

TINBERGEN INSTITUTE

Machine Learning | Block Two

Risk at First Sight

*Default-Risk Prediction & Screening at Loan Origination in P2P
Consumer Lending, with a Double Machine Learning Extension of the
Effects of Longer Terms and High Interest Rates*

Student Number: **756486**

December 15, 2025

Disclosure

In the preparation of this research project, AI has been used for the following purposes:

1. **Idea development and discussion:** to bounce off ideas, explore alternative formulations, and stress-test the logic behind the research design.
2. **Validity checks (reasoning & implementation):** to help check the consistency of mathematical steps, identify possible mistakes, and diagnose issues in code implementation (e.g., debugging suggestions and sanity checks).
3. **Clarity and presentation:** to improve English phrasing, structure, and readability, and to refine or simplify L^AT_EX and mathematical notation when needed.
4. **Code quality and maintainability:** to improve code structure, readability, and organization (e.g., refactoring, clearer naming, comments, and modularization) without changing the intended analysis.

All written content, methodological choices, analyses, and interpretations reflect my own reasoning and final decisions. AI outputs were treated as suggestions, reviewed critically, and edited or discarded when necessary. Any remaining errors are my responsibility. AI was not used to fabricate data, results, or references; all sources and claims are verified by the author.

1 Introduction

Online peer-to-peer (P2P) lending has emerged as a significant alternative to traditional banking, directly connecting borrowers and lenders and shifting credit risk from intermediaries to investors (Lenz, 2016). Because these loans are typically unsecured and investors often lack the expertise to evaluate creditworthiness, P2P markets face severe information asymmetry and a heightened risk of adverse selection (Byanjankar et al., 2015; Emekter et al., 2015; Jin & Zhu, 2015). In practice, default-risk assessment becomes a *screening* problem under limited attention: investors and platforms must decide which applications to flag first, not merely estimate probabilities. Moreover, evidence suggests that charging higher interest rates to riskier borrowers may be insufficient to offset expected losses, reinforcing the need for effective screening at origination (Emekter et al., 2015).

In this paper, we ask: *How well can we predict default risk at loan origination using borrower and loan characteristics, and what does that imply for screening decisions?* We study LendingClub loans and restrict the feature set strictly to origination-time variables to prevent data leakage. We benchmark standard statistical approaches against commonly used machine learning methods that may better capture non-linear patterns in borrower characteristics (Byanjankar et al., 2015; Jin & Zhu, 2015). Unlike much of the existing P2P credit-scoring literature, which focuses mainly on generic predictive metrics, we explicitly link model performance to screening outcomes and to the distribution of access to credit. Our evaluation therefore emphasizes not only predictive performance but also ranking quality, which determines how effectively a model supports screening decisions.

As a complementary extension, we examine how default risk differs across two salient contract features: loan maturity (60 vs. 36 months) and interest rates. Because these terms may be endogenous to borrower risk, we first report naive differences in default rates as descriptive benchmarks and then use double machine learning (DML) to estimate *conditional* treatment effects while flexibly controlling for origination-time covariates, interpreting these estimates as adjusted associations and remaining cautious about causal claims (Emekter et al., 2015; Serrano-Cinca et al., 2015).

The remainder of the paper is organized as follows: Section 2 describes the data, Section 3 outlines the modeling approach, Section 4 presents the results, and Section 5 concludes.

2 Data

To study default-risk prediction and screening *at loan origination*, we use a widely used LendingClub loan-level dataset distributed via Kaggle. LendingClub is a major U.S. online marketplace for consumer credit in which investors fund unsecured installment loans; it was incorporated in 2006 and launched its marketplace in 2007 (LendingClub Corporation, 2014). The raw file contains 2,260,701 loan records with 151 variables spanning issuance years 2007–2018, covering borrower attributes, contract terms, platform risk measures, and loan outcomes.

Because our question is explicitly about *screening at origination*, we restrict features to variables observable at issuance and construct a clean default outcome from the loan-status field. Table 4 summarizes the processing steps (create issuance year, harmonize status labels, drop invalid records). To avoid right-censoring, we focus on loans issued in 2012–2015; as Figure 5 shows, later cohorts contain many active loans, which would mechanically bias default rates downward and contaminate screening evaluation.

Even after this restriction, the data remain high-dimensional (over 100 variables), and one-hot expansion of categorical features would exceed local computational limits. We therefore use a stratified sample (by issuance year and default outcome) and literature-guided feature selection, remove identifiers and post-origination leakage fields (e.g., realized payments), standardize factor levels, and engineer a parsimonious origination-time predictor set. Our binary outcome equals one

for loans that end in default and zero otherwise under the harmonized status definition. Predictors are origination-time borrower and contract characteristics used in screening, including credit-quality proxies (e.g., credit grade and FICO band), pricing (interest rate), maturity (36 vs. 60 months), and affordability measures (e.g., income and debt-to-income).

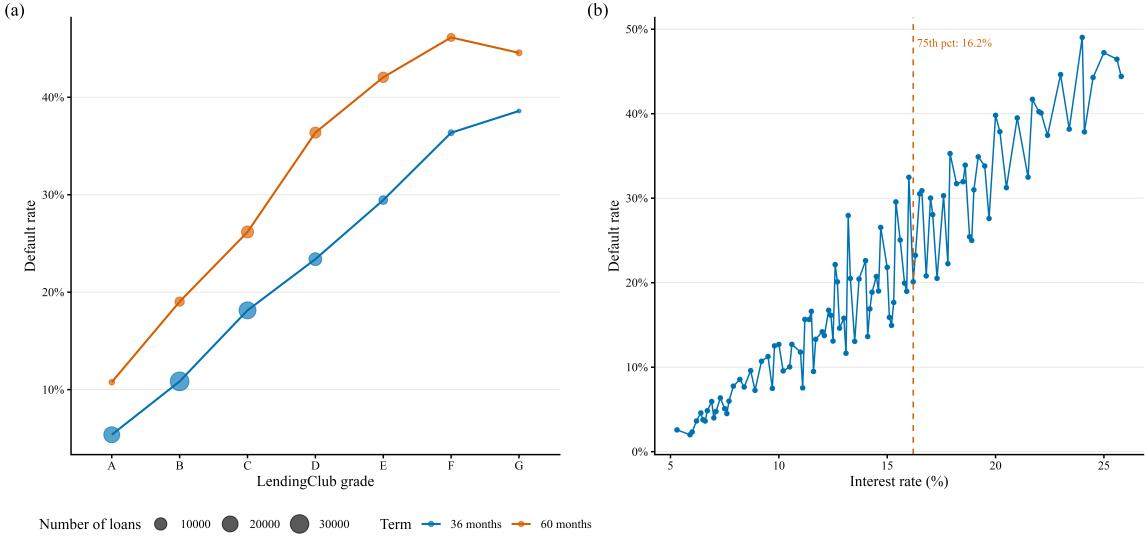


Figure 1: Default rates by LendingClub grade and loan term, and by interest rate.

The cohort exhibits strong monotone risk gradients in the platform’s internal risk signals and in observed pricing. Figure 1(a) shows default rates rising sharply from grade A to grade G, with 60-month loans defaulting more often than 36-month loans within each grade, consistent with borrower-risk sorting and a maturity-risk component. Figure 1(b) shows a steep but noisy positive relationship between interest rates and default rates, consistent with risk-based pricing; because rates are assigned to observably riskier applicants, raw differences across rate levels are not immediately causal.

Risk also varies along joint credit dimensions. Figure 4 reports default rates by credit grade and FICO band separately for 36- and 60-month loans: within each term, default generally rises as grade worsens and FICO weakens, and the 60-month panel is uniformly riskier, especially in lower-grade segments. Together, these descriptives establish two empirical facts that guide the remainder of the paper: (i) default risk is highly predictable from origination-time information in economically meaningful ways, which makes screening a natural objective; and (ii) maturity and interest rate are strongly related to default in the raw data, but are also plausibly chosen *because* of borrower risk, which motivates a cautious, adjustment-based treatment-effect extension.

