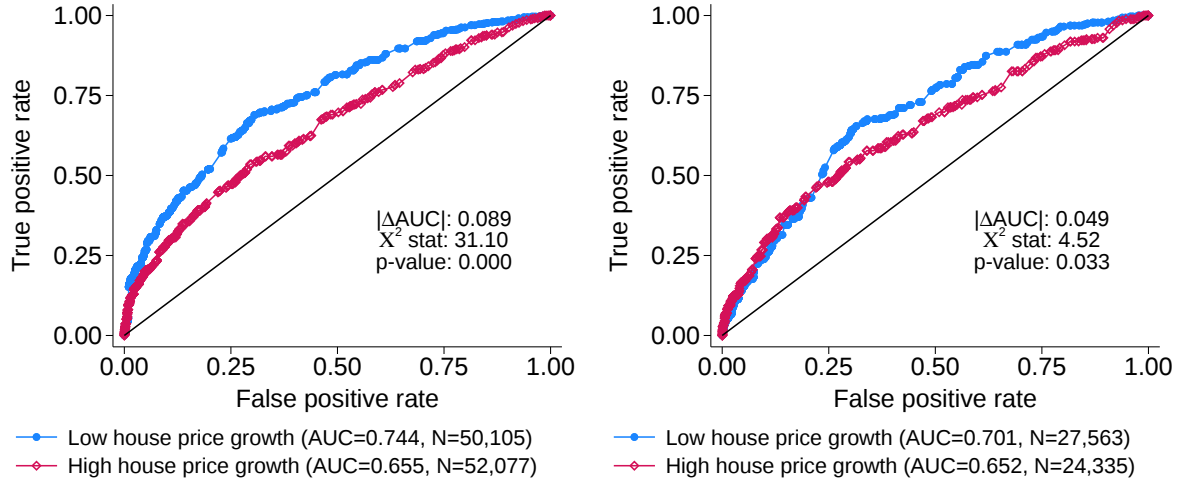


Figure 9: ROC and local house prices

This figure shows ROC curves for new loans split by county-level house price growth. The left panel includes our entire sample of new loans. The right panel includes only new loans in counties where the local unemployment rate was below its county-level median during our sample period. High (low) house price growth is determined by whether the quarterly change in house prices was above (below) that county's median during our sample period. The area under each ROC curve (AUC) is reported along with the number of observations in the legend. $|\Delta AUC|$ reports the difference between the two AUCs. Below $|\Delta AUC|$, the [DeLong, DeLong, and Clarke-Pearson \(1988\)](#) statistics are reported: the χ^2 test statistic and its corresponding p-value, which tests the null hypothesis that the difference between the two AUCs equals zero.

Table 1: Summary statistics

Panel A reports loan characteristics calculated as simple averages across all new loans. Panel B reports firm characteristics, which are calculated after collapsing all loan observations (including both new and existing loans) to the firm-quarter level using simple averages. Panel C reports county characteristics, which are calculated after collapsing all loan observations (including both new and existing loans) to the county-quarter level. Section 2 describes our sample.

	Mean	Median	5%	95%	SD	N
Panel A: Loan characteristics						
Interest rate (pp)	3.70	3.56	1.64	6.16	1.46	87,193
PD (pp)	1.63	0.91	0.14	5.07	2.73	104,111
LGD (ratio)	0.35	0.35	0.08	0.61	0.16	101,884
Realized default (pp)	1.13	0.00	0.00	0.00	10.57	104,111
Maturity (months)	47.49	58.00	7.00	96.00	31.39	104,111
Loan size (\$ mil)	13.32	3.58	1.00	50.00	41.54	104,111
Revolver (indicator)	0.38	0.00	0.00	1.00	0.49	104,111
Term loan (indicator)	0.41	0.00	0.00	1.00	0.49	104,111
Floating rate (indicator)	0.56	1.00	0.00	1.00	0.50	104,111
Panel B: Firm characteristics						
Sales	783.80	44.07	3.34	1,412.12	36,404.71	948,138
Assets	1,328.81	21.52	1.63	1,484.10	253,734.53	947,916
Leverage	0.31	0.26	0.00	0.81	0.27	919,616
Profitability	0.18	0.13	-0.04	0.56	0.24	937,776
Tangibility	0.89	0.99	0.39	1.00	0.20	936,201
Nontradeable	0.65	1.00	0.00	1.00	0.48	1,069,329
PD	2.49	0.91	0.13	9.82	6.87	1,070,295
Total number of loans	1.56	1.00	1.00	4.00	4.11	1,070,295
Number of new loans	0.10	0.00	0.00	1.00	0.47	1,070,295
Number of banks	1.15	1.00	1.00	2.00	0.65	1,070,295
Total loan volume (\$ mil)	20.77	4.16	1.00	85.00	119.29	1,070,295
Panel C: County characteristics						
Unemployment rate	5.01	4.53	2.53	9.03	2.26	52,832
Number of new loans	1.97	0.00	0.00	9.00	7.14	52,910
Number of total loans	31.58	5.00	1.00	138.00	102.11	52,910
Total new loan volume (\$ mil)	420.14	39.87	1.41	1,892.21	1,788.92	52,910

Table 2: Additional loan-level summary statistics

This table contains additional loan-level summary statistics (observations are at the same level as in panel A of Table 1). Panel A includes our entire sample, panel B includes new loans in low-UR areas, and Panel C includes new loans in high-UR areas. Because a small number of loans are made in counties with missing unemployment rates, the sample sizes in panel A are greater than the sum of panels B and C. Section 2 describes our sample.

	Mean	Median	5%	95%	SD	N
Panel A: All loans						
Sales	2,584.48	81.50	3.63	4,252.16	47,106.60	84,208
Assets	3,816.42	50.51	2.04	4,872.10	87,873.09	84,240
Leverage	0.34	0.31	0.00	0.81	0.26	82,291
Profitability	0.20	0.15	-0.00	0.61	0.24	84,240
Tangibility	0.85	0.97	0.31	1.00	0.23	84,040
Nontradeable	0.63	1.00	0.00	1.00	0.48	103,914
PD	1.63	0.91	0.14	5.07	2.73	104,111
Total number of loans	1.00	1.00	1.00	1.00	0.00	104,111
Number of new loans	1.00	1.00	1.00	1.00	0.00	104,111
Number of banks	0.67	1.00	0.00	1.00	0.47	104,111
Total loan volume (\$ mil)	13.32	3.58	1.00	50.00	41.54	104,111
Panel B: Loans in low-UR areas						
Sales	2,410.35	85.49	5.30	4,252.16	28,254.39	42,667
Assets	3,245.06	53.20	2.42	4,829.70	35,858.18	42,721
Leverage	0.34	0.31	0.00	0.80	0.26	41,761
Profitability	0.20	0.15	0.00	0.60	0.23	42,721
Tangibility	0.85	0.97	0.31	1.00	0.23	42,618
Nontradeable	0.63	1.00	0.00	1.00	0.48	52,571
PD	1.61	0.93	0.14	4.66	2.58	52,679
Total number of loans	1.00	1.00	1.00	1.00	0.00	52,679
Number of new loans	1.00	1.00	1.00	1.00	0.00	52,679
Number of banks	0.67	1.00	0.00	1.00	0.47	52,679
Total loan volume (\$ mil)	13.11	3.50	1.00	50.00	35.90	52,679
Panel C: Loans in high-UR areas						
Sales	2,740.06	78.04	2.15	4,175.67	60,795.20	41,306
Assets	4,371.10	47.59	1.65	4,739.00	120,080.33	41,284
Leverage	0.35	0.31	0.00	0.81	0.26	40,299
Profitability	0.20	0.15	-0.01	0.62	0.24	41,284
Tangibility	0.85	0.97	0.31	1.00	0.23	41,187
Nontradeable	0.63	1.00	0.00	1.00	0.48	51,052
PD	1.66	0.90	0.14	5.37	2.88	51,141
Total number of loans	1.00	1.00	1.00	1.00	0.00	51,141
Number of new loans	1.00	1.00	1.00	1.00	0.00	51,141
Number of banks	0.67	1.00	0.00	1.00	0.47	51,141
Total loan volume (\$ mil)	13.49	3.70	1.00	50.00	46.59	51,141

Table 3: PDs predict default

This table shows the results of estimating (1), which tests whether PDs predict default after controlling for observables. The dependent variable in each regression is a dummy variable indicating whether each loan defaults within eight quarters after origination, multiplied by 100. Interest rates and PDs are measured in percentage points. Standard errors are clustered at the county level and are shown below the parameter estimates in parentheses. *, **, and *** indicate statistical significance at the 10%, 5%, and 1% levels, respectively.

	Default				
	(1)	(2)	(3)	(4)	(5)
PD	0.407*** (0.049)	0.444*** (0.064)			0.445*** (0.073)
Interest rate			0.508*** (0.068)	0.452*** (0.078)	0.227*** (0.073)
Leverage		1.027*** (0.310)		1.458*** (0.355)	0.981*** (0.360)
Profitability		0.009 (0.248)		-0.698*** (0.233)	0.029 (0.249)
Tangibility		-0.359 (0.345)		-0.274 (0.417)	-0.289 (0.410)
Log firm size		-0.139*** (0.041)		-0.151*** (0.049)	-0.118** (0.047)
Log loan amount		0.251*** (0.063)		0.267*** (0.071)	0.252*** (0.070)
LGD		0.709* (0.396)		0.130 (0.492)	0.673 (0.487)
Log maturity		-0.058 (0.064)		-0.203** (0.080)	-0.100 (0.074)
Controls	N	Y	N	Y	Y
Bank-quarter FE	Y	Y	Y	Y	Y
Industry-quarter FE	Y	Y	Y	Y	Y
Bank-county FE	Y	Y	Y	Y	Y
Observations	100,368	77,226	83,602	64,566	64,566
R ²	0.174	0.196	0.185	0.206	0.214

Table 4: PDs predict default (random forest)

This table tests whether PDs contain additional information regarding future realized default beyond a predicted default variable estimated using a random forest regression. The dependent variable in each regression is an indicator that equals one if a loan defaults within two years of origination. “RF predicted PD” is the random forest estimate of the loan’s default probability measured in percentage points. The details of these estimates are described in Appendix A. Each column corresponds to a different set of controls used to estimate the random forest. All specifications include our default set of firm and loan controls; other specifications include indicator variables for industry, bank, and time and are indicated in the rows below the table. “PD” is the probability of default reported by the bank and is measured in percentage points. All regressions exclude the half of the baseline sample used to train the random forest. Standard errors are clustered at the county level and are shown below the parameter estimates in parentheses. *, **, and *** indicate statistical significance at the 10%, 5%, and 1% levels, respectively.

	Default				
	(1)	(2)	(3)	(4)	(5)
RF predicted PD	1.087*** (0.075)	1.279*** (0.078)	1.038*** (0.058)	1.153*** (0.070)	1.207*** (0.067)
PD	0.344*** (0.045)	0.295*** (0.041)	0.315*** (0.040)	0.330*** (0.042)	0.270*** (0.039)
Industry controls	N	Y	N	N	Y
Bank controls	N	N	Y	N	Y
Time controls	N	N	N	Y	Y
Observations	51,913	51,913	51,913	51,913	51,913
R ²	0.071	0.104	0.093	0.096	0.150

Table 5: Loan characteristics over the business cycle

This table analyzes the relationship between the local unemployment rate and loan characteristics. “High UR” is a dummy variable equal to one if that county’s unemployment rate was above its median during our sample period. The dependent variable in each regression is shown at the top of each column. All regressions include our default firm controls. The unemployment rate (UR), PD, default, and interest rate are measured in percentage points. Maturity is measured in log months and loan size is measured in log dollars. Standard errors are clustered at the county level and are shown below the parameter estimates in parentheses. *, **, and *** indicate statistical significance at the 10%, 5%, and 1% levels, respectively.

	Loan size	Interest rate	Maturity	Default	PD
High UR	1.437 (1.369)	0.027 (0.017)	1.911 (1.351)	-0.076 (0.143)	0.060 (0.046)
Firm controls	Y	Y	Y	Y	Y
Bank-quarter FE	Y	Y	Y	Y	Y
Industry-quarter FE	Y	Y	Y	Y	Y
Bank-county FE	Y	Y	Y	Y	Y
Observations	79,078	66,149	79,050	79,078	79,078
R ²	0.541	0.552	0.419	0.187	0.293

Table 6: Firm characteristics over the business cycle

This table analyzes the relationship between the local unemployment rate and firm characteristics. “High UR” is a dummy variable equal to one if that county’s unemployment rate was above its median during our sample period. The dependent variable in these regressions is shown at the top of each column. Firm size is the log of total assets, while profitability, leverage, and tangibility are ratios taking on values between zero and one. All regressions include our default loan controls. Standard errors are clustered at the county level and are shown below the parameter estimates in parentheses. *, **, and *** indicate statistical significance at the 10%, 5%, and 1% levels, respectively.

	Firm size	Profitability	Leverage	Tangibility
High UR	-1.780 (2.913)	-0.004 (0.003)	-0.003 (0.004)	-0.004 (0.004)
Loan controls	Y	Y	Y	Y
Bank-quarter FE	Y	Y	Y	Y
Industry-quarter FE	Y	Y	Y	Y
Bank-county FE	Y	Y	Y	Y
Observations	79,128	79,128	77,268	78,937
R ²	0.646	0.295	0.354	0.455

Table 7: County-level lending over the business cycle

This table analyzes the relationship between the unemployment rate and county-level loan volume. “High UR” is a dummy variable equal to 1 if that county’s unemployment rate was above its median during our sample period. Data are aggregated at the county level. The dependent variable in each regression is shown at the top of each column and both are in logs. Standard errors are clustered at the county level and are shown below the parameter estimates in parentheses. *, **, and *** indicate statistical significance at the 10%, 5%, and 1% levels, respectively.

	(1) Loan count	(2) Loan volume
High UR	-0.067*** (0.009)	-0.089*** (0.018)
County FE	Y	Y
Observations	18,396	18,396
R ²	0.749	0.587

Table 8: Collateral values and house prices

This table shows the effects of changes in county-level house prices on loans' loss given default (Columns (1)-(6)) and log expected losses (Columns (7)-(12)). "House price growth" is the quarterly change in house prices from Zillow. "High UR" is a dummy variable equal to 1 if that county's unemployment rate was above its median during our sample period. Standard errors are clustered at the county level and are shown below the parameter estimates in parentheses. *, **, and *** indicate statistical significance at the 10%, 5%, and 1% levels, respectively.