3 Methods

Prediction problem and information set. Our goal is to support *screening at loan origination*. For each loan i we observe origination-time features \mathbf{x}_i and a binary default outcome $y_i \in \{0, 1\}$ (default = 1) from the harmonized status field. With $\mathcal{Z} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, models output scores

$$\hat{p}_i = \hat{p}(\mathbf{x}_i) \approx \Pr(Y_i = 1 \mid \mathbf{x}_i),$$

interpreted as estimated default probabilities and used to rank borrowers by risk. A threshold rule at $\tau \in (0, 1)$ is $\hat{y}(\tau) = \mathbb{I}(\hat{p} \geq \tau)$.

Training, testing, and progressive evaluation logic. All models are compared on a common held-out test set. Any data-dependent preprocessing and tuning (including cross-validation and bootstrap aggregation) is performed on the training set only; the final fitted model is then applied once to the test set to obtain out-of-sample scores \hat{p}_i . We proceed from linear benchmarks (logit/elastic net) to single trees, then to bagging and random forests, and end with a DML extension. Diagnostics in Section 4 are defined once and applied uniformly; for bagging and forests we also use out-of-bag (OOB) diagnostics as an internal training check. All procedures are implemented in R and documented in the [code appendix](#); the bagging routine and the DML estimator are coded from scratch (no high-level implementations).

Linear probability models as stable baselines: logit and elastic net. Logistic regression links default risk to a linear index,

$$\Pr(Y_i = 1 \mid \mathbf{x}_i) = \sigma(\eta_i), \quad \eta_i = \beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta}, \quad \sigma(t) = \frac{1}{1 + e^{-t}},$$

so features enter additively in log-odds. This is interpretable and stable, but can underfit with many correlated predictors and nonlinearities.

We therefore also fit elastic-net penalized logistic regression. Writing $\hat{p}_i(\beta_0, \boldsymbol{\beta}) = \sigma(\beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta})$, the estimator solves

$$\min_{\beta_0, \boldsymbol{\beta}} \frac{1}{n} \sum_{i=1}^n [-y_i \log \hat{p}_i - (1 - y_i) \log(1 - \hat{p}_i)] + \lambda \left(\alpha \|\boldsymbol{\beta}\|_1 + \frac{1-\alpha}{2} \|\boldsymbol{\beta}\|_2^2 \right),$$

where $\lambda \geq 0$ controls shrinkage and $\alpha \in [0, 1]$ mixes ℓ_1 (sparsity/selection) with ℓ_2 (stability under collinearity). Predictors are standardized prior to estimation. Elastic net is typically effective when signals are spread across many correlated variables and the decision boundary is approximately linear.

Classification trees (CART) as flexible but high-variance learners. Trees partition the feature space into regions $\{R_m\}_{m=1}^M$ and predict a region-specific default probability (the within-node default rate). Splits are chosen greedily to reduce impurity, e.g. the Gini index

$$G = 1 - p_0^2 - p_1^2 = 2p_1(1 - p_1),$$

where p_1 is the node default share. Trees capture nonlinearities and interactions and handle mixed variable types, but are high-variance: small data changes can alter early splits and the entire tree.

Bagging (bootstrap aggregating). Bagging reduces variance by averaging many trees. Fix B bootstraps. For $b = 1, \dots, B$: (i) draw a bootstrap sample \mathcal{Z}_b^* of size n from \mathcal{Z} ; (ii) fit a classification tree to obtain $\hat{p}_b^*(\mathbf{x})$; (iii) store bootstrap membership indicators for OOB prediction. The bagged probability (soft voting) is

$$\hat{p}_{\text{bag}}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B \hat{p}_b^*(\mathbf{x}), \quad \hat{y}_{\text{bag}}(\tau) = \mathbb{I}(\hat{p}_{\text{bag}}(\mathbf{x}) \geq \tau),$$

which preserves probabilistic information and typically improves ranking/calibration relative to hard voting.

For observation i , let \mathcal{B}_i be the set of trees where i is out-of-bag. The OOB score is

$$\hat{p}_i^{\text{OOB}} = \frac{1}{|\mathcal{B}_i|} \sum_{b \in \mathcal{B}_i} \hat{p}_b^*(\mathbf{x}_i),$$

providing an internal performance estimate without a separate validation set. We also track OOB error as a function of the number of trees to confirm stabilization.

Random forests. If a few strong predictors dominate, bagged trees can become too similar. Random forests inject additional randomness by selecting m candidate predictors at each split (rather than all p). With B bootstrap trees, the forest probability is

$$\hat{p}_{\text{rf}}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B \hat{p}_{b,\text{rf}}^*(\mathbf{x}),$$

with OOB diagnostics defined as above. For screening, forests retain tree flexibility while typically improving variance reduction under correlated predictors, at the cost of interpretability.

Diagnostics and performance measures. Let \mathcal{T} be the common test set of size n_{test} . For threshold τ , define $\hat{y}_i(\tau) = \mathbb{I}(\hat{p}_i \geq \tau)$ and confusion-matrix counts $\text{TP}(\tau), \text{FP}(\tau), \text{TN}(\tau), \text{FN}(\tau)$ in the usual way. We report:

- *Misclassification error and accuracy,*

$$\text{Error}(\tau) = \frac{\text{FP}(\tau) + \text{FN}(\tau)}{n_{\text{test}}}, \quad \text{Accuracy}(\tau) = 1 - \text{Error}(\tau),$$

which are intuitive but can mislead under class imbalance.

- *Brier score* for probability accuracy and calibration,

$$\text{Brier} = \frac{1}{n_{\text{test}}} \sum_{i \in \mathcal{T}} (y_i - \hat{p}_i)^2,$$

which penalizes both poor ranking and miscalibration (confidently wrong scores are costly).

- *ROC curve and AUC.* For $\tau \in (0, 1)$,

$$\text{TPR}(\tau) = \frac{\text{TP}(\tau)}{\text{TP}(\tau) + \text{FN}(\tau)}, \quad \text{FPR}(\tau) = \frac{\text{FP}(\tau)}{\text{FP}(\tau) + \text{TN}(\tau)}.$$

The ROC curve is $\{(\text{FPR}(\tau), \text{TPR}(\tau))\}$ and

$$\text{AUC}_{\text{ROC}} = \int_0^1 \text{TPR}(u) du,$$

interpretable as the probability a randomly chosen default receives a higher score than a randomly chosen non-default.

- *Precision–Recall (PR) curve.* At threshold τ :

$$\text{Precision}(\tau) = \frac{\text{TP}(\tau)}{\text{TP}(\tau) + \text{FP}(\tau)}, \quad \text{Recall}(\tau) = \text{TPR}(\tau),$$

useful under imbalance because precision answers: among flagged loans, what fraction actually default?

- *Screening curves (cumulative gains and lift).* Sort test loans by decreasing \hat{p}_i and let $u \in (0, 1]$ be the screened population share:

$$G(u) = \frac{\sum_{j=1}^{\lfloor un_{\text{test}} \rfloor} y_{(j)}}{\sum_{j=1}^{n_{\text{test}}} y_{(j)}}, \quad L(u) = \frac{G(u)}{u},$$

which quantify how many defaults are captured by screening the top u fraction and how much this improves over random screening.

Double/Debiased Machine Learning (DML). As a causal add-on, we study how default risk is associated with binary origination-time treatments D (60-month vs. 36-month term; and a high-interest indicator) after flexibly adjusting for remaining covariates \mathbf{X} . Because contract terms may be endogenous to borrower risk, naive differences can be misleading (Emekter et al., 2015; Serrano-Cinca et al., 2015). Following (Chernozhukov et al., 2016, 2018), we use the partially linear model

$$Y = \theta_0 D + g_0(\mathbf{X}) + \varepsilon, \quad D = m_0(\mathbf{X}) + v,$$

where $g_0(\mathbf{X})$ captures baseline default risk and $m_0(\mathbf{X})$ the propensity of treatment given \mathbf{X} ; θ_0 summarizes the average adjusted association between D and Y .

The key idea is to estimate θ_0 by “regressing residuals on residuals.” We estimate g_0 and m_0 via cross-fitted elastic-net logistic regression: split the sample into K folds; for each fold k , fit nuisance models on I_k^c and predict $\hat{g}^{(-k)}(\mathbf{X}_i)$ and $\hat{m}^{(-k)}(\mathbf{X}_i)$ for $i \in I_k$. This yields orthogonalized residuals

$$\tilde{Y}_i = Y_i - \hat{g}(\mathbf{X}_i), \quad \tilde{D}_i = D_i - \hat{m}(\mathbf{X}_i),$$

where each observation is predicted by models trained without that observation. The DML estimator is the least-squares slope

$$\hat{\theta}_{\text{DML}} = \frac{\sum_{i=1}^n \tilde{D}_i \tilde{Y}_i}{\sum_{i=1}^n \tilde{D}_i^2}.$$

Because \tilde{D}_i is approximately uncorrelated with functions of \mathbf{X}_i , first-stage errors affect $\hat{\theta}_{\text{DML}}$ only to second order (Neyman orthogonality), reducing regularization bias relative to naive plug-in approaches.

For comparison we also report the naive difference in default rates,

$$\hat{\theta}_{\text{naive}} = \bar{Y}_{D=1} - \bar{Y}_{D=0},$$

so the gap between $\hat{\theta}_{\text{naive}}$ and $\hat{\theta}_{\text{DML}}$ indicates how strongly observable characteristics mediate the relationship between contract terms and default.

4 Results

Having set up the prediction problem and modelling framework, we now turn to the empirical performance of the six classifiers and to the treatment-effect extension on contract terms.

Overall predictive performance. Table 1 reports test-set error, accuracy, AUC, and Brier scores on the common hold-out sample. All models achieve similar accuracy (around 82%) in a setting where the unconditional default rate is roughly 18.5%. The carefully tuned logit and elastic-net models provide a very strong baseline: they attain the highest or near-highest AUC (about 0.73) and the lowest Brier scores. Random forests perform comparably, while both bagging variants and especially CART lag somewhat, with lower AUC and slightly higher error. Thus, there is clear signal in origination-time characteristics, but moving from a well-specified penalized linear model to more complex tree ensembles yields only modest gains in standard predictive metrics.

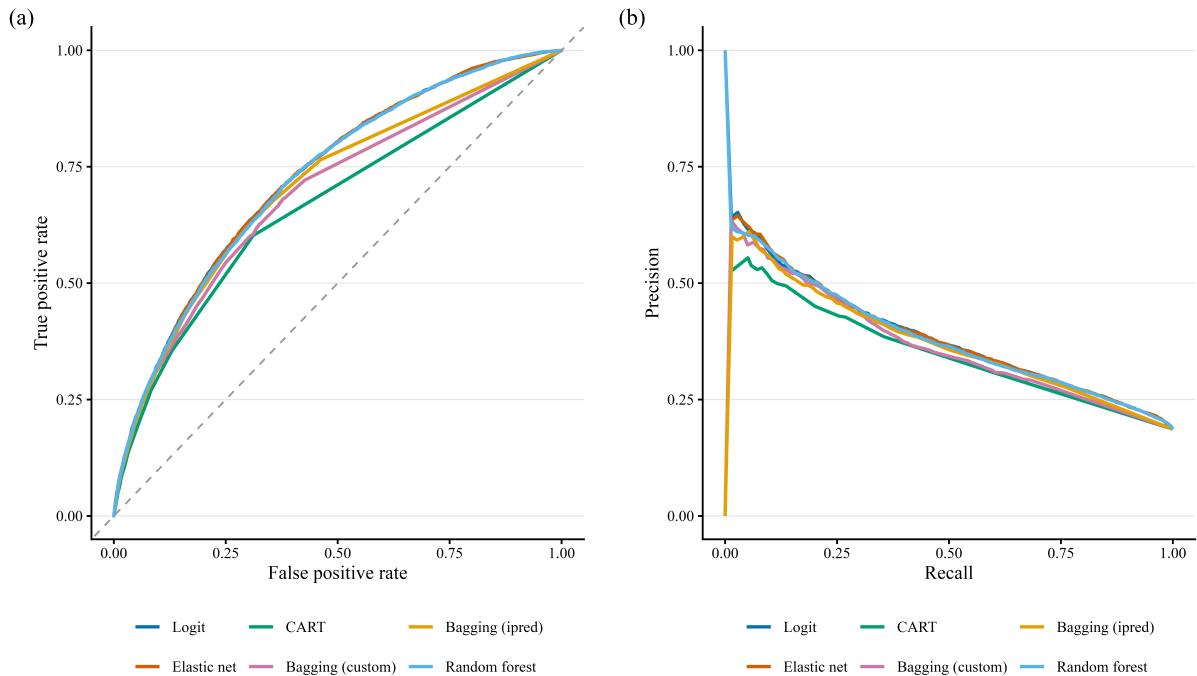
These patterns align with prior findings that logistic regression can match more complex techniques in P2P lending and that model performance is often similar across methods (Jin & Zhu, 2015; Serrano-Cinca et al., 2015).

Ranking quality and screening implications. Because screening is fundamentally about ranking rather than choosing a single threshold, Figure 2 and Figure 3 focus on ordering performance. The ROC and precision-recall curves in Figure 2 are broadly similar for logit, elastic net, and random forest, which dominate bagging and CART over most thresholds.

Table 1: Test-set performance comparison (common test set)

Model	Error	Accuracy	AUC	Brier	TP	FP	TN	FN
Logit	0.182	0.818	0.727	0.135	486	322	28,969	6,223
Elastic net	0.182	0.818	0.727	0.135	488	314	28,977	6,221
Random forest	0.182	0.818	0.725	0.136	415	273	29,018	6,294
Bagging (ipred)	0.183	0.817	0.704	0.152	526	394	28,897	6,183
Bagging (custom)	0.183	0.817	0.690	0.138	402	288	29,003	6,307
CART	0.184	0.816	0.666	0.140	552	484	28,807	6,157

This means that for a given tolerance for false positives, these models recover a larger share of eventual defaulters and maintain higher precision when targeting the riskiest slice of applicants.

**Figure 2:** ROC and precision–recall curves by model.

The cumulative gains curve in Figure 3 translates this into a screening scenario. If an investor can manually review only the top 20–30% of applicants, the best-performing models capture roughly half of all future defaults, versus about 20–30% under random screening (the 45-degree line). Random forests sit slightly above logit and elastic net, but the gap is small; the main contrast is between any of these models and naive screening. The practical message is therefore twofold: even an interpretable, well-tuned linear model substantially improves who to flag first, while tree ensembles offer at most incremental gains in this dataset.

Predicted risk distributions. Figure 6 plots predicted default probabilities by true outcome for each classifier. Under logit and elastic net, non-defaults concentrate at low predicted risk while defaults shift toward the right tail, yielding clear but overlapping separation. Tree-based models produce more irregular or multi-modal shapes because they average discrete terminal-node probabilities, but

the qualitative ranking remains: defaulters are systematically assigned higher scores. The bagging models place more mass at very low probabilities, consistent with conservative predictions and their slightly worse Brier scores. Overall, the distributions confirm meaningful risk sorting, with calibration differences across algorithms.

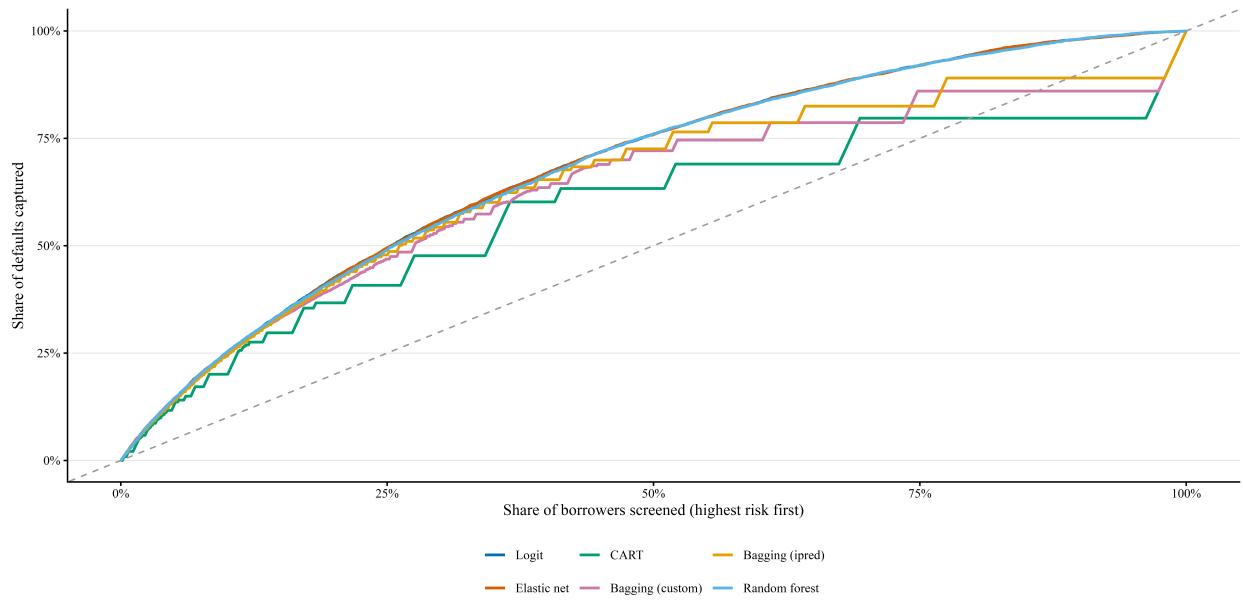


Figure 3: Cumulative gains / lift curves by model.