	Δ LGD						Δ EL					
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)
House price growth	-0.021*** (0.006)	-0.027*** (0.006)	-0.020*** (0.006)	-0.027*** (0.006)	-0.034*** (0.007)	-0.026*** (0.007)	-1.143*** (0.099)	-1.257*** (0.117)	-1.196*** (0.131)	-1.102*** (0.096)	-1.207*** (0.115)	-1.130*** (0.126)
High UR				0.001*** (0.000)	0.001*** (0.000)	0.001*** (0.000)				-0.008*** (0.002)	-0.008*** (0.002)	-0.011*** (0.002)
County FE	N	Y	N	N	Y	N	N	Y	N	N	Y	N
Loan FE	N	N	Y	N	N	Y	N	N	Y	N	N	Y
Observations	1,386,280	1,386,247	1,359,216	1,386,280	1,386,247	1,359,216	1,355,384	1,355,350	1,328,971	1,355,384	1,355,350	1,328,971
R ²	0.000	0.002	0.087	0.000	0.002	0.087	0.001	0.003	0.094	0.001	0.003	0.094

Appendix A. Random Forest Regression Methodology and Additional Analysis

This section describes the details behind the random forest regression estimates we use as a nonlinear benchmark for predicting default from observables in Sections 3 and 4. The fundamental building block of a random forest is a *regression tree*. A regression tree generates predictions for an outcome variable by sequentially partitioning observations into K regions (also called “leaves” or terminal nodes) based on similarly valued explanatory variables and then calculating average outcomes within each region. The number of partitions L is called the “depth” of the tree. A simple illustrative example based on a similar figure shown in Gu, Kelly, and Xiu (2020) with $K = 3$ and $L = 2$ is shown below in the left panel of Figure A.1.

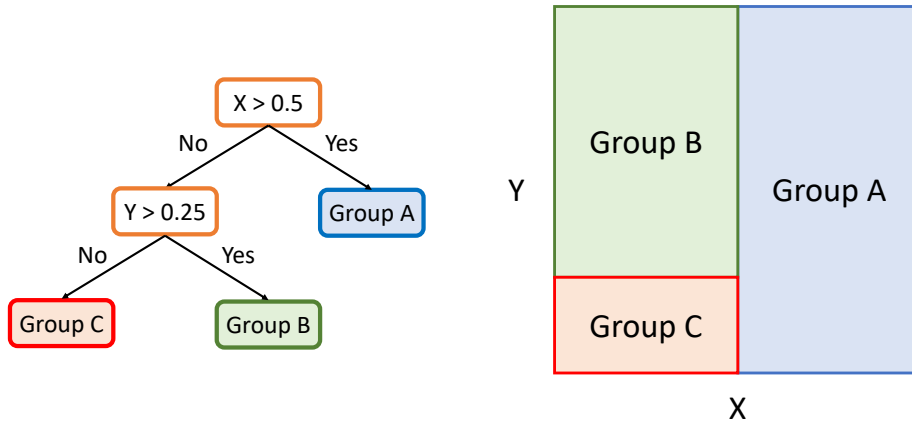


Figure A.1: Example Regression Tree

There are two predictors, X and Y , each of which takes on values between 0 and 1. To obtain a prediction for a given observation i , the node at the top of the tree first partitions the data based on X_i . If $X_i > 0.5$, the observation is assigned to Group A. If $X_i \leq 0.5$, the remaining observations are further partitioned by Y_i ; observations with $Y_i > 0.25$ are assigned to Group B, while observations with $Y_i \leq 0.25$ are assigned to Group C.²⁸ An equivalent graphical representation of this partition is shown in the right panel. A prediction is obtained by taking the sample average of the outcome variable for all observations in a group.

Because regression trees can easily accommodate outliers and idiosyncrasies in the training data, they tend to perform very well in sample, but are prone to over-fitting. Breiman (2001) showed that the random forest model, which averages predictions across a large number of regression trees estimated from bootstrapped data, substantially improves out-of-sample forecast accuracy. We generate random forest predictions of loan defaults using the `rforest` Stata command developed in Schonlau and Zou (2020). We use the default settings for the number of trees (100) and the number of explanatory variables used in each tree (3), and we do not specify a maximum depth or minimum number of observations per node.

²⁸The threshold values shown here at each node are chosen arbitrarily, but in practice, they are usually determined as the solution to an optimization problem such as minimizing the mean squared forecast error.

Following their recommendation, 50% of the data are used for training, while the other 50% are used for validation. We assign observations to each set by first sorting our baseline sample of newly originated loans by loan ID number, and then alternatively placing each loan in the training set or the validation set. The training set is used to estimate the parameters of the random forest, and with these parameters, we generate predictions for the validation set. For consistency, we maintain the same training and validation sets across the range of specifications we estimate. All specifications shown throughout the paper include the observable characteristics described in Section 3: for loans, these are log size, LGD, and log maturity, and for firms, these are leverage, profitability, tangibility, and log assets. Several of the specifications shown in Table 4 also include separate dummy variables for each industry, bank, or time period.

Figure A.2 uses this random forest approach to evaluate the marginal discriminatory power of PD using ROC curves. The blue curve uses default estimates from a random forest regression that uses only the baseline set firm and loan controls, while the red curve includes these same controls plus PD. Neither specification includes any firm, bank, or time dummies. As before, we train each random forest on 50% of the sample and compare the out-of-sample predictive power using the remaining 50%. The difference between the two is statistically significant and larger than the difference between periods of high and low unemployment reported in our baseline results in Figure 4. This provides further evidence that PD contains information useful for predicting default that is not fully captured by nonlinear combinations of other observable characteristics.

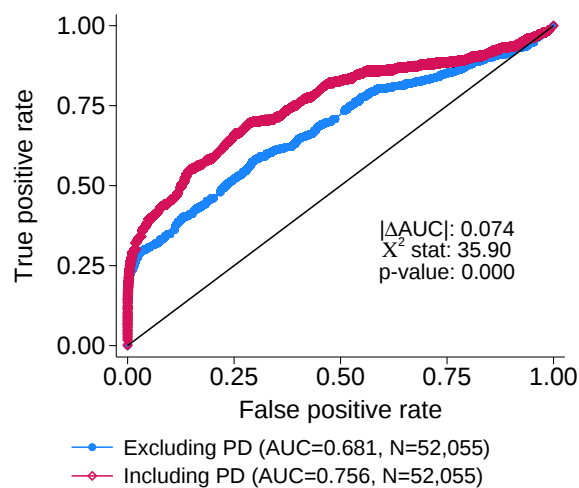


Figure A.2: Comparing random forest predictions with and without actual PDs

This figure compares default predictions from two random forest regressions. The blue curve shows estimates using the following controls: log loan size, LGD, log maturity, firm leverage, firm profitability, firm tangibility, and firm size. The red curve uses the same controls plus PD. The area under each ROC curve (AUC) is reported along with the number of observations in the legend. $|\Delta AUC|$ reports the difference between the two AUCs. Below $|\Delta AUC|$, the DeLong, DeLong, and Clarke-Pearson (1988) statistics are reported: the χ^2 test statistic and its corresponding p-value, which tests the null hypothesis that the difference between the two AUCs equals zero.

Appendix B. Simple Theoretical Framework

In this section, we develop a simple model, from first principles, which shows that increased information production by banks in bad times leads to a higher area under the receiver operating characteristics curve (AUC) in bad times, consistent with our empirical results.

Before proceeding, we emphasize two important points. First, we keep the model as simple as possible to make our main argument—information production in bad times leads to increased forecast accuracy of PDs—as clearly as possible. Hence, the model will have little to say about ancillary issues we analyze in the paper, such as collateral and real estate prices, though the model can easily be adapted to incorporate these features. The main goal of the model is not to offer a new explanation for why banks produce more information in bad times but rather to show how information production affects the discrimination ability of PDs. Nonetheless, the way in which we generate higher information production incentives in bad times—lower expected cash flows for low-quality borrowers—seems fairly reasonable. Second, there are some ingredients in the model that we include to fit it into our empirical framework that would not be necessary in a model purely meant to generate intuitions (e.g., different classes of borrowers and three types of borrowers rather than two).

B.1. Setup

There is a single borrower seeking funds from a bank at $t = 0$ for a project that yields a random payoff at $t = 1$. The borrower and bank are risk-neutral, and there is no discounting. At $t = 0$, there is a publicly observable aggregate state $\omega \in \{H, L\}$ (High or Low) which represents current economic conditions. The borrower belongs to a publicly observable class $\alpha \in [\underline{\alpha}_\omega, \bar{\alpha}_\omega]$ where $\underline{\alpha}_\omega > 0$ and $\bar{\alpha}_\omega < 1$, which is distributed according to the density function $f_\omega(\alpha)$, where the distribution potentially depends on the aggregate state. The borrower's class can be thought of as the borrower's public credit score based on observable characteristics.²⁹ Within each class of borrowers, there are three types of borrowers $\theta \in \{G, M, B\}$ (Good, Medium or Bad) where θ is initially unknown to all, and each type is equally likely.³⁰

The borrower has an investment opportunity that requires an initial investment of I , which we normalize to 1, at $t = 0$ and yields a cash flow at $t = 1$ of $R_\omega^\theta(\alpha) > 0$ if it succeeds and 0 if it fails, where the cash flow in the case of success can depend on the state, class of the borrower, and its unobservable type. We assume that the cash flows in the case of success are increasing in the quality of the borrower, i.e., $R_\omega^G(\alpha) \geq R_\omega^M(\alpha) \geq R_\omega^B(\alpha)$, for $\omega \in \{H, L\}$ and $\alpha \in [\underline{\alpha}, \bar{\alpha}]$. The average probability of success within class α of borrowers is α . If the borrower

²⁹As will become clear below, the purpose of having a borrower class is such that we obtain multiple probabilities of default to derive a receiver operating characteristics curve based on the model. We assume the class is continuously distributed for analytic tractability; however, we have also analyzed a discrete version which is available upon request.

³⁰As shown below, by having three borrower types, the bank produces information and screens out the bad borrower, but still has improved information among those that it ultimately lends to.

is good ($\theta = G$), the probability of success is $\alpha + \epsilon$; if the borrower is medium ($\theta = M$), the probability of success is α , and if the borrower is bad ($\theta = B$), the probability of success is $\alpha - \epsilon$.

The borrower offers the bank a loan contract that raises 1 at $t = 0$ to finance the investment and promises to repay F , which is endogenously determined, at $t = 1$. The firm has limited liability, so if the project's realized cash flow is lower than F , the firm defaults, and the bank collects the realized cash flow. Although the borrower's type θ is initially unknown, the bank can pay a cost $c > 0$ to learn θ . The potential value of information for the bank is to screen out lower-quality borrowers.³¹

B.2. Information Production

There are several ways we could model higher information production incentives in the low state. However, one that is natural and convenient for modeling purposes is simply to assume that the project is always ex-ante NPV positive in the high state and information production is unprofitable for the bank, while the project is ex-ante NPV negative in the low state absent information production. Formally,

Assumption 1. *The project is ex-ante NPV positive in the high state and NPV negative in the low state regardless of its class. It is unprofitable for the bank to produce information in the high state.*

1. $\frac{1}{3} [(\alpha + \epsilon)R_H^G(\alpha) + \alpha R_H^M(\alpha) + (\alpha - \epsilon)R_H^B(\alpha)] \geq 1 \quad \forall \alpha \in [\underline{\alpha}, \bar{\alpha}]$,
2. $\frac{1}{3} [(\alpha + \epsilon)R_L^G(\alpha) + \alpha R_L^M(\alpha) + (\alpha - \epsilon)R_L^B(\alpha)] < 1 \quad \forall \alpha \in [\underline{\alpha}, \bar{\alpha}]$,
3. $\frac{\epsilon}{3\alpha} < c \quad \forall \alpha \in [\underline{\alpha}, \bar{\alpha}]$.

We also make the following assumptions:

Assumption 2.

1. $\frac{1}{3} ((\alpha + \epsilon)R_L^G(\alpha) + \alpha R_L^M(\alpha)) - \frac{2}{3} \geq c \quad \forall \alpha \in [\underline{\alpha}, \bar{\alpha}]$,
2. $c \geq \frac{1}{3} (\alpha R_L^M(\alpha) - \frac{\alpha + \epsilon}{\alpha}) \quad \forall \alpha \in [\underline{\alpha}, \bar{\alpha}]$,
3. $R_L^M(\alpha) \geq \frac{1}{\alpha} \quad \forall \alpha \in [\underline{\alpha}, \bar{\alpha}]$,
4. $R_H^B(\alpha) \geq \frac{1}{\alpha} \quad \forall \alpha \in [\underline{\alpha}, \bar{\alpha}]$.

The purpose of this second set of assumptions will become clearer in the proof below. However, Assumption 2.1 ensures that producing information and lending to the good and medium borrower is NPV positive in the low state, Assumption 2.2 ensures that if the bank

³¹In a more complicated bargaining game, there could be value from being able to adjust interest rates. This force is not present here because the firm simply makes a take-it-or-leave-it offer.

produces information in the low state, it lends to both the good type and the medium type, not just the good type, and Assumption 2.3 ensures that this is feasible.³² Finally, Assumption 2.4 is not strictly necessary but simplifies the expressions by guaranteeing the borrower never defaults in the high state when the project succeeds.³³

We next characterize the equilibrium.

Proposition 1. *The bank does not produce information in the high state and lends to the borrower regardless of the type and class of the borrower. The bank always produces information in the low state and lends to the good and medium borrowers regardless of their class. Expected lending volume is lower in the low state than in the high state.*

Proof. Throughout the proof, we can consider a fixed α ; however, the proof applies to all $\alpha \in [\underline{\alpha}, \bar{\alpha}]$. First, consider the case in which the state is high. Assuming the bank does not produce information, the participation constraint $\alpha F \geq 1$ must hold. Because the firm has all of the bargaining power, it can offer the bank a zero profits contract with face-value $F_H = \frac{1}{\alpha}$. We can then check whether the bank would have incentives to produce information and only lend to the medium and good borrower if offered this contract:

$$\frac{2}{3} \left(\left(\alpha + \frac{\epsilon}{2} \right) F_H \right) - c > 0, \quad \implies \quad \frac{\epsilon}{3\alpha} > c. \quad (2)$$

However, Assumption 1.3 implies that (2) is violated. Hence, the bank would not produce information and only lend to the medium and good borrower if offered the zero-profits contract. We can also rule out the possibility that the bank produces information and only lends to the good borrower:

$$\frac{1}{3} ((\alpha + \epsilon) F_H - 1) > c, \quad \implies \quad \frac{\epsilon}{3\alpha} > c.$$

Hence, the firm offers the bank a contract with face-value F_H and the bank does not produce information and lends to the firm regardless of its type. Notice also that Assumption 2.4 implies that this contract is always feasible because the firm will never default if the project succeeds.