Tree-ensemble diagnostics. Table 2 and the out-of-bag (OOB) error path in Figure 7 provide additional checks for the ensemble methods. For the custom bagging implementation, OOB misclassification error declines with the number of trees and stabilizes around 18.4% at roughly 70 trees, slightly improving on the single-value OOB error of the off-the-shelf `ipred` bagging model. Random forests attain the lowest OOB error and highest OOB AUC among the tree-based methods, consistent with their held-out test performance. These diagnostics support our modelling choices: the custom bagging routine behaves as expected and the random forest is a strong nonlinear benchmark, even though its gains over simpler models remain limited.

Table 2: Out-of-bag (OOB) performance summary

Method	Number of Trees	OOB error (misclassification)	OOB AUC	OOB Brier
Bagging (custom)	80	0.184	0.677	0.140
Bagging (ipred)	80	0.185		
Random forest (ranger)	300	0.137		

Double machine learning extension. Finally, we examine how default risk differs across contract terms once we adjust for observed borrower characteristics. In Table 3, the naive difference-in-means implies that 60-month loans default about 18.3 percentage points more often than 36-month loans, with a tight confidence interval. By contrast, the orthogonalised DML estimator is numerically unstable: the point estimate becomes very large and the 95% confidence interval is extremely wide, spanning negative values and implausibly large positives. The reason is overlap: conditional on the rich origination covariates, term length is almost perfectly predictable from X , leaving too little

residual variation in treatment. Econometrically, this near-violation of overlap means the term-length DML estimate is effectively not identified in our design, so we avoid causal claims about moving from 36 to 60 months.

Table 3: *Naive vs DML estimates for two treatments*

Treatment	Estimator	Effect (pp) [95% CI]
60-month term (vs 36-month)	Naive diff-in-means	18.28 [17.59, 18.96]
60-month term (vs 36-month)	DML (elastic net, K-fold)	67.69 [-343.62, 479.01]
High interest (\geq 75th pct)	Naive diff-in-means	18.19 [17.50, 18.87]
High interest (\geq 75th pct)	DML (elastic net, K-fold)	1.74 [-4.87, 8.34]

The picture is clearer for pricing. High-interest loans (above the 75th percentile) have about 18.2 percentage points higher default rates in the raw data, consistent with the interest-rate gradient in [Figure 1](#). After flexible adjustment via DML, the estimated effect shrinks to roughly 1.7 percentage points and the confidence interval includes zero. Hence most of the naive association appears driven by selection on observables: borrowers who look riskier are both more likely to receive high rates and more likely to default. Put differently, conditional on observables, LendingClub appears primarily to *price* risk rather than *create* it through high interest rates.

5 Conclusion

This paper examined how well default risk can be predicted *at loan origination* in a large P2P consumer lending platform, and what this implies for screening and contract terms. Using LendingClub loan-level data with an origination-only information set, we compared logistic regression and elastic net with tree-based methods. Across [Table 1](#) and [Figure 2](#), default risk is highly predictable from standard borrower and loan characteristics: a well-tuned (penalized) logit already attains strong accuracy, AUC, and calibration. Operationally, [Figure 3](#) shows that reviewing the riskiest 20–30% of applicants captures about half of all subsequent defaults, far above random selection, so simple, transparent models already provide substantial “*risk at first sight*” signals.

Tree-based models behave as expected but add only modest improvements. Bagging and random forests deliver slightly better ranking (and stable OOB diagnostics in [Table 2](#)), yet their gains over elastic net are small relative to the increase in model complexity, implying only marginal benefits from moving to more opaque ensembles.

The double machine learning extension highlights the limits of causal interpretation. While naive comparisons link 60-month terms and high interest rates to much higher default, DML shrinks the high-interest association to a small, statistically inconclusive estimate and yields an unstable, very wide interval for term length ([Table 3](#)). This is consistent with strong selection on observables: riskier borrowers are more likely to receive longer terms or higher rates and to default, so contract terms largely *price* risk rather than clearly *cause* it. Our analysis is limited to one platform, cohort, and observables; future work could use richer data and quasi-experimental variation in pricing or term policies to sharpen causal claims. Within these bounds, the main message is that in P2P consumer lending, much of default risk is visible at origination, and exploiting this information through well-calibrated, transparent models can meaningfully improve screening performance while keeping distributional impacts on access to credit more interpretable.

References

- Byanjankar, A., Heikkilä, M., & Mezei, J. (2015). Predicting credit risk in peer-to-peer lending: A neural network approach. *2015 IEEE Symposium Series on Computational Intelligence*, 719–725.
- Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., & Robins, J. (2016). Double/debiased machine learning for treatment and causal parameters. <https://arxiv.org/abs/1608.00060>
- Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W. K., & Robins, J. (2018). Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1), C1–C68. <https://doi.org/10.1111/ectj.12097>
- Emekter, R., Tu, Y., Jirasakuldech, B., & Lu, M. (2015). Evaluating credit risk and loan performance in online peer-to-peer (p2p) lending. *Applied Economics*, 47(1), 54–70.
- Jin, Y., & Zhu, Y. (2015). A data-driven approach to predict default risk of loan for online peer-to-peer (p2p) lending [Also appearing in 2015 Fifth International Conference on Communication Systems and Network Technologies]. *Procedia Computer Science*.
- LendingClub Corporation. (2014). Form s-1 registration statement [As filed Aug. 27, 2014. U.S. Securities and Exchange Commission].
- Lenz, R. (2016). Peer-to-peer lending: Opportunities and risks. *European Journal of Risk Regulation*, 7(4), 688–700.
- Serrano-Cinca, C., Gutiérrez-Nieto, B., & López-Palacios, L. (2015). Determinants of default in p2p lending. *PloS one*, 10(10), e0139427.

Appendix

A Figures

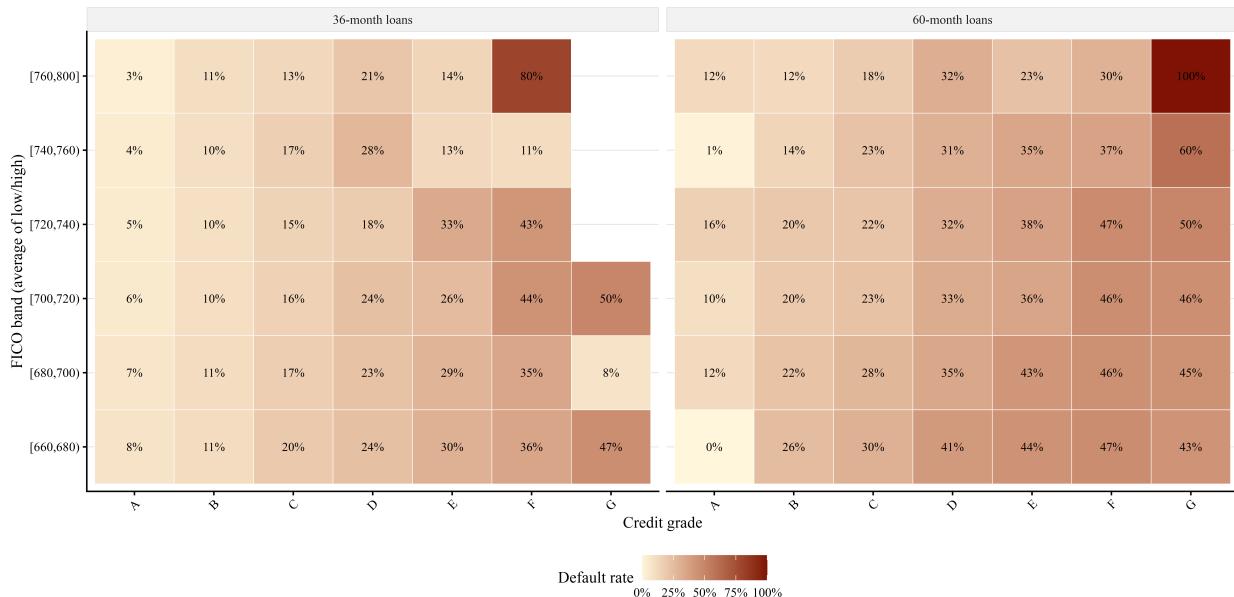


Figure 4: Default-rate heatmaps by credit grade and FICO band, shown for 36- and 60-month loans.

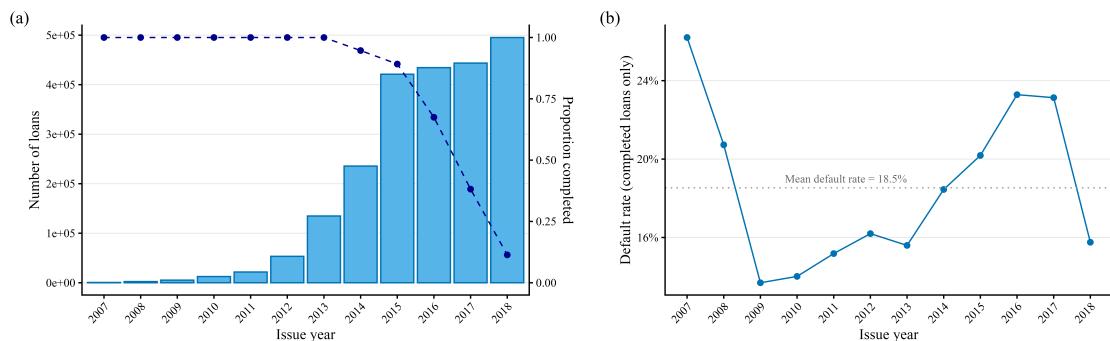


Figure 5: Loan issuance volume and performance trends (2007-2018)

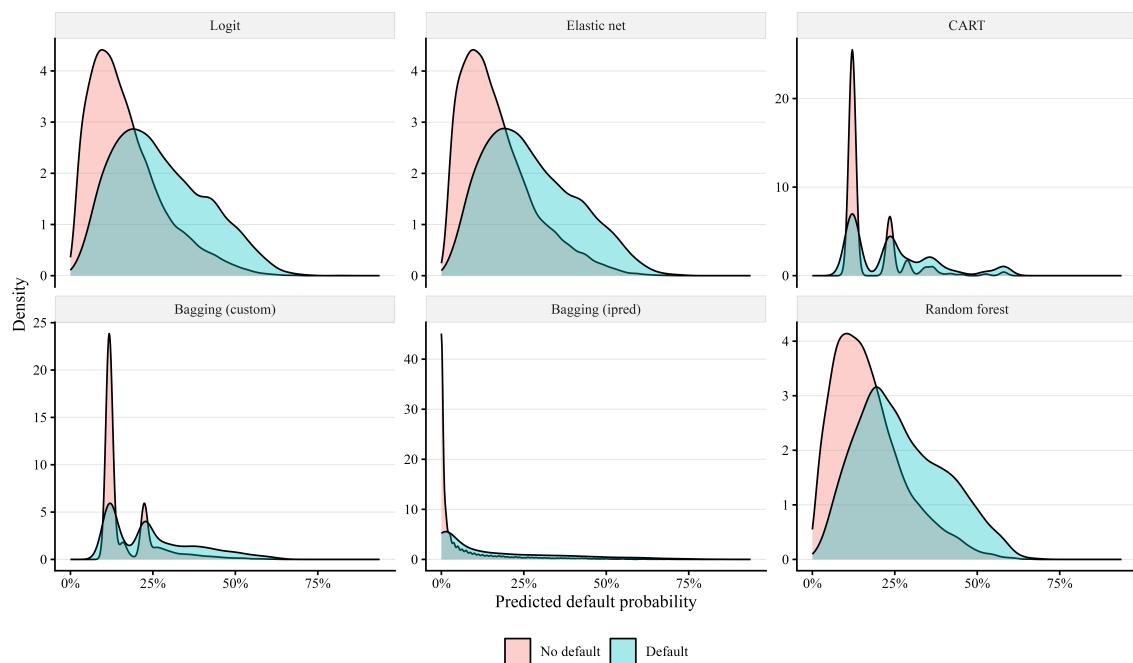


Figure 6: Predicted default-probability distributions by true outcome across models

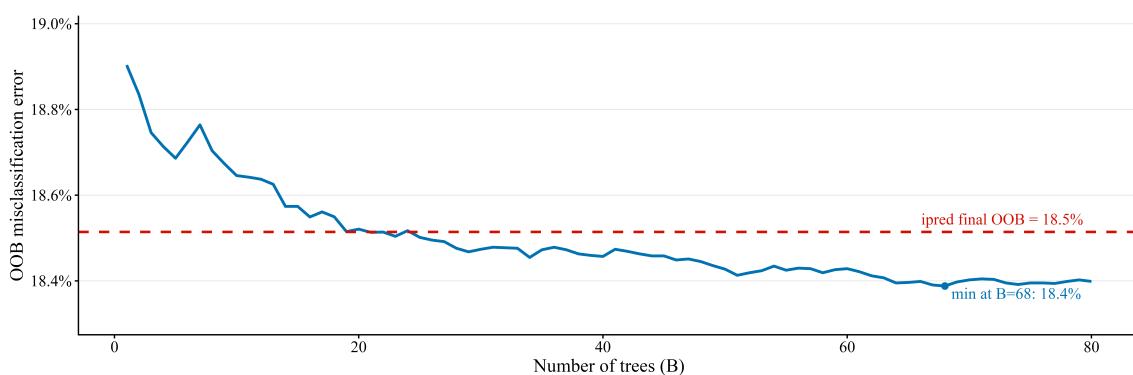


Figure 7: Custom bagging out-of-bag misclassification error versus number of trees

B Tables

Table 4: Data preparation audit trail

Action	Rows	Rows removed	Cols	Cols dropped	Cols added
Load raw data	2,260,701		151		
Create issue year + status	2,260,668	33	154	0	3
Cohort filter (2012–2015)	786,818	1,473,850	156	0	2
Stratified sample by (year × default)	119,995	666,823	156	0	0
Drop ID/text/leakage vars	119,995	0	112	44	0
Feature engineering + factor cleanup	119,995	0	113	0	1
Final modelling table	119,995	0	32	81	0

C Full R Implementation

Below is the complete R code used in this study. It imports the raw LendingClub data, constructs the origination-only modelling cohort (including cleaning the loan-status outcome, restricting to 2012–2015 completed loans to mitigate right-censoring, dropping identifiers and post-origination leakage fields, and applying a stratified downsample for computational tractability), and builds the final design matrices used for estimation. The code then fits and evaluates the full set of classifiers reported in the paper (logit, elastic net, CART, bagging, and random forest), reproducing all performance tables and screening-focused figures (ROC, precision–recall, cumulative gains, and predicted-risk distributions), and implements the double machine learning extensions for maturity and high-interest treatments alongside naive benchmarks. The script is organised in clearly labelled blocks that mirror the paper’s workflow (data preparation, descriptives, modelling, diagnostics, and DML), so results can be replicated or extended with minimal modification.

```

1 ## =====
2 ## Risk at First Sight - ML Final Assignment
3 ## -----
4 ## Student Number: 756486 - Dec 2025
5 ## =====
6
7 ## =====
8 ## 0. Set Up
9 ## =====
10
11 ## Packages used
12 ## stargazer : exporting LaTeX tables
13 ## readr      : fast CSV import (read_csv)
14 ## dplyr     : data wrangling (mutate/filter/summarise)
15 ## tidyverse  : reshaping (pivot_longer/pivot_wider/unnest)
16 ## stringr    : string handling (str_detect/str_sub/str_to_upper)
17 ## forcats   : factor utilities (lump levels, NA level)
18 ## tibble     : tidy tables (tibble)
19 ## ggplot2    : plotting
20 ## patchwork  : combine ggplots (plot1 / plot2)
21 ## scales     : axis formatting (percent_format)
22 ## glmnet     : elastic net / ridge / lasso (cv.glmnet)
23 ## pROC       : ROC/AUC + DeLong tests (roc, auc, ci.auc, roc.test)
24 ## rpart      : CART trees
25 ## ipred      : bagging benchmark (ipred::bagging)
26 ## ranger     : random forest (probability=TRUE + OOB error)