Now consider the low state. Because of Assumption 1.2, the bank will not lend without producing information. There are two potential alternatives. First, the bank could produce information and only lend to the good type. If this were the case, the firm would offer the following zero-profits face-value of debt $F'_L = \frac{1+3c}{\alpha+\epsilon}$, and the firm's profits would be:

$$\frac{\alpha + \epsilon}{3} \left(R_L^G(\alpha) - \frac{1 + 3c}{\alpha + \epsilon} \right). \quad (3)$$

³²While this assumption is not strictly necessary from an economic intuition perspective, we need the bank to lend to at least two types of borrowers within each class to generate an improvement in discrimination from information production.

³³Note that this assumption directly implies the project is NPV positive in the high state (i.e., Assumption 1.1).

Notice that from Assumption 2.3 this contract is feasible since $R_L^G(\alpha) > R_L^M(\alpha) \quad \forall \alpha \in [\underline{\alpha}, \bar{\alpha}]$.

Alternatively, the firm could offer a contract that induces the bank to lend to both the medium and good borrower. This contract cannot have zero profits, as this would mean the bank earns negative profits from the medium type and would not be willing to lend. Here, it must be incentive-compatible for the bank to lend to the medium type, i.e., $\alpha F \geq 1$. Again, since the firm has all of the bargaining power, we can consider a face-value of debt such that this constraint binds, i.e., $F_L'' = \frac{1}{\alpha}$. Notice again that Assumption 2.3 ensures that this loan contract is feasible. Under this contract, the firm's profits would be:

$$\frac{1}{3} \left((\alpha + \epsilon)(R_L^G(\alpha) - F_L'') + \alpha(R_L^M(\alpha) - F_L'') \right). \quad (4)$$

To check which contract the firm offers the bank, we need to compare (4) to (3). Subtracting (3) from (4) we have:

$$c + \frac{1}{3} \left(\alpha R_L^M(\alpha) - \frac{\alpha + \epsilon}{\alpha} \right),$$

which is positive from Assumption 2.2. Hence, the firm earns higher profits from offering a loan with face value F_L'' than F_L' , and so the bank produces information and lends to both the good and medium borrowers. Finally, based on these lending decisions, the expected lending volume is 1 in the high state and $\frac{2}{3}$ in the low state. □

B.3. Area Under the Curve (AUC) Derivation

In this section, we derive the AUC from the model and analyze its properties. First, it will also be useful to convert the class of borrower α into a corresponding probability of default: $p \equiv 1 - \alpha$. The corresponding density function is then $g_\omega(p) \equiv f_\omega(1 - \alpha)$ with support $[\underline{p}_\omega, \bar{p}_\omega]$ where $\underline{p}_\omega \equiv 1 - \bar{\alpha}_\omega$ and $\bar{p}_\omega \equiv 1 - \underline{\alpha}_\omega$.

The bank's perceived probability of default is:

$$\hat{PD} = \begin{cases} p & \text{if } \omega = H, \\ p - \epsilon & \text{if } \omega = L, \theta = G, \\ p & \text{if } \omega = L, \theta = M, \\ p + \epsilon & \text{if } \omega = L, \theta = B. \end{cases}$$

Notice that because the bank produces information in the low state, its perceived probability of default is more precise than in the high state.

Henceforth, we will assume that the bank interacts with an infinite number of borrowers independently drawn from both the class distribution and the ex-ante unobservable type distribution within each class. The receiver operating characteristics curve (ROC) plots the true

positive rate against the false positive rate. Specifically, for a given threshold t , the ROC considers any probability of default larger than t , i.e., $\hat{PD} > t$ a predicted positive and any probability of default less than t , i.e., $\hat{PD} \leq t$, a predicted negative. A predicted positive is a true positive if the borrower actually defaults and a false positive if it does not. A predicted negative is a true negative if the borrower does not default and a false negative if it does default. The true positive rate is equal to the ratio of true positives to total positives and in the high state is:

$$\begin{aligned} TPR_H(t) &= \frac{\int_t^{\bar{p}_H} \left(\frac{1}{3}(p - \epsilon) + \frac{1}{3}p + \frac{1}{3}(p + \epsilon) \right) f_H(p) dp}{\int_{\underline{p}_H}^{\bar{p}_H} \left(\frac{1}{3}(p - \epsilon) + \frac{1}{3}p + \frac{1}{3}(p + \epsilon) \right) f_H(p) dp} \\ &= \frac{\int_t^{\bar{p}_H} p f_H(p) dp}{\int_{\underline{p}_H}^{\bar{p}_H} p f_H(p) dp}. \end{aligned}$$

The false positive rate is equal to the ratio of false positives over true negatives, which in the high state is:

$$FPR_H(t) = \frac{\int_t^{\bar{p}_H} (1 - p) f_H(p) dp}{\int_{\underline{p}_H}^{\bar{p}_H} (1 - p) f_H(p) dp}.$$

Because the probability of a true positive or false positive only depends on the average realized default rate for a given PD, the ϵ term drops out in both expressions given that it has a mean of zero.

The receiver operating characteristics curve plots the true positive rate against the false positive rate, and the AUC is the area under this curve. Given the continuous distribution, we can write the AUC in the high state as:

$$AUC_H = \int_{\underline{p}_H}^{\bar{p}_H} TPR_H(t) |FPR'_H(t)| dt, \quad (5)$$

where $|FPR'_H(t)|$ denotes the absolute value of the derivative of FPR_H with respect to t . In the low state, the true positive rate is as follows:

$$TPR_L(t) = \begin{cases} \frac{\frac{1}{2} \int_{t+\epsilon}^{\bar{p}_L} (p - \epsilon) f_L(p) dp + \frac{1}{2} \int_{\underline{p}_L}^{\bar{p}_L} p f_L(p) dp}{\int_{\underline{p}_L}^{\bar{p}_L} \left(p - \frac{\epsilon}{2} \right) f_L(p) dp} & \text{if } t < \underline{p}_L, \\ \frac{\frac{1}{2} \int_{t+\epsilon}^{\bar{p}_L} (p - \epsilon) f_L(p) dp + \frac{1}{2} \int_t^{\bar{p}_L} p f_L(p) dp}{\int_{\underline{p}_L}^{\bar{p}_L} \left(p - \frac{\epsilon}{2} \right) f_L(p) dp} & \text{if } t \in [\underline{p}_L, \bar{p}_L - \epsilon], \\ \frac{\frac{1}{2} \int_t^{\bar{p}_L} p f_L(p) dp}{\int_{\underline{p}_L}^{\bar{p}_L} \left(p - \frac{\epsilon}{2} \right) f_L(p) dp}, & \text{if } t > \bar{p}_L - \epsilon, \end{cases}$$

and the false positive rate is as follows:

$$FPR_L(t) = \begin{cases} \frac{\frac{1}{2} \int_{t+\epsilon}^{\bar{p}_L} (1-p+\epsilon) f_L(p) dp + \frac{1}{2} \int_{\underline{p}_L}^{\bar{p}_L} (1-p) f_L(p) dp}{\int_{\underline{p}_L}^{\bar{p}_L} (1-p+\frac{\epsilon}{2}) f_L(p) dp} & \text{if } t < \underline{p}_L, \\ \frac{\frac{1}{2} \int_{t+\epsilon}^{\bar{p}_L} (1-p+\epsilon) f_L(p) dp + \frac{1}{2} \int_t^{\bar{p}_L} (1-p) f_L(p) dp}{\int_{\underline{p}_L}^{\bar{p}_L} (1-p+\frac{\epsilon}{2}) f_L(p) dp} & \text{if } t \in [\underline{p}_L, \bar{p}_L - \epsilon], \\ \frac{\frac{1}{2} \int_t^{\bar{p}_L} (1-p) f_L(p) dp}{\int_{\underline{p}_L}^{\bar{p}_L} (1-p+\frac{\epsilon}{2}) f_L(p) dp} & \text{if } t > \bar{p}_L - \epsilon. \end{cases}$$

In contrast to the high state, the ϵ term is present in the expressions for the TPR and FPR because the bank is only lending to good and medium borrowers. Finally, the area under the curve in the low state is then:

$$AUC_L = \int_{\underline{p}_L - \epsilon}^{\bar{p}_L} TPR_L(t) |FPR'_L(t)| dt.$$

B.4. Analysis of AUC

In general, the expressions for AUC are not easily expressed analytically. However, we can solve for the AUC in closed form for both the low and high states assuming a uniform distribution of the class of borrowers. Specifically, suppose that $p \sim U[\underline{p}_\omega, \bar{p}_\omega]$, where $\underline{p}_\omega - \epsilon \geq 0$, $\bar{p}_\omega + \epsilon \leq 1$ and $\bar{p}_\omega - \epsilon \geq \underline{p}_\omega$ for $\omega \in \{H, L\}$.³⁴

First, we analyze the case in which the distribution of borrower class is the same in the high and low state, i.e., $f_H(\alpha) = f_L(\alpha)$. Under the uniform distribution, this amounts to $\bar{p}_H = \bar{p}_L = \bar{p}$ and $\underline{p}_H = \underline{p}_L = \underline{p}$. In this case, so long as the average probability of default is below 50%, the AUC is always higher in the low state. Formally,

Proposition 2. *Assume the class of borrowers is the same in the high and low state, i.e., $f_H(\alpha) = f_L(\alpha)$ and follows a uniform distribution as described above. If the average probability of default is less than 50%, the AUC is always larger in the low state.*

Proof. If we subtract AUC_H from AUC_L and remove the ϵ term we have

$$\frac{4(\bar{p} - \underline{p})^3(1 - \underline{p} - \bar{p}) + (\underline{p} - \bar{p})(\underline{p}^2 - (6 - \bar{p})\bar{p} + 2\underline{p}(5\bar{p} - 3))\epsilon - (2 - \underline{p} - \bar{p})(\underline{p} + \bar{p})\epsilon^2}{6(\underline{p} - \bar{p})^2(2 - \underline{p} - \bar{p})(\underline{p} + \bar{p})(2 - \underline{p} - \bar{p} + \epsilon)(\underline{p} + \bar{p} - \epsilon)}$$

Since the denominator is clearly positive, it suffices to show the numerator is positive. The first term of the numerator is positive because $\bar{p} - \underline{p} > 0$ and $\underline{p} + \bar{p} < 1$. Thus, it is sufficient to show that the remaining terms are positive. Hence, we need:

$$-(\bar{p} - \underline{p})(\underline{p}^2 + \bar{p}^2 - 6(\underline{p} + \bar{p}) + 10\underline{p}\bar{p}) > (2 - (\underline{p} + \bar{p}))(\underline{p} + \bar{p})\epsilon$$

³⁴Note this is equivalent to $\alpha \sim U[1 - \bar{p}_\omega, 1 - \underline{p}_\omega]$.

Since $\bar{p} - \epsilon > \underline{p}$, then $\bar{p} - \underline{p} > \epsilon$ and since $(2 - (\underline{p} + \bar{p}))(\underline{p} + \bar{p}) > 0$, it suffices to show

$$\begin{aligned} & -(\bar{p} - \underline{p})(\underline{p}^2 + \bar{p}^2 - 6(\underline{p} + \bar{p}) + 10\underline{p}\bar{p}) > (2 - (\underline{p} + \bar{p}))(\underline{p} + \bar{p})(\bar{p} - \underline{p}) \\ \iff & (8\underline{p}\bar{p} - 4(\underline{p} + \bar{p}))(\bar{p} - \underline{p}) < 0 \end{aligned}$$

This term is positive because $\bar{p} - \underline{p} > 0$ and both \bar{p} and \underline{p} are less than one. Hence, the entire term is positive. \square

Proposition 2 says that so long as the average perceived PDs are below 50%, the AUC is always higher in the low state. Why does an average of 50% matter? The reason is that while the improved discrimination across borrowers from information production raises the AUC, there is also a second effect due to the change in the average probability of default. To best understand this, notice that even in the absence of information production, the AUC in the high state depends on the average risk of borrowers:

$$AUC_H = \frac{1}{6} \left(3 + \frac{2(1 - \underline{p})}{2 - \underline{p} - \bar{p}} - \frac{2\underline{p}}{\underline{p} + \bar{p}} \right). \quad (6)$$

Suppose we fix the range between the upper and lower bound of the uniform distribution, i.e., $\delta \equiv \bar{p} - \underline{p}$, and substitute \underline{p} with $\bar{p} - \delta$ into the expression for AUC_H , then we have:

$$\frac{1}{6} \left(3 - \frac{2\bar{p}}{\delta - 2\bar{p}} - \frac{2(1 - \bar{p})}{2 - 2\bar{p} + \delta} \right). \quad (7)$$

Differentiating (7) with respect to \bar{p} and substituting back in \underline{p} we have:

$$-\frac{4(\bar{p} - \underline{p})(1 - \bar{p} - \underline{p})}{3(2 - \bar{p} - \underline{p})^2(\bar{p} + \underline{p})^2},$$

which is positive if $\bar{p} + \underline{p} > 1$ and negative otherwise. Hence, shifting the distribution of PDs upwards while fixing the distance between the upper and lower bound of the distribution increases the AUC if the average probability of default is above one-half. Because the bank produces information and screens out bad borrowers in the low state, the distribution of PDs shifts downwards, causing a reduction in the average PD. When $\bar{p} + \underline{p} > 1$, it is possible for information production to decrease the AUC. The intuition for this is as follows: when the average PD is close to 50%, uncertainty is highest.³⁵ Probabilities close to 0 and 1 are more certain and hence provide clearer separation. For example, a difference in PDs of 1% versus 2% is much more discriminative than a difference of 49% versus 50%.

While this effect is theoretically possible, it is extremely unlikely to be driving our results for several reasons. First, this is a necessary, but not sufficient, condition for the AUC to decrease from information production. There are still parameters (particularly when ϵ is large

³⁵As discussed below, this is also true, although to a much larger extent, when assessing forecast errors.

enough) such that even when the average PD is above 50%, the AUC is higher in the low state. Second, PDs in the data average less than 2% unconditionally, and hence, are very far from 50%. Lastly, in the data, we do not actually observe lower PDs in bad times. Rather, they are quite similar and, if anything, slightly larger (1.66% versus 1.61%), which is likely due to the distribution of underlying borrowers also changing in bad times. Indeed, we next show that if the distribution of potential borrowers has higher PDs in the low state, but the average PD among borrowers who actually receive loans is the same in the low and high state, which is approximately what we see in the data, the AUC is always higher.