```

```

27 pkgs <- c(                                     # list of packages required
28   → for the full analysis
29   "readr", "dplyr", "tidyverse", "stringr", "forcats", "tibble",
30   "ggplot2", "patchwork", "scales",
31   "glmnet", "pROC", "rpart", "ipred", "ranger"
32 )
33
34 to_install <- pkgs[!pkgs %in% rownames(installed.packages())]      # check which required
35   → packages are not yet installed
36 if (length(to_install) > 0) install.packages(to_install)           # install missing packages so
37   → the script runs on a clean machine
38
39 invisible(lapply(pkgs, library, character.only = TRUE))            # load all packages quietly
40   → for use throughout the script
41
42 base_seed <- 1363                                                 # set the master seed for
43   → reproducible results across the project
44 set.seed(base_seed)                                              # fix the random number
45   → generator state using the master seed
46
47 theme_academic <- function(base_size = 11, base_family = "Times New Roman") { # define a
48   → consistent ggplot theme for all figures
49   theme_classic(base_size = base_size, base_family = base_family) +          # use a clean
50     → classic theme as the base style
51   theme(
52     plot.title      = element_text(face = "bold", size = base_size + 1, hjust = 0),    #
53       → left-aligned bold title
54     plot.subtitle   = element_text(size = base_size, hjust = 0),                      #
55       → left-aligned subtitle
56     plot.caption    = element_text(size = base_size - 1, hjust = 0, color = "grey40"), # small
57       → caption
58     axis.title      = element_text(face = "plain"),                                     # standard
59       → axis titles
60     axis.text       = element_text(color = "black"),                                    # readable
61       → axis labels
62     axis.ticks      = element_line(color = "black"),                                 # visible
63       → axis ticks
64     panel.grid.major.y = element_line(color = "grey90", linewidth = 0.3),             # subtle
65       → horizontal grid for readability
66     panel.grid.major.x = element_blank(),                                         # remove
67       → vertical grid lines
68     panel.grid.minor = element_blank(),                                         # remove
69       → minor grid lines
70     legend.position = "bottom",                                              # place
71       → legend under the plot
72     legend.title    = element_text(face = "plain"),                                # keep
73       → legend title simple
74     legend.key      = element_blank(),                                         # remove
75       → legend key background
76     strip.background = element_rect(fill = "grey95", color = "grey80", linewidth = 0.3), # light
77       → facet strip background
78     strip.text      = element_text(face = "plain"),                                # simple
79       → facet labels
80     plot.margin     = margin(10, 10, 10, 10)                                     # add
81       → consistent outer margins
82   )
83 }
84
85 theme_set(theme_academic())                                         # apply the academic theme
86   → globally to all ggplot figures

```

```

64 ## =====#
65 ## 1. Load full Kaggle LendingClub file
66 ## =====#
67 ## =====#
68
69 data_dir <- "data"                                # folder name where the raw
   → Kaggle file is stored
70 if (!dir.exists(data_dir)) dir.create(data_dir)      # create the folder if it does
   → not already exist
71
72 lc_file <- file.path(data_dir, "accepted_2007_to_2018Q4.csv")    # full path to the expected
   → LendingClub CSV file
73
74 lc_raw <- readr::read_csv(lc_file, show_col_types = FALSE)      # read the full dataset and
   → suppress column type printing
75
76 cat("\nLoaded LendingClub data\n")                      # print a small header so logs
   → are easy to scan
77 cat("Rows : ", nrow(lc_raw), "\n")                      # print number of rows in the
   → raw dataset
78 cat("Columns:", ncol(lc_raw), "\n\n")                  # print number of columns in the
   → raw dataset
79
80 # cat("Top 20 distinct loan_status values:\n")          # print a header for the loan
   → status frequency table
81 # print(head(sort(table(lc_raw$loan_status), decreasing = TRUE), 20)) # show the most common loan
   → status labels for auditing
82
83 ## =====#
84 ## 2. Derive issue year and coarse outcome categories
85 ## =====#
86
87 status_completed <- c(                                # loan status labels treated
   → as completed outcomes
     "Fully Paid",
     "Charged Off",
     "Default",
     "Does not meet the credit policy. Status:Fully Paid",
     "Does not meet the credit policy. Status:Charged Off"
)
88
89 status_in_progress <- c(                            # loan status labels treated
   → as still active or delinquent
     "Current",
     "In Grace Period",
     "Late (16-30 days)",
     "Late (31-120 days")
)
90
91 lc_year <- lc_raw %>%
   → LendingClub file and create key cohort variables
92 mutate(
93   issue_year = str_sub(issue_d, -4L),                # extract the calendar year
   → from the issue date string
94   issue_year = na_if(issue_year, ""),                 # convert empty year strings
   → to missing values
95   status_group = case_when(                          # map detailed loan_status
   → into completed, in_progress, or other
     loan_status %in% status_completed ~ "completed",
     loan_status %in% status_in_progress ~ "in_progress",
     TRUE ~ "other"
)
96
97
98
99
100
101
102
103
104
105
106
107
108
109

```

```

110 ),
111 default_prelim = case_when(
112   ↪ indicator for completed outcomes only
113   loan_status %in% c(
114     "Charged Off",
115     "Default",
116     "Does not meet the credit policy. Status: Charged Off",
117     "Does not meet the credit policy. Status:Charged Off"
118   ) ~ 1L,                                     # set a preliminary default
119   ↪ charge-offs and defaults
120   loan_status %in% c(
121     "Fully Paid",
122     "Does not meet the credit policy. Status:Fully Paid"
123   ) ~ 0L,                                     # default equals one for
124   ↪ fully paid loans
125   TRUE ~ NA_integer_                          # default equals zero for
126   ↪ non-completed statuses
127 )
128 ) %>%
129 filter(!is.na(issue_year))                  # drop rows with missing or
130   ↪ malformed issue dates
131
132 ## =====
133 ## 3. Year-by-year summary: volume, censoring, default rates
134 ## -----
135 ## I have done this part and the following to see what years to keep
136 ## -----
137
138 year_summary <- lc_year %>%
139   ↪ with issue_year and status_group          # start from the dataset
140   group_by(issue_year) %>%
141   ↪ separately for each issuance year        # compute summaries
142   summarise(
143     n_total      = n(),                         # total number of loans
144     ↪ issued in that year
145     n_completed = sum(status_group == "completed"),    # number of loans with
146     ↪ completed outcomes
147     n_in_prog   = sum(status_group == "in_progress"),  # number of loans still
148     ↪ active or delinquent
149     n_other     = sum(status_group == "other"),        # number of loans with
150     ↪ other or ambiguous statuses
151     prop_completed = n_completed / n_total,           # fraction of loans with
152     ↪ completed outcomes
153     prop_in_prog = n_in_prog / n_total,              # fraction of loans still
154     ↪ in progress
155     default_rate_completed = mean(
156       ↪ completed loans only
157       default_prelim[status_group == "completed"],    # default rate among
158       ↪ outcomes for comparability
159       na.rm = TRUE                                    # restrict to completed
160       ↪ from non-completed statuses
161     ),                                              # ignore missing values
162     .groups = "drop"                                # return an ungrouped data
163     ↪ frame after summarising
164   ) %>%
165   arrange(issue_year)                            # order rows
166   ↪ chronologically by issuance year
167
168 ## =====
169 ## 4. Visual inspection: volume & default rate by year