Proposition 3. *Suppose that $p \sim U[\underline{p}, \bar{p}]$ in the high state and $p \sim U[\underline{p} + \frac{\epsilon}{2}, \bar{p} + \frac{\epsilon}{2}]$ in the low state. While the average PD of potential borrowers is higher in the low state, the average PD of firms that actually receive financing is the same, and the AUC is always higher in the low state.*

Proof. First, note that the bank lends to all borrowers in the high state from Proposition 1. Hence, the mean of both the true PD and perceived PD is $\frac{1}{2}(\underline{p} + \bar{p})$, which is the same as in the low state because the borrower only lends to the medium and good borrowers. If we subtract AUC_H from AUC_L we have:

$$\frac{\epsilon^2(3\bar{p} - 3\underline{p} + \epsilon)}{6(\bar{p} - \underline{p})^2(2 - \bar{p} - \underline{p} + \epsilon)(\bar{p} + \underline{p} - \epsilon)},$$

which is strictly positive. Hence, the AUC is always higher in the low state than the high state when the average PD of granted loans is the same in both states. \square

This result suggests that it is the improved discrimination in bad times that is driving the increase in the AUC in our results, not the fact that information production changes the average quality of borrowers.³⁶ Finally, since PDs are slightly higher in bad times in our data, if anything, this should go against our main result.

A higher variance of underlying PDs generally raises the AUC, even in the absence of differences in information production. To see this, suppose we fix the average PD in the high state for the uniform case, i.e., $\sigma \equiv \bar{p} + \underline{p}$ and replace \underline{p} with $\sigma - \bar{p}$ into (6), then we have: $\frac{4\bar{p}+4\sigma-3\sigma^2}{12\sigma-6\sigma^2}$, which is increasing in \bar{p} . This implies that a higher variance of PDs increases the AUC in the uniform case. In our data, the standard deviation of PD for new loans is slightly higher in periods of high unemployment (2.88pp versus 2.58pp). However, it is highly unlikely that this effect can quantitatively explain our results. For instance, if we assume PDs follow a uniform distribution, we can solve for \bar{p} and \underline{p} given the mean and variance that we observe in the data in low and high unemployment periods. We can then plug these values into the AUC in the high state without information production (5). Doing so gives us an AUC of 0.542 and

³⁶Relatedly, under the uniform case it can easily be shown that if the class distribution is the same in both states, but the bank exogenously receives information in the low state and lends to all borrowers, the AUC is always higher in the low state.

0.551 in the high and low states (a difference of 0.009). Hence, it is unlikely that this small difference in variance mechanically explains all of our results. For example, in Figure 4, we find a difference in AUCs of over four times as large (0.043) among newly issued loans.

Although we cannot solve the AUC analytically for other distributions, we can do a similar exercise as above and simulate the data assuming the class of borrowers follows a lognormal distribution³⁷, matching the mean and variance of PDs to what we observe in the data in periods of high and low unemployment. Figure B.1 displays the distribution of the difference between the low and high state AUC over 100 simulations. The average difference is 0.011, and the maximum difference is 0.042, which is less than the 0.043 difference we find in the data.³⁸ Hence, mechanical differences in the mean and variance of PD in periods of high unemployment are unlikely to explain our main results. In Online Appendix Section C.4, we also show how mechanical differences in the distribution of PD are unlikely to explain our cross-sectional results comparing information quality across different types of loans.

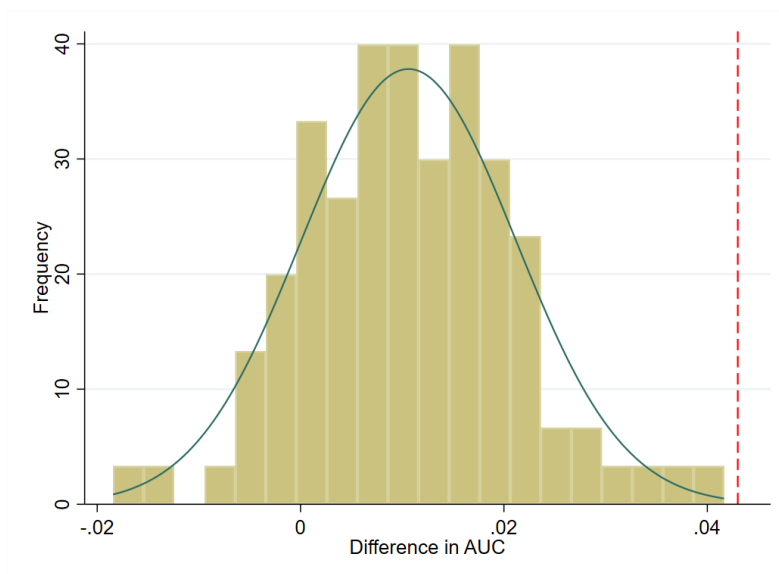


Figure B.1: Difference in AUCs for simulated lognormally distributed PDs

This figure plots the simulations of the difference in AUCs across periods of high and low unemployment, assuming no information production and PDs are lognormally distributed. The number of observations in each simulation is equal to that in our main sample (104,111). The parameters of each distribution are estimated to match the means and variances we observe in the data. The sample average of the difference in AUCs is 0.011. The dashed red line at 0.043 marks the difference in AUC between high and low unemployment periods that we estimate in the data.

³⁷In Figure 1, $\log(\text{PD})$ approximates a normal distribution.

³⁸It is also comforting that the average difference under a lognormal distribution (0.011) is fairly similar to the difference when the class of borrowers follows a uniform distribution (0.009).

Appendix C. Additional Analysis For Online Publication.

In this section, we include additional robustness tests, analyses, and discussions regarding our main results.

C.1. Additional Descriptive Statistics and Collinearity Tests

In Table OA.1, we report summary statistics for the broader sample of loans used when analyzing the origination unemployment rate in Section 5.1. The Y-14Q data include quarterly information about all eligible loans on banks' balance sheets starting in 2014Q4, regardless of when the loan was originated. This expanded sample thus includes not only subsequent observations of the new loans included in Table 1 but also post-2014 observations of loans that were originated prior to the collection of the Y-14Q data.

Table OA.2 reports correlations between PD, interest rate, and firm or loan characteristics for our baseline sample of new loans. The first column shows that the magnitudes of the correlations between PD and control variables are small, suggesting collinearity is not a major issue. Table OA.3 formally supports this interpretation using a statistical measure of collinearity known as the variance inflation factor (VIF). For each explanatory variable, we calculate the VIF as follows. First, we estimate a regression where the independent variable is the two-year default indicator used throughout our analysis, and the independent variables are shown in the first column of Table OA.3 (the “full” model). Next, we estimate separate regressions including only one of the explanatory variables at a time (the “separate” models). For each independent variable, the VIF is calculated as the ratio of the variance of the parameter estimate in the full model to the variance of the parameter estimate in the separate model. If an independent variable has zero correlation with every other independent variable, then the VIF will be one. High VIF values—a common rule of thumb (see O'Brien (2007)) is to use a VIF threshold of 10—indicate that the uncertainty around a parameter estimate is much greater when other explanatory variables are included, indicating the presence of collinearity. Since the VIF for PD is just 1.11, this suggests that PD cannot simply be replicated by a linear combination of other firm or loan characteristics.

C.2. Further Evidence that PD Contains Private Information

In this section, we discuss a range of additional robustness checks to verify that 1) PD contains information useful for predicting default and 2) banks use this information to set the price and terms of their loans.

Throughout the main paper, we test the ability of PD to predict default within the following two years. In Table OA.4, we show that PD also predicts other measures of loan nonperformance. In each column, we regress the indicator reported in the top row on PD along with the same set of firm/loan controls and fixed effects as in our baseline regression (1). Columns

(1) - (3) use alternative definitions of default with different default horizons. Column (1) uses whether a loan defaults at any point in our sample, Column (2) uses the annualized average default rate over the life of the loan, and Column (3) uses the one year ahead default rate. Columns (4) and (5) use indicators that equal one if the loan is reported as delinquent or the bank records a charge-off for that loan within the following two years as dependent variables. Across these alternative measures, PD remains a statistically and economically significant predictor, suggesting that our main results do not depend on how loan performance is measured.

Next, we test how the information contained in PD is used to set loan terms. Holding all else equal, a bank should charge a higher interest rate on loans for which it believes the borrower to be at higher risk of default. We test this formally in Table OA.5, which regresses the interest rate on different combinations of firm and loan controls and fixed effects. Column (1) includes only observable characteristics, Column (2) includes only banks' risk assessments (PD and LGD), and Column (3) includes both. The results suggest that a one percentage point increase in PD leads to a statistically significant increase of about 9 basis points in the interest rate, confirming that banks use the information contained in PD when pricing their loans.

To further verify that the information contained in PD is used to set interest rates, we perform a modified version of the exercise shown in Table 4 of the main paper. We first use several different random forest regression specifications to generate predicted interest rates based on observables using the training sample. Then, using the validation sample, we regress actual interest rates on both 1) the interest rates predicted by the random forest regressions and 2) PD. The results are shown in Table OA.6. Across all specifications, a higher PD is estimated to increase interest rates by a similar magnitude to the estimates shown in Table OA.5.

Having shown that the information contained in PD is reflected in interest rates charged at origination, we next test whether changes in PD for existing loans are accompanied by changes in other loan terms. We follow Bidder et al. (2023) and focus on changes in either the maturity date or the interest rate,³⁹ as they show that changes in other loan terms are generally accompanied by changes in one of these two measures. In Table OA.7, we regress dummy variables indicating whether a loan had a change in its maturity date (first four columns), its interest rate (middle four columns), or any combination of the two (last four columns) relative to the prior quarter. The primary coefficient of interest, which we report in the top row, is a dummy variable equal to one if a loan's reported PD changed from the prior quarter. The top row implies that a bank changing a loan's PD increases the probability of a modification at that same time by around 0.1pp. The effects we estimate are small relative to the unconditional quarterly probability of a loan modification reported in Bidder et al. (2023) of roughly 10% because we include bank-time, industry-time, and bank-county fixed effects in all specifications.

³⁹Because the interest rate for floating rate loans will mechanically change with the underlying reference rate, we define the indicator variable for a change in the interest rate as follows. For floating rate loans, we say that the interest rate was modified at time t if the reported interest rate spread was not equal to the spread reported at time $t - 1$. For fixed-rate loans, we say that the interest rate was modified at time t if the interest rate was different from both the initially reported rate and the rate reported at time $t - 1$.

Figure OA.1 shows the dynamic effects of changes in PD over the following year for a firm-bank relationship. For these regressions, we use all new and existing loans and collapse the data to the firm-bank-quarter level. While many loan terms (such as the committed amount) are fixed at the loan level, banks can use newly produced information to affect the terms of *new* loans to the same borrower. We use firm-time fixed effects (e.g., Khwaja and Mian (2008)) to compare differences in bank information and lending amounts across banks for the same firm at the same time. The top panel shows that a 1pp increase in the PD assessment of bank b for firm i decreases total lending across all loans from b to i by about 0.1% for the current and subsequent quarters. Consistent with this, the bottom-left panel shows that a firm becomes roughly 0.05pp less likely to receive a new loan from a bank over this same horizon. Lastly, the bottom-right panel shows that the probability of a firm-bank match no longer being observed in the data—which we view as a proxy for the dissolution of that firm-bank relationship—increases in subsequent quarters. Taken together, these results provide further confirmation that the information contained in banks’ PDs also affects their lending behavior.

Finally, we compare the ability of PD and the interest rate to predict default in Figure OA.2. The first three panels plot ROC curves comparing the discriminatory ability of PD to interest rates. The top panels compare the level of PD to the level of the interest rate (left) or interest rate spread (right), while the bottom-left panel compares within-quarter percentile ranks of PD and IR. In all three cases, the AUC is much higher for PD, and these differences are all statistically significant. The bottom-right performs a similar exercise to Figure A.2 in Appendix B by testing the marginal predictive power of including the interest rate in a random forest regression model using only firm and loan characteristics. Unlike in Figure A.2, where the addition of PD leads to a large and statistically significant improvement in the AUC, the inclusion of the interest rate has only a small effect on the ability of other observables to predict default.⁴⁰

C.3. Robustness of Countercyclical Information Quality

This section includes additional robustness tests that support the results in Section 4.2, where we show that bank information quality is countercyclical.

One potential concern is that the median number of new loans in a county-quarter is zero, indicating the presence of a large number of counties that receive only a small number of loans during our sample period (see Panel C of Table 1). To the extent that we do not observe sufficient variation in the economic conditions at origination for these counties, this could potentially distort our results. To verify that this is not the case, in Figure OA.3, we repeat our baseline analysis using only counties with at least 5 (left panel) or 15 (right panel) new loans throughout our sample period. In both cases, the differences in information quality across periods of high and low unemployment remain statistically significant, and the magnitudes are

⁴⁰This may seem surprising given that interest rates are affected by PDs; however, interest rates are also affected by other factors unrelated to borrower risk such as market power and banks’ cost of capital.

even larger than our baseline results shown in Figure 4. This suggests that our main findings are not driven by the inclusion of small counties with few loans.

To verify that our results are representative across banks, in Table OA.8, we separately estimate differences in AUC for each of the 22 banks in our sample (out of the 33 total banks in the data) with at least 10,000 total loan-quarter observations. The left panel splits loans by current UR, and the right panel splits them by origination UR. In both cases, the mean of the difference in AUCs across banks is similar to the estimates obtained from pooling loans across all banks shown in Figure 6.⁴¹ In Table OA.8, we report the number of banks that are estimated to have statistically significant differences in AUC across periods of high and low unemployment, along with the share of loans and loan volume coming from these banks. This table shows that roughly half of the banks in our sample demonstrate statistically significant countercyclical information quality and that these banks account for more than half of all lending volume. These exercises also provide further reassurance that our main findings are not distorted by our decision to pool together loans across banks.⁴²

Another potential concern is that our sample period includes loans made during the Covid pandemic. While the turbulent conditions observed during this period can provide potentially useful identifying variation in county-level unemployment rates, the loans made during this time may not be representative given the widespread financial market interventions in place at the time. To show that our results do not depend on the inclusion of this period, Figure OA.4 reports our main result using only loans made through the end of 2019. The difference in AUCs is statistically significant and even larger than our baseline results, suggesting that the cyclicity of information quality is not purely driven by the exceptional economic circumstances following Covid.

C.4. Further Evidence of Information Production

In this section, we provide further evidence that our results are driven by countercyclical information production by banks. In Section 5, we show that the cyclicity of information quality is concentrated in new loans, larger loans, and loans with higher expected losses. However, as we discuss in Appendix B.4, differences in the AUC can result from mechanical changes in the underlying distributions of PDs. For example, if the variance of the true underlying PDs increases more for large loans during recessions than for small loans, this would mechanically increase the AUC for large loans relative to small loans during downturns, even if banks did not change their information production. To address this concern, we report the differences in the distributions of PDs over the business cycle for different types of loans in Figures OA.5 - OA.7. We then perform a similar exercise as in Section B.4, where we estimate the model-implied AUC without any information production, assuming the distribution of PDs has the

⁴¹In that figure, ΔAUC is 0.024 when split by current UR and 0.038 when split by origination UR.

⁴²For example, changes in PDs that do not alter the relative ordering of PDs at the bank level could cause a change in the relative ordering of the pool of PDs across banks.

same mean and variance that we observe in the specific subsample of the data.⁴³ We estimate the AUCs analytically, assuming PDs follow a uniform distribution, and based on simulations assuming PDs follow a lognormal distribution. For the lognormal case, we estimate each AUC 100 times based on simulations of 25,000 PDs and report the sample average AUC.

Our main cross-sectional results in Figures 4, 6, 7, and 8 can be summarized as follows. We find that 1) information quality is more cyclically sensitive for new loans, and 2) information quality is higher and more cyclically for loans with high EADs and expected losses. While there are meaningful differences in the underlying distribution of PDs for new versus existing loans and high-EL versus low EL-loans, the bottom panels of Figures OA.5 - OA.7 show that none of the abovementioned cross-sectional results can be explained mechanically by differences in the underlying distribution of PDs.⁴⁴ Taken together, these results provide further reassurance that our main results are not driven by mechanical changes in the distributions of banks' reported PDs over the business cycle.

Next, we test whether realized default rates respond differently to the business cycle for different types of loans. In Table OA.9, we regress a default indicator on our unemployment indicator interacted with indicators for whether a loan is new and whether its exposure at default (EAD), loss given default (LGD), or expected losses (EL) were above the median for that bank-quarter. The interactions with the high-UR dummy are not statistically significant for the new loan indicator, LGD, and EL, suggesting that the default rates for new loans and loans with larger potential losses do not differentially respond when the local unemployment rate increases. And while the interaction coefficient between the high-UR dummy and EAD is statistically significant, the sign of the coefficient is negative, suggesting that default rates for larger loans are *less* cyclically sensitive than default rates for smaller loans.

Next, as we describe in Section 2, we exclude government guaranteed loans from our baseline analysis. Here, however, we use these loans to provide further evidence in support of theories of endogenous bank information production. A bank should have diminished incentives for producing information about a borrower with a government guarantee. Intuitively, that information is not going to be particularly valuable because it will have little effect on the bank's expected losses. Thus, information quality for loans with guarantees should be both lower on average and less cyclically sensitive.

We test this in Figure OA.8. The left panel shows that the AUC for the new loans used in our main analysis (blue line) is larger than the government guaranteed loans that we exclude (red line). The difference of 0.104 is statistically significant and represents more than twice the magnitude of the difference in our baseline results in Figure 4, confirming that information quality is lower for guaranteed loans. The right panel splits this set of guaranteed loans based on whether they are originated in periods of high or low unemployment and finds no statisti-

⁴³The AUC without information production is equivalent to the AUC in the high state in the model (5).

⁴⁴That we observe differences in the underlying distribution of PDs for these subsamples of loans is not surprising because older loans are more likely to become close to default and high-EL loans have higher PDs to begin with.

cally significant difference between the two. While we believe that the sample of government guaranteed loans is too small to draw far-reaching conclusions, these results are consistent with endogenous information production.

Banks' information quality may depend not only on the conditions of borrowers but also on the bank's own condition. In particular, banks with lower capital ratios may be more concerned about borrower default, leading them to produce more information. We test this prediction in Figure OA.9, which splits our baseline sample of new loans based on whether the total risk-based capital ratio for each issuing bank was above or below the median of all banks in our data in each quarter. We find that information quality is higher for new loans made by banks with lower capital ratios; however, this difference is only marginally statistically significant.

C.5. *Discussion of Why Regression Coefficient Deviates from One*

In this section, we discuss why the coefficient we estimate from regressing realized default on PD yields a coefficient of around 0.4 in Table 3 instead of one. Under rational expectations, the coefficient from regressing an outcome variable on its forecast should be equal to one (Muth (1961)); however, empirically, this is often not the case (e.g., Mincer and Zarnowitz (1969)). We offer several potential explanations for this in our setting.

The first is systematic biases in banks' PD assessments. For example, banks could be conservative in their estimates of PDs to avoid regulatory scrutiny. In fact, banks are encouraged by regulators to be conservative in their forecasts, particularly when they have little information:

Given the difficulties in forecasting future events and the influence they will have on a particular borrower's financial condition, a bank must take a conservative view of projected information. Furthermore, where limited data are available, a bank must adopt a conservative bias to its analysis.⁴⁵

A conservative bias can cause the average realized default rate to be lower than the average ex-ante PD, resulting a regression coefficient less than one. Importantly, and as discussed in more detail in Appendix Section B, a systematic conservative bias of this type would not affect the AUC so long as it did not change the relative ordering of PDs. More broadly, banks may have asymmetric loss functions that cause them to report forecasts that deviate systematically from their conditional expectation (e.g., Granger (1969)).

It is also possible that banks have incorrect priors or use a misspecified forecasting model. Moreover, if banks overweight their private information (e.g., Tversky and Kahneman (1974) and Grether (1980)), this can introduce noise in the forecast, which causes the coefficient to be attenuated due to measurement error. Finally, PD is a "long-run" average default rate, while we measure realized default only over a two-year window. Hence, banks' PDs do not necessarily

⁴⁵Source: Bank for International Settlements instructions for calculation of RWA for credit risk.

reflect their true conditional expectation over the following two years.⁴⁶

C.6. Biases in PD Reporting

Our model in Appendix B assumes that the AUC is calculated based on banks' actual beliefs about their borrowers and that these beliefs are unbiased. However, one concern could be that banks' reported PDs are biased. For example, banks may have incorrect priors, incentives to misreport PDs (e.g., Behn, Haselmann, and Vig (2016) and Plosser and Santos (2018)), or behavioral biases that cause them to underreact to information. One useful aspect of the AUC is that it is unaffected by these biases so long as they do not affect the relative ordering of loans' PDs. This can easily be seen by inspecting the terms for the TPR and FPR in the high and low states (see Section B.3). For a given threshold t , the true positive rate and false positive rate are determined by the actual probability of default of loans with PDs higher than threshold t . Hence, any transformation of the reported PDs that maintains the original ordering of the PDs does not affect the AUC.

Although it is certainly possible that biases may change the relative ordering of PDs, we do not see an obvious reason why that would entirely explain our empirical results. Moreover, in Figure OA.9, we find that the AUC is actually *higher* for banks with lower regulatory capital, suggesting their PDs are actually more accurate despite potentially higher incentives to manipulate them.

C.7. The Effect of Aggregate Shocks on AUC

In this section, we address the concern that defaults are driven by aggregate rather than idiosyncratic components in bad times and that this may have mechanical effects on the discriminatory ability of banks' PDs. To do so, we analyze the model developed in Appendix B.

Suppose that each firm's probability of default within its class is equal to $p + \xi$ where ξ is an aggregate shock distributed according to the density function $f_\omega(\xi)$, where ξ is unobservable ex-ante and independent of p .

For simplicity, we focus on the high state; however, the exact same logic also applies to the

⁴⁶PDs are also annual default rates, while we measure realized default over a two-year horizon. However, in Online Appendix Table OA.4 column (3), we find a coefficient of 0.33 when we use a one-year default indicator as the dependent variable instead.

low state. The true positive rate in the high state is as follows:

$$\begin{aligned}
TPR_H(t) &= \frac{\int \int_t^{\bar{p}_H} \frac{1}{2} \left(\frac{1}{3}(p + \xi) - \epsilon \right) + \frac{1}{3}(p + \xi) + \frac{1}{3}(p + \xi + \epsilon) f_H(p) f_\omega(\xi) dp d\xi}{\int \int_{\underline{p}_H}^{\bar{p}_H} \left(\frac{1}{3}(p + \xi) - \epsilon \right) + \frac{1}{3}(p + \xi) + \frac{1}{3}(p + \xi + \epsilon) f_H(p) f_\omega(\xi) dp d\xi} \\
&= \frac{\int_t^{\bar{p}_H} \frac{1}{2} \left(\frac{1}{3}(p + \mathbb{E}[\xi]) - \epsilon \right) + \frac{1}{3}(p + \mathbb{E}[\xi]) + \frac{1}{3}(p + \mathbb{E}[\xi] + \epsilon) f_H(p) dp}{\int_{\underline{p}_H}^{\bar{p}_H} \left(\frac{1}{3}(p + \mathbb{E}[\xi]) - \epsilon \right) + \frac{1}{3}(p + \mathbb{E}[\xi]) + \frac{1}{3}(p + \mathbb{E}[\xi] + \epsilon) f_H(p) dp} \\
&= \frac{\int_t^{\bar{p}_H} (p + \mathbb{E}[\xi]) f_H(p) dp}{\int_{\underline{p}_H}^{\bar{p}_H} (p + \mathbb{E}[\xi]) f_H(p) dp}.
\end{aligned}$$

Similarly, the false positive rate is:

$$FPR_H(t) = \frac{\int_t^{\bar{p}_H} (1 - p - \mathbb{E}[\xi]) f_H(p) dp}{\int_{\underline{p}_H}^{\bar{p}_H} (1 - p - \mathbb{E}[\xi]) f_H(p) dp}.$$

Hence, the aggregate shock only affects the AUC to the extent that it raises or lowers the average default rate. This is a direct implication of the fact that the probability of a true positive or false positive only depends on the average realized default rate for a given PD.⁴⁷

To see how the expected size of the aggregate shock affects the AUC, we calculate the AUC in the high state, assuming borrower class follows a uniform distribution, as follows:

$$\frac{1}{6} \left(3 + \frac{2(\bar{p} - \underline{p})}{(2 - \underline{p} - \bar{p} - 2\mathbb{E}[\xi])(\underline{p} + \bar{p} + 2\mathbb{E}[\xi])} \right)$$

Differentiating this with respect to $\mathbb{E}[\xi]$ we have:

$$-\frac{4(\bar{p} - \underline{p})(1 - \underline{p} - \bar{p} - 2\mathbb{E}[\xi])}{3(2 - \underline{p} - \bar{p} - 2\mathbb{E}[\xi])^2(\underline{p} + \bar{p} + 2\mathbb{E}[\xi])^2},$$

which is negative when the average PD is below 0.5, i.e., $1 - \underline{p} - \bar{p} - 2\mathbb{E}[\xi] > 0$. Hence, when PDs are distributed uniformly, the AUC decreases in $\mathbb{E}[\xi]$ for any reasonable parameters. Note that this follows exactly from the analysis in Section B.4 where we consider how changes in the average PD affect the AUC. Hence, if anything, higher average aggregate shocks should put downward pressure on the AUC we observe in the data in periods of high unemployment.

Finally, one may also wonder what happens if the bank observes ξ . As we show in Section B.4, a higher variance of PDs leads to a higher AUC. Hence, if ξ exhibits a higher variance in bad times, this would cause the AUC to be higher by increasing the variance of PDs.⁴⁸ However, as discussed earlier, in the data, there is only a slightly larger variance in PDs in bad times, and thus, this is highly unlikely to explain our results.

⁴⁷This is the same reason the ϵ terms drop out of the TPR and FPR in the high state.

⁴⁸This would also be true, although to a lesser extent if the bank observed the aggregate shock with noise.

C.8. Alternate Measure of Forecast Accuracy: Brier Score

In this section, we discuss an alternative measure of forecast accuracy, the Brier Score, which is the average of the squared forecast errors. For a sample of N borrowers indexed by i the Brier score (BS) is equal to:

$$BS = \frac{1}{N} \sum_i^N (\hat{P}D_i - d_i)^2,$$

where d_i is an indicator that equals one if firm i defaults. A lower Brier Score indicates a better probability prediction.

Following [Murphy \(1973\)](#), we can decompose the Brier Score as follows:

$$\begin{aligned} E[BS] &= \mathbb{E} \left[(\hat{P}D_i - d_i)^2 \right] \\ &= \mathbb{E} \left[(\hat{P}D_i - \mathbb{E}[d_i | \hat{P}D_i] + \mathbb{E}[d_i | \hat{P}D_i] - d_i)^2 \right] \\ &= \mathbb{E} \left[(\hat{P}D_i - \mathbb{E}[d_i | \hat{P}D_i])^2 + 2(\hat{P}D_i - \mathbb{E}[d_i | \hat{P}D_i])(\mathbb{E}[d_i | \hat{P}D_i] - d_i) + (\mathbb{E}[d_i | \hat{P}D_i] - d_i)^2 \right] \\ &= \mathbb{E} \left[(\hat{P}D_i - \mathbb{E}[d_i | \hat{P}D_i])^2 \right] + 2\mathbb{E} \left[(\hat{P}D_i - \mathbb{E}[d_i | \hat{P}D_i])(\mathbb{E}[d_i | \hat{P}D_i] - d_i) \right] + \mathbb{E} \left[(\mathbb{E}[d_i | \hat{P}D_i] - d_i)^2 \right] \\ &= \mathbb{E} \left[(\hat{P}D_i - \mathbb{E}[d_i | \hat{P}D_i])^2 \right] + \mathbb{E} \left[(\mathbb{E}[d_i | \hat{P}D_i] - d_i)^2 \right] \\ &= \mathbb{E} \left[(\hat{P}D_i - \mathbb{E}[d_i | \hat{P}D_i])^2 \right] + \mathbb{E} \left[\text{Var}(d_i | \hat{P}D_i) \right] \\ &= \underbrace{\mathbb{E} \left[(\hat{P}D_i - \mathbb{E}[d_i | \hat{P}D_i])^2 \right]}_{\text{Reliability}} - \underbrace{\text{Var}(\mathbb{E}[d_i | \hat{P}D_i])}_{\text{Resolution}} + \underbrace{\text{Var}(d_i)}_{\text{Uncertainty}}, \end{aligned}$$

where the second term in the fourth line equals zero because $\mathbb{E} \left[(\hat{P}D_i - \mathbb{E}[d_i | \hat{P}D_i]) \right]$ does not vary with d_i and the term it is multiplied by, $\mathbb{E}[d_i | \hat{P}D_i] - d_i$, is mean zero due to the law of total expectation and the last line comes from the law of total variance.

The first term is referred to as reliability and measures how close the bank's average PD is to the average realized default rate. In the model above, this will always be zero because banks' perceived PDs are, on average, unbiased, i.e., $\hat{P}D = \mathbb{E}[d | \hat{P}D]$. However, in practice, this may not always be the case. For example, banks may have incorrect priors, biased beliefs, or incentives to report PDs that do not match their beliefs. The second term is referred to as resolution and measures the discriminatory ability of the model. Finally, the last term is equal to the variance of default rates, which is maximized when the probability of default is closest to 50%. Hence, the Brier Score captures multiple aspects of the forecast at once and is not a pure measure of discrimination.