```

```

153 ## =====
154
155 year_summary <- year_summary %>%
156   # start from the year-level
157   # summary table
158   mutate(issue_year = factor(issue_year, levels = issue_year))      # convert issue_year to an
159   # ordered factor for a clean x-axis
160
161 mean_default_rate <- mean(year_summary$default_rate_completed, na.rm = TRUE) # compute the
162   # average completed-loan default rate across years
163 x_label_year <- year_summary$issue_year[ceiling(nrow(year_summary) / 2)]       # choose a
164   # mid-sample year for placing the reference-line label
165
166 ## -----
167 ## 4.1 Total loans per year + proportion completed (dual axis)
168 ## -----
169
170 plot_year_counts <- year_summary %>%
171   # use the year-level table
172   # to plot annual volumes and completion share
173   ggplot(aes(x = issue_year)) +          # set issue year on the
174     # x-axis
175   geom_col(                                # draw bars for total number
176     # of loans per year
177     aes(y = n_total),                      # draw bars for total number
178     fill = "#56B4E9", color = "#0072B2")
179   ) +
180   geom_line(                                # overlay the completed
181     # share scaled to the count axis
182     aes(y = prop_completed * max(n_total), group = 1),
183     linetype = "dashed",
184     color = "#00008B")
185   ) +
186   geom_point(                               # add points to highlight
187     # each yearly completed share
188     aes(y = prop_completed * max(n_total)),
189     size = 1.7,
190     color = "#00008B")
191   ) +
192   scale_y_continuous(                      # define the primary count
193     # axis and a secondary proportion axis
194     name = "Number of loans",
195     sec.axis = sec_axis(
196       ~ . / max(year_summary$n_total),      # rescale back to
197       # proportions using the maximum annual volume
198       name = "Proportion completed"
199     )
200   ) +
201   labs(                                    # label the x-axis
202     x = "Issue year"
203   ) +
204   theme(
205     axis.text.x = element_text(angle = 45, hjust = 1)      # rotate year labels for
206     # readability
207   )
208
209 ## -----
210 ## 4.2 Default rate among completed loans by year
211 ## -----
212
213 plot_year_default <- year_summary %>%
214   # use the year-level table
215   # to plot default rates among completed loans

```

```

200  ggpplot(aes(x = issue_year, y = default_rate_completed, group = 1)) + # connect points in
201    geom_line(color = "#0072B2") +                                     # draw the default-rate
202    geom_point(size = 1.7, color = "#0072B2") +                      # add points at each year
203    geom_hline(                                                       # add a reference line at
204      yintercept = mean_default_rate,                                the mean completed-loan default rate
205      linetype = "dotted",                                          
206      color = "#999999"                                           
207    ) +
208    annotate(                                                       # label the reference line
209      "text",
210      x = x_label_year,
211      y = mean_default_rate,
212      label = paste0(
213        "Mean default rate = ",
214        scales::percent(mean_default_rate, accuracy = 0.1)
215      ),
216      vjust = -0.7,
217      hjust = 0.7,
218      size = 3,
219      color = "#666666"
220    ) +
221    scale_y_continuous(                                              # format the y-axis as
222      "percentages"
223      labels = scales::percent_format(accuracy = 1),
224      name = "Default rate (completed loans only)"
225    ) +
226    labs(                                                       # label the x-axis
227      x = "Issue year"
228    ) +
229    theme(
230      axis.text.x = element_text(angle = 45, hjust = 1)           # rotate year labels for
231      # readability
232
233  ## -----
234  ## 4.3 Combine as horizontal subplots (a) and (b)
235  ## -----
236  plot_year_panel <- plot_year_counts + plot_year_default +          # combine the two plots into
237    patchwork::plot_layout(ncol = 2) +                                    # a single panel
238    patchwork::plot_annotation(                                         # arrange the plots
239      # referencing in the paper
240      tag_levels = "a",
241      tag_prefix = "(",
242      tag_suffix = ")"
243
244  print(plot_year_panel)                                              # add panel tags for figure
245
246  ggplot2::ggsave(                                                    # display the combined panel
247    # disk at high resolution
248    filename = "plot_year_panel.png",
249    plot      = plot_year_panel,

```

```

249   width     = 12, height = 4, units = "in",
250   dpi       = 1000
251 )
252
253 ## =====
254 ## 5. Cohort selection: 2012-2015, completed loans, individual apps
255 ## =====
256
257 years_keep <- c("2012", "2013", "2014", "2015")                      # restrict to years with
258   ↪  large volumes and moderate censoring
259
260 cat("\napplication_type by year (for 2012-2015):\n")                  # print application_type
261   ↪  counts to confirm the prevalence of joint loans
262 print(
263   lc_year %>%
264     dplyr::filter(issue_year %in% years_keep) %>%
265       ↪ cohort years
266     dplyr::count(issue_year, application_type) %>%
267       ↪ by year
268     dplyr::arrange(issue_year, dplyr::desc(n))
269       ↪ frequency
270 )
271
272 lc_cohort <- lc_year %>%                                         # build the main modeling
273   ↪ cohort from the cleaned year-level dataset
274   dplyr::filter(
275     issue_year %in% years_keep,
276       ↪ target years
277     loan_status %in% status_completed,
278       ↪ outcomes for comparable default labeling
279     is.na(application_type) |
280       ↪ types rather than dropping observations
281     stringr::str_to_upper(application_type) != "JOINT"
282       ↪ applications based on case-insensitive matching
283   ) %>%
284   dplyr::mutate(
285     default = default_prelim,                                         # set the final binary
286       ↪ default outcome used in modeling
287     term_60 = ifelse(stringr::str_detect(term, "60"), 1L, 0L)        # create a 60-month term
288       ↪ indicator from the term string
289   ) %>%
290   dplyr::filter(
291     !is.na(default),                                              # require a valid default
292       ↪ label
293     !is.na(loan_amnt),                                            # require loan amount at
294       ↪ origination
295     !is.na(int_rate),                                             # require interest rate at
296       ↪ origination
297     !is.na(annual_inc),                                           # require annual income at
298       ↪ origination
299     !is.na(dtif)
300       ↪ origination
301   )
302
303 # cat("\nCohort (2012-2015, completed, non-joint) summary:\n")          # print a cohort summary for
304   ↪ manual verification
305 # cat("Rows in lc_cohort:", nrow(lc_cohort), "\n")                         # print number of
306   ↪ observations in the final cohort
307 # cat("Default distribution:\n")                                         # print default counts and
308   ↪ shares for a quick sanity check
309 # print(table(lc_cohort$default))                                         # show default counts

```

```

290 # print(prop.table(table(lc_cohort$default)))
291 # cat("\nCohort by year and default (counts):\n")
292   → for stratified sampling checks
293 # print(
294 #   lc_cohort %>%
295 #     dplyr::count(issue_year, default) %>%
296   → and default
297 #     tidyr::pivot_wider(
298 #       names_from = default,
299 #       values_from = n,
300 #       values_fill = 0,
301 #       names_prefix = "default_"
302 #     )
303 # )
304 ## =====
305 ## 6. Stratified sample (by issue_year × default)
306 ## =====
307 target_n <- 120000
308   → for computational tractability
309 N_cohort <- nrow(lc_cohort)
310   → observations in the full modeling cohort
311 sample_frac <- target_n / N_cohort
312   → within each year by default stratum
313
314 # cat("\nTotal cohort size:", N_cohort, "\n")
315   → auditing
316 # cat("Target N:", target_n, "\n")
317   → size
318 # cat("Sampling fraction (per issue_year by default stratum):",
319   → sampling fraction used in stratified sampling
320 #   round(sample_frac, 4), "\n")
321
322 set.seed(base_seed + 1)
323   → reproducible stratified draw
324
325 lc_sample <- lc_cohort %>%
326   → preserving cohort composition
327   dplyr::group_by(issue_year, default) %>%
328     → and default outcome
329   dplyr::slice_sample(prop = sample_frac) %>%
330     → from each stratum
331   dplyr::ungroup()
332     → frame for downstream steps
333
334 # cat("\nActual lc_sample size:\n")
335   → size after stratification
336 # print(nrow(lc_sample))
337   → working sample
338 #
339 # cat("\nDefault distribution in lc_sample:\n")
340   → prevalence is close to the cohort rate
341 # print(table(lc_sample$default))
342 # print(prop.table(table(lc_sample$default)))
343 #
344 # cat("\nYear distribution in lc_sample:\n")
345   → composition is preserved
346 # print(table(lc_sample$issue_year))
347 # print(prop.table(table(lc_sample$issue_year)))