Even beyond the abovementioned issues, the Brier score is less commonly used by practitioners to evaluate PD models for several reasons. First, there is a strong mechanical relationship between average forecast errors and the underlying risk of the borrower, particularly when

PDs are close to zero. To see this, assume that the bank observes the PD of a loan perfectly. We can then write the expected squared forecast error as follows:

$$\hat{PD}(1 - \hat{PD})^2 + (1 - \hat{PD})\hat{PD}^2.$$

Differentiating this w.r.t to \hat{PD} we have $1 - 2\hat{PD}$. Hence, so long as the PD is less than 50%, a higher PD leads to a higher average squared forecast error. Notice also that the second derivative is equal to -2 ; hence, the effect is stronger the closer the PD is to zero. For example, suppose the bank perfectly observes the borrower's PD. Increasing the PD from 1% to 2% roughly doubles the expected squared forecast error (0.99pp to 1.96pp) even though the bank's information is not any better or worse.⁴⁹

Second, even in the absence of the issue above, the Brier score does not distinguish forecast ability well for rare events. For example, suppose that if a bank does not produce information, it assigns a PD of 1% to a firm, whereas if the bank does produce information, it will realize the firm's PD is 0.9%. The improvement in the Brier score from producing information would be 10^{-6} even though the bank's PD is substantially more accurate.⁵⁰

These exercises suggest that the Brier score is a fairly unreliable measure of the discriminatory power of PDs and illustrate why the AUC approach is more commonly used by practitioners and regulators (e.g., [Engelmann and Rauhmeier \(2011\)](#) and [Basel Committee on Banking Supervision \(2005\)](#)).

⁴⁹Note the expected absolute forecast error also almost doubles (1.98pp to 3.92pp).

⁵⁰See [Benedetti \(2010\)](#) for a similar discussion.

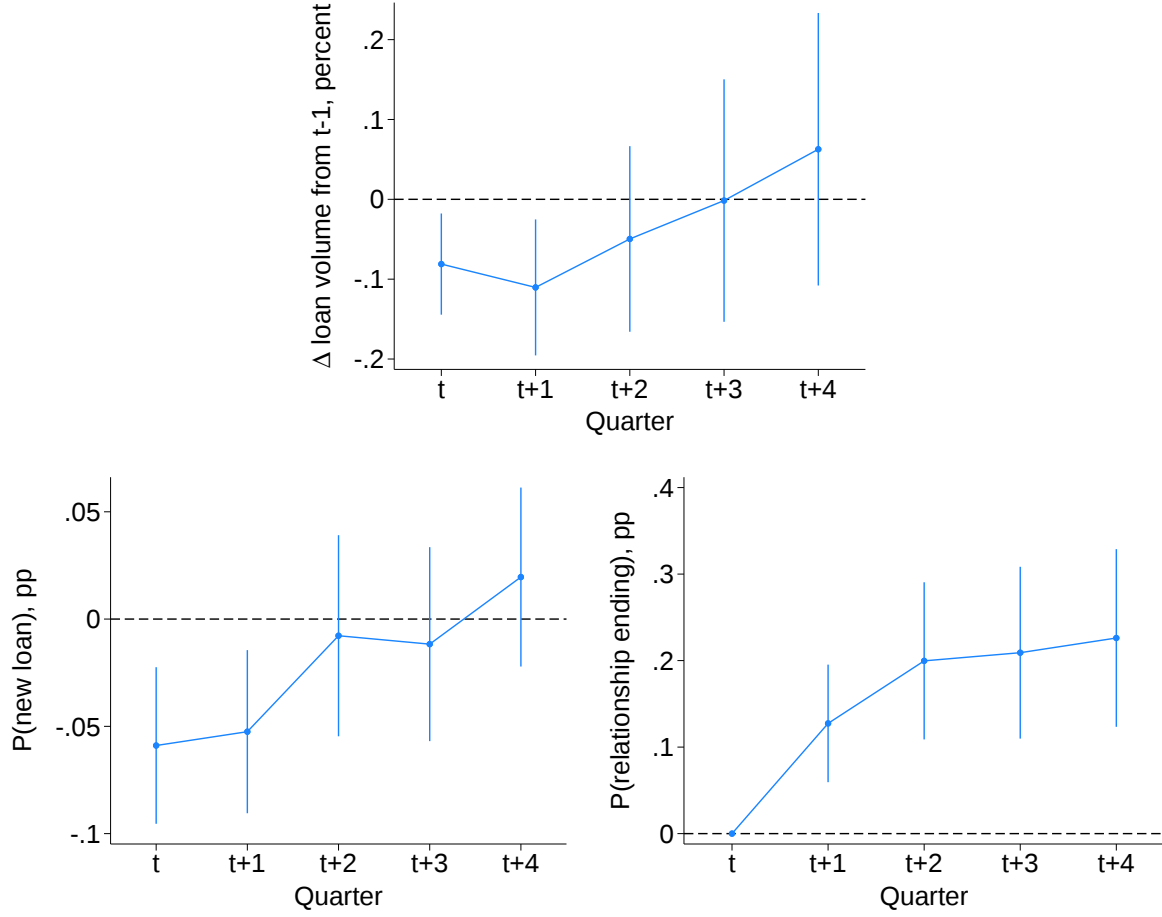


Figure OA.1: Dynamic responses to changes in PD

This figure shows dynamic responses of a change in bank b 's assessed PD for firm i at time t on outcome $y_{i,b,t+j}$ for $j \in \{0, 1, 2, 3, 4\}$ after collapsing all observations of new and existing loans into a firm-bank-quarter panel. Each plot shows coefficient estimates β^j from the following regression specification: $y_{i,b,t+j} = \beta^j PD_{i,b,t} + \alpha_{i,t}^j + \gamma_{i,b}^j + \epsilon_{i,b,t+j}$, where $PD_{i,b,t}$ is bank b 's PD for firm i at time t , $\alpha_{i,t}$ is a firm-time fixed effect, and $\gamma_{i,b}$ is a firm-bank fixed effect. 95% confidence intervals are shown and calculated using standard errors that are two-way clustered by firm and bank-time. The top panel shows the change in log committed lending volume from time $t - 1$. The bottom-left panel shows the probability of a firm having at least one new loan. The bottom-right panel shows the probability that the firm-bank match is not observed in our data.

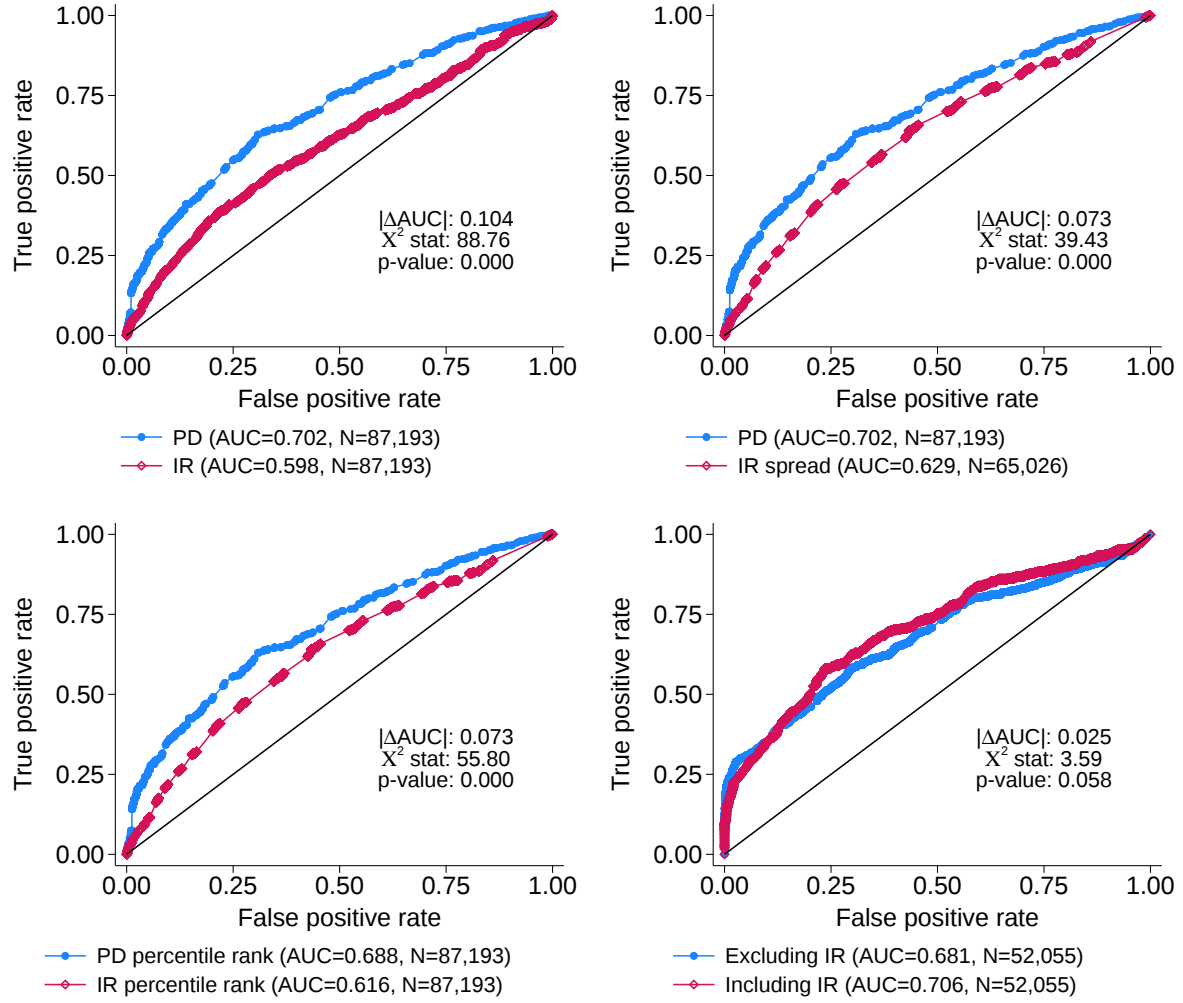


Figure OA.2: ROC comparing PD and interest rates

The top-left panel compares ROC curves calculated using PD and interest rates. The top-right panel uses PD and interest rate spreads. The bottom-left panel uses percentile ranks for PD and the interest rate rank across all observations within each quarter. The bottom-right panel shows the marginal predictive power from including the interest rate in a random forest regression estimate of the probability of default. The blue curve plots the ROC curve for default estimates using the following controls: log loan size, LGD, log maturity, firm leverage, firm profitability, firm tangibility, and firm size. The red curve plots estimates using the same controls plus interest rate. The area under each ROC curve (AUC) is reported along with the number of observations in the legend. All figures use our baseline sample of new loans. The area under each ROC curve (AUC) is reported along with the number of observations in the legend. $|\Delta AUC|$ reports the difference between the two AUCs. Below $|\Delta AUC|$, the DeLong, DeLong, and Clarke-Pearson (1988) statistics are reported: the χ^2 test statistic and its corresponding p-value, which tests the null hypothesis that the difference between the two AUCs equals zero.

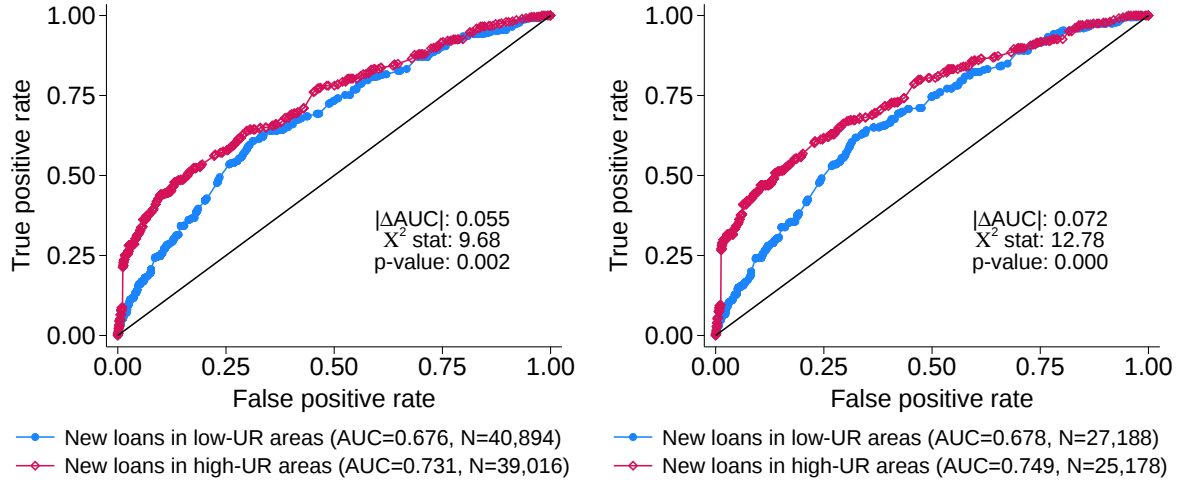


Figure OA.3: ROC excluding counties with few loans

This figure shows ROC curves for new loans split by the local unemployment rate at origination. The left panel restricts the sample to county-quarters that have at least five total new loans from 2014Q4-2021Q4, which drops about 25% of all loan-quarter observations. The right panel restricts the sample to county-quarters that have at least 15 new loans from 2014Q4-2021Q4, which drops roughly half of all loan-quarter observations. The area under each ROC curve (AUC) is reported along with the number of observations in the legend. $|\Delta AUC|$ reports the difference between the two AUCs. Below $|\Delta AUC|$, the DeLong, DeLong, and Clarke-Pearson (1988) statistics are reported: the χ^2 test statistic and its corresponding p-value, which tests the null hypothesis that the difference between the two AUCs equals zero.

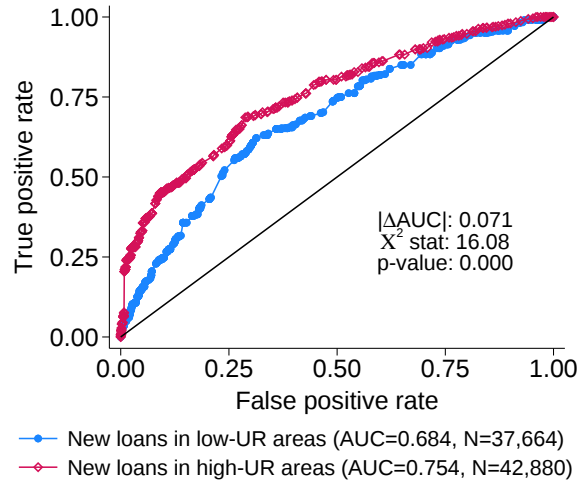
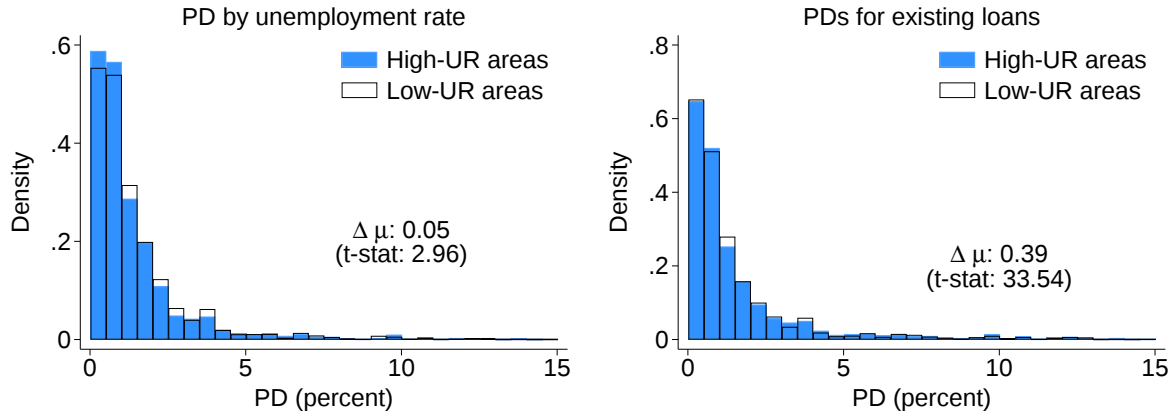


Figure OA.4: ROC for pre-2020 loans

This figure shows ROC curves for new loans in our sample originated from 2014Q4-2019Q4 split by whether the unemployment rate at origination was above or below its county-level median during this period. The area under each ROC curve (AUC) is reported along with the number of observations in the legend. $|\Delta AUC|$ reports the difference between the two AUCs. Below $|\Delta AUC|$, the DeLong, DeLong, and Clarke-Pearson (1988) statistics are reported: the χ^2 test statistic and its corresponding p-value, which tests the null hypothesis that the difference between the two AUCs equals zero.

	New loans		Existing loans	
	Low UR	High UR	Low UR	High UR
5th percentile	0.14	0.14	0.12	0.12
25th percentile	0.45	0.44	0.37	0.40
Median	0.93	0.90	0.88	0.90
Mean	1.61	1.66	2.33	2.72
75th percentile	1.91	1.87	1.91	2.00
95th percentile	4.66	5.37	9.06	11.00
Standard deviation	2.58	2.88	6.54	7.71

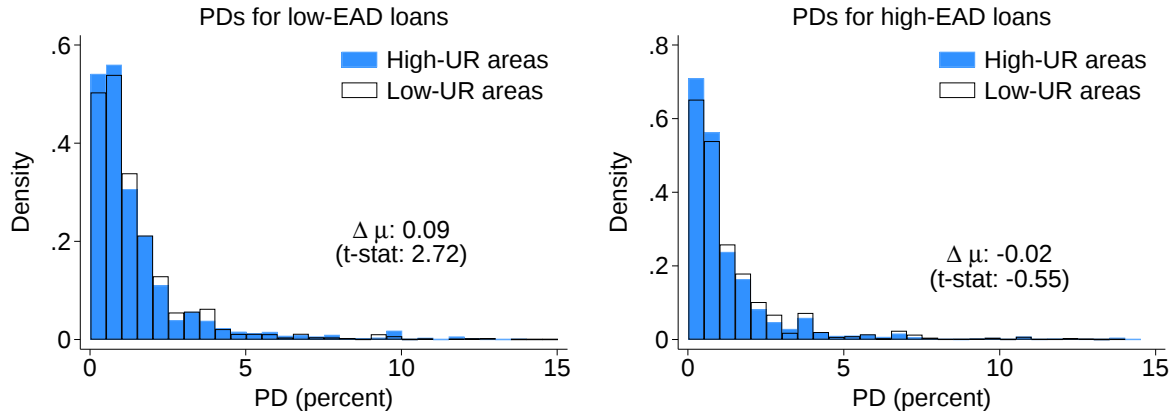


	New loans			Existing loans		
	Low UR	High UR	Difference	Low UR	High UR	Difference
AUC (Data)	0.681	0.724	0.043	0.806	0.830	0.024
AUC (Uniform)	0.542	0.551	0.009	0.688	0.725	0.037
AUC (Lognormal)	0.792	0.801	0.009	0.854	0.855	0.001

Figure OA.5: PD distributions by UR for new (left) and existing (right) loans

The top table reports summary statistics for PD split by periods of high versus low UR for new loans (left columns) or existing loans (right columns). The bottom figures show the overlaid distributions of PDs for periods of high and low UR for new loans (bottom-left) and existing loans (bottom-right). Both distributions are truncated at 15 percent for readability. Sample sizes for all figures and summary statistics are reported in Table 1 (for new loans) and Online Appendix Table OA.1 (for existing loans). The bottom panel compares the AUCs we estimate in the data in Figures 4 and 6 with the model-implied AUCs assuming there is no information production. We estimate the AUCs analytically assuming PDs follow a uniform distribution, and based on simulations assuming PDs follow a lognormal distribution. In both cases, we assume the distribution of PD has the same mean and variance as the specific subsample of the data displayed in the top panel. For the lognormal case, we estimate each AUC 100 times based on simulations of 25,000 PDs and report the sample average AUC.

	Low-EAD loans		High-EAD loans	
	Low UR	High UR	Low UR	High UR
5th percentile	0.16	0.16	0.11	0.10
25th percentile	0.51	0.49	0.36	0.35
Median	1.00	0.93	0.88	0.75
Mean	1.62	1.71	1.56	1.54
75th percentile	1.91	1.90	1.81	1.53
95th percentile	4.54	5.67	5.34	5.34
Standard deviation	2.43	2.76	2.59	2.80

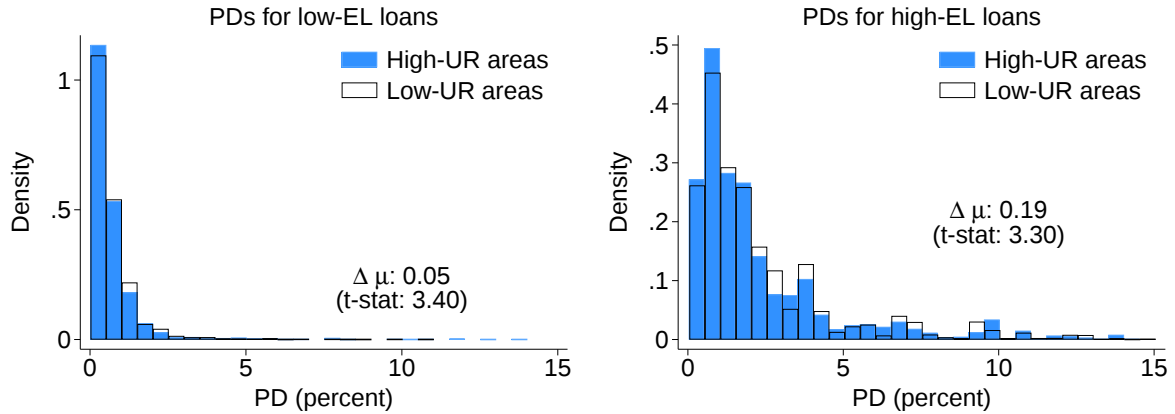


	Low-EAD loans			High-EAD loans		
	Low UR	High UR	Difference	Low UR	High UR	Difference
AUC (Data)	0.686	0.646	-0.040	0.676	0.755	0.079
AUC (Uniform)	0.537	0.545	0.008	0.544	0.552	0.008
AUC (Lognormal)	0.783	0.792	0.009	0.796	0.806	0.010

Figure OA.6: PD distributions by UR for low-EAD (left) and high-EAD (right) loans

The top table reports summary statistics for PD split by periods of high versus low UR for loans with below-median (left columns) or above-median (right columns) exposure at default (EAD) calculated within each bank-quarter. The middle panel displays the overlaid distributions of PDs for periods of high and low UR for bottom quartile (bottom-left) or top quartile (bottom-right) EAD. Both distributions are truncated at 15 percent for readability. Sample sizes for all figures and summary statistics are the same as shown in the bottom panels of Figure 7. The bottom panel compares the AUCs we estimate in the data in Figure 7 with the model-implied AUCs assuming there is no information production. We estimate the AUCs analytically assuming PDs follow a uniform distribution, and based on simulations assuming PDs follow a lognormal distribution. In both cases, we assume the distribution of PD has the same mean and variance that we observe in the specific subsample of the data displayed in the top panel. For the lognormal case, we estimate each AUC 100 times based on simulations of 25,000 PDs and report the sample average AUC.

	Low-EL loans		High-EL loans	
	Low UR	High UR	Low UR	High UR
5th percentile	0.09	0.08	0.28	0.24
25th percentile	0.22	0.20	0.87	0.75
Median	0.44	0.43	1.68	1.50
Mean	0.67	0.72	2.95	3.14
75th percentile	0.87	0.81	3.07	3.33
95th percentile	2.00	2.00	10.00	12.40
Standard deviation	0.95	1.28	4.27	4.79



	Low-EL loans			High-EL loans		
	Low UR	High UR	Difference	Low UR	High UR	Difference
AUC (Data)	0.552	0.592	0.040	0.655	0.761	0.106
AUC (Uniform)	0.514	0.523	0.009	0.564	0.575	0.012
AUC (Lognormal)	0.771	0.800	0.029	0.781	0.788	0.007

Figure OA.7: PD distributions by UR for low-EL (left) and high-EL (right) loans

The top table reports summary statistics for PD split by periods of high versus low UR for loans with bottom quartile (left columns) or top quartile (right columns) expected loss (EL) calculated within each bank-quarter. The bottom figures show the overlaid distributions of PDs for periods of high and low UR for low-EL loans (bottom-left) and high-EL loans (bottom-right). Both distributions are truncated at 15 percent for readability. Sample sizes for all figures and summary statistics are the same as shown in the bottom panels of Figure 8. The bottom panel compares the AUCs we estimate in the data in Figure 8 with the model-implied AUCs assuming there is no information production. We estimate the AUCs analytically assuming PDs follow a uniform distribution, and based on simulations assuming PDs follow a lognormal distribution. In both cases, we assume the distribution of PD has the same mean and variance that we observe in the specific subsample of the data displayed in the top panel. For the lognormal case, we estimate each AUC 100 times based on simulations of 25,000 PDs and report the sample average AUC.

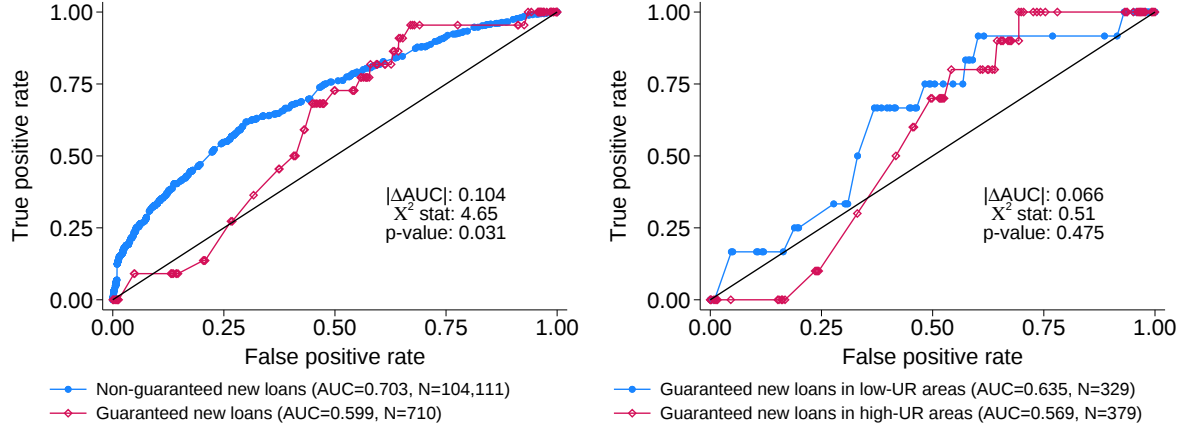


Figure OA.8: ROC for new loans with government guarantees

This figure shows ROC curves for new loans split by the local unemployment rate at origination. The left panel compares loans with and without government guarantees. The right panel compares loans with government guarantees depending on whether the unemployment rate is above or below its county-level median. The area under each ROC curve (AUC) is reported along with the number of observations in the legend. $|\Delta AUC|$ reports the difference between the two AUCs. Below $|\Delta AUC|$, the [DeLong, DeLong, and Clarke-Pearson \(1988\)](#) statistics are reported: the χ^2 test statistic and its corresponding p-value, which tests the null hypothesis that the difference between the two AUCs equals zero.

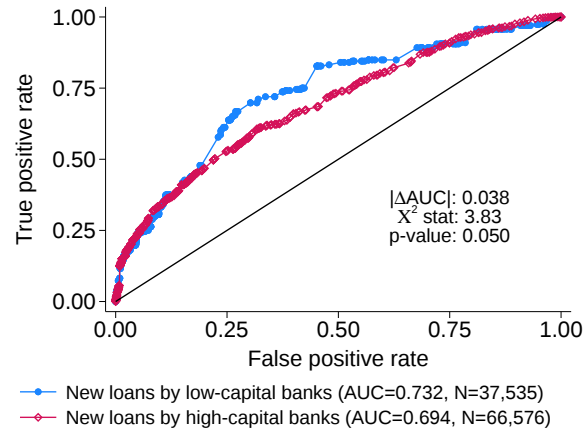


Figure OA.9: ROC by bank capitalization

This figure shows ROC curves for new loans split by whether the bank's total risk-based capital ratio was above or below the median of all banks in our data in each quarter. The area under each ROC curve (AUC) is reported along with the number of observations in the legend. $|\Delta AUC|$ reports the difference between the two AUCs. Below $|\Delta AUC|$, the [DeLong, DeLong, and Clarke-Pearson \(1988\)](#) statistics are reported: the χ^2 test statistic and its corresponding p-value, which tests the null hypothesis that the difference between the two AUCs equals zero.

Table OA.1: Summary statistics for both new and existing loans

This table contains summary statistics for both new and existing loans. Because a small number of loans are in counties with missing unemployment rates, the sample sizes in panels B and C do not add up to panel A. Section 2 describes our sample.

	Mean	Median	5%	95%	SD	N
Panel A: All loans						
Interest rate (pp)	3.63	3.50	1.56	6.00	1.67	1,230,296
PD (pp)	2.47	0.90	0.12	9.82	6.97	1,670,940
LGD (ratio)	0.34	0.35	0.08	0.60	0.16	1,655,331
Realized default (pp)	1.64	0.00	0.00	0.00	12.68	1,670,940
Maturity (months)	72.52	60.00	12.00	175.00	114.98	1,670,940
Loan size (\$ mil)	13.30	4.00	1.00	55.00	28.32	1,670,940
Revolver (indicator)	0.58	1.00	0.00	1.00	0.49	1,670,940
Term loan (indicator)	0.27	0.00	0.00	1.00	0.44	1,670,940
Floating rate (indicator)	0.60	1.00	0.00	1.00	0.49	1,670,940
Panel B: All loans in high-UR areas						
Interest rate (pp)	3.21	3.03	1.40	5.59	1.54	630,404
PD (pp)	2.66	0.90	0.12	10.99	7.52	850,786
LGD (ratio)	0.34	0.35	0.08	0.57	0.16	842,864
Realized default (pp)	1.62	0.00	0.00	0.00	12.61	850,786
Maturity (months)	71.92	60.00	12.00	171.00	113.94	850,786
Loan size (\$ mil)	13.28	4.00	1.00	54.45	29.54	850,786
Revolver (indicator)	0.58	1.00	0.00	1.00	0.49	850,786
Term loan (indicator)	0.26	0.00	0.00	1.00	0.44	850,786
Floating rate (indicator)	0.60	1.00	0.00	1.00	0.49	850,786
Panel C: All loans in low-UR areas						
Interest rate (pp)	4.07	3.99	1.98	6.31	1.70	595,469
PD (pp)	2.29	0.89	0.12	8.43	6.36	812,183
LGD (ratio)	0.34	0.35	0.09	0.60	0.16	804,516
Realized default (pp)	1.67	0.00	0.00	0.00	12.80	812,183
Maturity (months)	73.10	60.00	12.00	178.00	116.52	812,183
Loan size (\$ mil)	13.14	4.00	1.00	55.00	26.52	812,183
Revolver (indicator)	0.57	1.00	0.00	1.00	0.49	812,183
Term loan (indicator)	0.28	0.00	0.00	1.00	0.45	812,183
Floating rate (indicator)	0.60	1.00	0.00	1.00	0.49	812,183

Table OA.2: Correlations between PD and controls

This table shows the correlation between PD, interest rate (IR), leverage, profitability, tangibility, firm size, loan size, LGD, and maturity for our baseline sample of new loans shown in panel A of Table 1.

Variables	PD	IR	Lev.	Prof.	Tang.	Firm size	Loan size	LGD	Maturity
PD	1.000								
IR	0.224	1.000							
Lev.	0.103	0.063	1.000						
Prof.	-0.143	-0.025	-0.015	1.000					
Tang.	-0.043	-0.137	-0.182	0.084	1.000				
Firm size	-0.021	-0.122	0.034	-0.285	-0.220	1.000			
Loan size	0.010	-0.057	0.062	-0.135	-0.206	0.499	1.000		
LGD	-0.056	0.047	-0.034	0.114	-0.088	-0.058	-0.075	1.000	
Maturity	-0.031	0.075	0.096	0.037	-0.150	0.046	0.131	-0.020	1.000

Table OA.3: Variance inflation factors

This table reports variance inflation factors for PD and our set of firm and loan characteristics. Using an indicator that equals one if the loan defaults within two years of origination as the dependent variable, we first estimate a regression using all independent variables shown in the first column (the “full” model). Next, we estimate separate regressions including only one of the explanatory variables at a time (the “separate” models). For each independent variable, the VIF is calculated as the ratio of the variance of the parameter estimate in the full model to the variance of the parameter estimate in the separate model. The right column shows the tolerance, which is simply defined as $1/\text{VIF}$.

Variables	VIF	Tolerance
Firm size	1.47	0.679
Loan size	1.38	0.724
Tangibility	1.15	0.867
Profitability	1.13	0.885
Interest rate	1.11	0.904
PD	1.11	0.905
Maturity	1.07	0.932
Leverage	1.05	0.948
LGD	1.04	0.960

Table OA.4: Alternate loan performance measures

This table shows the results of estimating (1) with alternative measures of loan performance. The dependent variable in each regression is a dummy variable corresponding to the column heading multiplied by 100. “Any default” measures whether a loan is recorded as defaulting at any point in our sample period. “Average default” divides the “Any default” measure by the number of years in which the loan is observed to generate an annual average; if a loan defaults within one quarter after origination, this variable will take on a value of 2, while if the loan defaults eight quarters after origination, this variable will take on a value of 0.5. “1Y default” is an indicator that equals one if the loan defaults within four quarters of origination. “Delinquency” is an indicator that equals one if the loan is reported as delinquent within eight quarters after origination. “Chargeoff” is an indicator representing whether a bank records a chargeoff for that loan within eight quarters after origination. Probability of default (PD) is measured in percentage points. Standard errors are clustered at the county level and shown in parentheses.

	Any default (1)	Average default (2)	1Y default (3)	Delinquency (4)	Chargeoff (5)
PD	0.519*** (0.066)	0.541*** (0.082)	0.331*** (0.055)	0.107*** (0.020)	0.065*** (0.019)
Controls	Y	Y	Y	Y	Y
Bank-quarter FE	Y	Y	Y	Y	Y
Industry-quarter FE	Y	Y	Y	Y	Y
Bank-county FE	Y	Y	Y	Y	Y
Observations	77,226	77,226	77,226	77,226	77,226
R ²	0.204	0.212	0.181	0.125	0.171

Table OA.5: PDs predict interest rates

This table shows that banks' PDs predict interest rates even after controlling for observable characteristics. The dependent variable in these regressions is the interest rate measured in percentage points. All regressions include our standard firm and loan controls. Standard errors are clustered at the county level and are shown below the parameter estimates in parentheses. *, **, and *** indicate statistical significance at the 10%, 5%, and 1% levels, respectively.

	Interest rate		
	(1)	(2)	(3)
PD		0.092*** (0.003)	0.088*** (0.004)
LGD		0.695*** (0.060)	0.647*** (0.057)
Leverage	0.354*** (0.041)		0.250*** (0.038)
Profitability	-0.294*** (0.039)		-0.154*** (0.038)
Tangibility	-0.936*** (0.054)		-0.859*** (0.055)
Log firm size	-0.145*** (0.008)		-0.131*** (0.008)
Log loan amount	-0.036*** (0.007)		-0.038*** (0.007)
Log maturity	0.023 (0.017)		0.046*** (0.017)
Bank-quarter FE	Y	Y	Y
Industry-quarter FE	Y	Y	Y
Bank-county FE	Y	Y	Y
Observations	66,121	81,752	64,566
R ²	0.552	0.518	0.578

Table OA.6: PDs predict interest rates (random forest)

This table shows that banks' PDs predict interest rates even after controlling for a predicted PD estimated from observable characteristics using a random forest. The dependent variable in each regression is each loan's interest rate reported at origination in percentage points. "RF predicted interest rate" is the random forest estimate of the loan's interest rate. Each column corresponds to a different set of controls used to estimate the random forest. All specifications include our default set of firm and loan controls; other specifications include indicator variables for industry, bank, and time and are indicated in the rows below the table. "PD" is the probability of default reported by the bank and is measured in percentage points. All regressions exclude the half of the baseline sample used to train the random forest. Standard errors are clustered at the county level and are shown below the parameter estimates in parentheses. *, **, and *** indicate statistical significance at the 10%, 5%, and 1% levels, respectively.

	Interest rate				
	(1)	(2)	(3)	(4)	(5)
RF predicted interest rate	1.426*** (0.032)	1.343*** (0.024)	1.240*** (0.021)	1.140*** (0.010)	1.183*** (0.011)
PD	0.096*** (0.004)	0.085*** (0.004)	0.083*** (0.004)	0.102*** (0.004)	0.081*** (0.004)
Industry controls	N	Y	N	N	Y
Bank controls	N	N	Y	N	Y
Time controls	N	N	N	Y	Y
Observations	43,383	43,383	43,383	43,383	43,383
R ²	0.250	0.282	0.324	0.404	0.491

Table OA.7: Probability of contemporaneous loan modifications accompanying changes in PD

The dependent variable in these regressions is the probability that a firm's maturity date or interest rate spread changed from the previous quarter, with "Any change" corresponding to a change in at least one of the two. "Change in PD" is an indicator variable equal to one if a loan's PD changed from the prior quarter. Standard errors are clustered at the county level and are shown below the parameter estimates in parentheses. *, **, and *** indicate statistical significance at the 10%, 5%, and 1% levels, respectively.

	Change in maturity				Change in IR				Any change			
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)
Change in PD	0.093*** (0.002)	0.092*** (0.002)	0.090*** (0.002)	0.012*** (0.002)	0.048*** (0.001)	0.045*** (0.001)	0.039*** (0.001)	0.049*** (0.002)	0.089*** (0.002)	0.098*** (0.002)	0.090*** (0.002)	0.011*** (0.003)
Leverage		-0.060*** (0.004)	-0.018*** (0.003)	0.004 (0.008)		0.057*** (0.003)	0.035*** (0.004)	-0.001 (0.011)		-0.009** (0.005)	0.006 (0.004)	-0.002 (0.013)
Profitability		-0.039*** (0.003)	-0.022*** (0.004)	0.002 (0.014)		0.001 (0.003)	-0.013*** (0.003)	-0.004 (0.015)		-0.032*** (0.004)	-0.030*** (0.005)	-0.003 (0.018)
Tangibility		0.071*** (0.003)	0.008 (0.005)	0.011 (0.014)		-0.040*** (0.004)	-0.032*** (0.006)	-0.024 (0.017)		0.015*** (0.005)	-0.007 (0.008)	-0.022 (0.024)
Log firm size		-0.015*** (0.001)	-0.000 (0.001)	0.002 (0.002)		0.005*** (0.001)	0.001 (0.001)	0.001 (0.002)		-0.008*** (0.001)	-0.000 (0.001)	0.002 (0.003)
Log loan amount		-0.014*** (0.001)	0.007*** (0.001)	0.010*** (0.001)		0.037*** (0.001)	0.032*** (0.001)	0.033*** (0.002)		0.023*** (0.001)	0.037*** (0.002)	0.039*** (0.002)
LGD		-0.019*** (0.005)	0.028*** (0.006)	0.072*** (0.013)		-0.040*** (0.006)	-0.027*** (0.005)	-0.049*** (0.012)		-0.053*** (0.007)	0.002 (0.008)	0.029 (0.018)
Log maturity		0.022*** (0.002)	0.028*** (0.003)	-0.023*** (0.003)		0.011*** (0.001)	0.017*** (0.001)	0.023*** (0.002)		0.027*** (0.002)	0.032*** (0.002)	-0.007** (0.003)
Controls	N	Y	Y	Y	N	Y	Y	Y	N	Y	Y	Y
Firm-bank FE	N	N	Y	Y	N	N	Y	Y	N	N	Y	Y
Firm-quarter FE	N	N	N	Y	N	N	N	Y	N	N	N	Y
Bank-quarter FE	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
Industry-quarter FE	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
Bank-county FE	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
Observations	1,413,034	1,228,694	1,214,205	625,903	1,647,212	1,412,737	1,396,866	741,691	1,392,922	1,228,694	1,214,205	625,903
R ²	0.158	0.177	0.314	0.705	0.253	0.277	0.455	0.691	0.370	0.209	0.370	0.696

Table OA.8: Δ AUCs estimated by bank

The top tables report summary statistics for Δ AUC when calculated separately for each bank for our sample of all loans (sample sizes are reported in Online Appendix Table OA.1) split by current UR (left panel) or origination UR (right panel). The “unweighted” column calculates summary statistics using equal weights across banks, while the “weighted” column weights the summary statistics by the number of loan-quarter observations for each bank. These calculations are based on the 22 banks with at least 10,000 loan-quarter observations in our sample (out of 33 total banks). The bottom table reports information about the banks used to calculate the summary statistics in the top tables. The left column calculates AUCs based on the current UR, while the right column uses the origination UR. The second row reports the total number of banks estimated to have Δ AUC > 0 , which means that their information quality is higher when the UR is high. The third row reports the number of banks for which Δ AUC > 0 and this difference is statistically significant. The fourth and fifth rows report the share of total loans or share of total loan volume, respectively, coming from the banks with Δ AUC > 0 and $p < 0.05$ (shown in the third row).

	Current UR			Origination UR	
	Unweighted	Weighted		Unweighted	Weighted
25th percentile	0.009	0.013	25th percentile	0.002	0.002
Median	0.033	0.026	Median	0.029	0.030
Mean	0.027	0.024	Mean	0.043	0.045
75th percentile	0.047	0.047	75th percentile	0.064	0.043
SD	0.050	0.049	SD	0.105	0.101

	Current UR	Origination UR
Total number banks	22	22
Banks with Δ AUC > 0	17	17
Banks with Δ AUC > 0 and $p < 0.05$	12	9
Volume share Δ AUC > 0 and $p < 0.05$	66.84%	46.25%
Loan share Δ AUC > 0 and $p < 0.05$	51.45%	44.91%

Table OA.9: Default sensitivity to the business cycle

This table compares the sensitivity of default to the unemployment rate for new versus old loans and high versus low EAD loans. The dependent variable in each regression is an indicator that equals one if the loan defaults within the next 8 quarters. “High UR” is an indicator that equals one if the unemployment rate is above its county-level median during our sample. High EAD/LGD/EL loan is an indicator that equals one if the loan’s exposure at default (EAD), loss given default (LGD), or expected loss ($EAD \times LGD \times PD$) is above the median within a given bank-quarter. Standard errors are clustered at the county level and are shown below the parameter estimates in parentheses. *, **, and *** indicate statistical significance at the 10%, 5%, and 1% levels, respectively.

	Default							
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
High UR	-0.105 (0.074)	-0.113 (0.072)	-0.005 (0.083)	-0.017 (0.081)	-0.102 (0.086)	-0.117 (0.083)	-0.086 (0.074)	-0.052 (0.068)
New loan	-0.704*** (0.099)	-1.212 (0.832)						
New loan \times High UR	0.109 (0.090)	0.125 (0.108)						
High EAD loan			0.340*** (0.062)	1.413* (0.857)				
High EAD loan \times High UR			-0.192** (0.080)	-0.183** (0.085)				
High LGD loan					0.301*** (0.070)	1.279 (0.844)		
High LGD loan \times High UR					0.007 (0.077)	0.026 (0.080)		
High EL loan							1.770*** (0.120)	6.175*** (1.449)
High EL loan \times High UR							-0.050 (0.088)	-0.132 (0.094)
Controls	N	Y	N	Y	N	Y	N	Y
Bank-quarter FE	Y	Y	Y	Y	Y	Y	Y	Y
Industry-quarter FE	Y	Y	Y	Y	Y	Y	Y	Y
Bank-county FE	Y	Y	Y	Y	Y	Y	Y	Y
Observations	1,647,139	1,425,014	1,647,139	1,425,014	1,647,139	1,425,014	1,647,139	1,425,014
R ²	0.132	0.141	0.132	0.141	0.132	0.141	0.135	0.145