```

show default proportions
show year-by-default counts
count observations by year
target working sample size
total number of
sampling fraction to apply
print the cohort size for
print the intended sample
print the per-cell
fix the random seed for a
draw a working sample while
stratify by issuance year
sample the same fraction
return to an ungrouped data
print the realized sample
show number of rows in the
verify that default
show default counts
show default proportions
verify that year
show year counts
show year proportions

```

334
335 ## =====#
336 ## 7. From lc_sample + drop ID/leakage + concept-driven lc_model_base
337 ## =====#
338
339 ## -----
340 ## 7.1 Flag ID/text and leakage variables using full lc_raw
341 ## -----
342
343 all_vars <- names(lc_raw)                                # store all raw variable
344   ↪  names from the full Kaggle file
345
346 id_text_vars <- c(                                         # define obvious identifiers
347   ↪  and free-text fields to exclude from modeling
348   "id",
349   "member_id",
350   "url",
351   "desc",
352   "title",
353   "emp_title"
354 )
355
356 quasi_id_vars <- c("zip_code")                         # define quasi-identifiers to
357   ↪  exclude to reduce re-identification risk
358
359 leakage_patterns <- c(                                  # define name patterns for
360   ↪  variables that reflect post-origination information
361   "~total_pymnt",
362   "~total_pymnt_inv",
363   "~total_rec_",
364   "~recoveries",
365   "~collection_recovery_fee",
366   "~last_pymnt_",
367   "~next_pymnt_",
368   "~last_credit_pull_d",
369   "~settlement_",
370   "~hardship_",
371   "~debt_settlement_"
372 )
373
374 is_id_text <- all_vars %in% c(id_text_vars, quasi_id_vars)    # flag variables that are
375   ↪  identifiers, URLs, or free-text fields
376
377 is_leakage <- vapply(                                       # flag variables that match
378   ↪  any leakage pattern in their name
379   X = all_vars,
380   FUN = function(v) any(stringr::str_detect(v, leakage_patterns)),
381   FUN.VALUE = logical(1L)
382 )
383
384 var_flags <- tibble::tibble(                                # create a table summarizing
385   ↪  which raw variables are flagged for exclusion
386   variable = all_vars,
387   id_or_text = is_id_text,
388   leakage = is_leakage
389 )
390
391 # cat("\nFlag summary on full lc_raw (ID/text vs leakage):\n")
392   ↪  ID/text flags and leakage flags
393 # print(table(
394   ↪  id_or_text = var_flags$id_or_text,

```

```

387 # leakage      = var_flags$leakage
388 # ))
389 #
390 # cat("\nVariables flagged as ID, text, or quasi ID:\n")          # list all variables
391 #→ flagged as identifiers or free-text fields
391 # print(var_flags$variable[var_flags$id_or_text])
392 #
393 # cat("\nVariables flagged as potential leakage:\n")                # list all variables
394 #→ flagged as potential post-origination leakage
394 # print(var_flags$variable[var_flags$leakage])
395 #
396 ## -----
397 ## 7.2 Build lc_model_base from lc_sample (drop ID / leakage / labels)
398 ## -----
399 #
400 id_text_vars_full <- intersect(                                     # keep only ID and text
401   → variables that actually exist in the working sample
402   c(id_text_vars, quasi_id_vars),
402   names(lc_sample)
403 )
404 #
405 leakage_vars_pattern <- intersect(                                    # keep only pattern-flagged
406   → leakage variables that exist in the working sample
406   var_flags$variable[var_flags$leakage],
407   names(lc_sample)
408 )
409 #
410 extra_leakage_vars <- intersect(                                     # add known post-origination
411   → variables that may not match the patterns
411   c("out_prncp", "out_prncp_inv",
412     "last_fico_range_low", "last_fico_range_high"),
413   names(lc_sample)
414 )
415 #
416 outcome_label_vars <- intersect(                                     # exclude outcome labels and
417   → helper fields used during cleaning
417   c("loan_status", "default_prelim", "status_group"),
418   names(lc_sample)
419 )
420 #
421 drop_vars <- unique(c(                                              # combine all exclusion sets
422   → into a single list of variables to drop
422   id_text_vars_full,
423   leakage_vars_pattern,
424   extra_leakage_vars,
425   outcome_label_vars
426 ))
427 #
428 cat("\nNumber of variables in lc_sample:", ncol(lc_sample), "\n")    # report the number of
428 #→ columns available before dropping anything
429 cat("Number of variables to drop (ID/text/leakage/labels):", length(drop_vars), "\n") # report
429 #→ the number of excluded variables
430 #
431 cat("\nDropping these ID/text/quasi ID vars:\n")                      # print the ID and text
431 #→ variables being removed
432 print(id_text_vars_full)
433 #
434 cat("\nDropping these leakage vars (pattern-based plus explicit extras):\n") # print the leakage
434 #→ variables being removed
435 print(sort(unique(c(leakage_vars_pattern, extra_leakage_vars))))
436

```

```

437 cat("\nDropping these outcome/label/helper vars:\n")
438   ↪ being removed
439 print(outcome_label_vars)
440
440 lc_model_base <- lc_sample %>%
441   ↪ after removing identifiers and leakage fields
442   dplyr::select(-dplyr::all_of(drop_vars))
443   ↪ retain origination-time predictors only
444
443 # cat(
444   ↪ columns after exclusions as a quick check
445 #   "\nNumber of variables in lc_model_base",
446 #   "(incl. default, term_60, issue_year):",
447 #   ncol(lc_model_base), "\n"
448 # )
449
450 ## =====
451 ## 8. Concept-driven feature set: lc_model_orig
452 ## -----
453 ## I needed to further reduce the dimension for computation sake
454 ## =====
455
456 ## -----
457 ## 8.1 Feature shortlists by economic block
458 ## -----
459
460 feat_terms_pricing_num <- c(
461   ↪ capturing loan size and price of credit
462   "loan_amnt",
463   "funded_amnt",
464   "funded_amnt_inv",
465   "int_rate",
466   "installment"
467 )
468
468 feat_terms_pricing_cat <- c(
469   ↪ capturing platform risk buckets and listing status
470   "term",
471   ↪ the binary indicator built earlier
472   "grade",
473   "sub_grade",
474   "initial_list_status"
475 )
475
475 feat_income_emp_num <- c(
476   ↪ affordability measures
477   "annual_inc",
478   "dti"
479 )
479
480 feat_income_emp_cat <- c(
481   ↪ employment and verification information
482   "emp_length",
483   "home_ownership",
484   "verification_status"
485 )
485
486 feat_credit_num <- c(
487   ↪ revolving debt indicators
488   "fico_range_low",
489 )

```

```

488 "fico_range_high",
489 "open_acc",
490 "total_acc",
491 "revol_bal",
492 "revol_util",
493 "inq_last_6mths",
494 "delinq_2yrs",
495 "mths_since_last_delinq",
496 "mths_since_last_record",
497 "mths_since_recent_inq",
498 "pub_rec",
499 "pub_rec_bankruptcies",
500 "mort_acc",
501 "collections_12_mths_ex_med",
502 "num_rev_tl_bal_gt_0"
503 )
504
505 feat_credit_cat <- character(0) # no additional categorical
506   ↳ credit-history variables in this shortlist
507
508 feat_purpose_region_cat <- c( # categorical purpose and
509   ↳ application metadata
510   "purpose",
511   "application_type" # mostly individual and may
512   ↳ become near-constant after filtering
513 )
514
515 feat_numeric_cand <- c( # stack numeric candidates
516   ↳ from the economic blocks
517   feat_terms_pricing_num,
518   feat_income_emp_num,
519   feat_credit_num
520 )
521
522 feat_categorical_cand <- c( # stack categorical
523   ↳ candidates from the economic blocks
524   feat_terms_pricing_cat,
525   feat_income_emp_cat,
526   feat_credit_cat,
527   feat_purpose_region_cat
528 )
529
530 feat_all_cand <- c(feat_numeric_cand, feat_categorical_cand) # combine all shortlisted
531   ↳ variables into one candidate set
532
533 feat_present <- intersect(feat_all_cand, names(lc_model_base)) # keep only shortlisted
534   ↳ features that actually exist in the data
535
536 # cat("\nNumber of candidate features in shortlist (present in data):", # print the number of
537   ↳ shortlisted features available in the dataset
538 #   length(feat_present), "\n")
539 # cat("Candidate numeric vars present:\n") # print the numeric
540   ↳ shortlist that exists in the dataset
541 # print(intersect(feat_numeric_cand, feat_present))
542 # cat("\nCandidate categorical vars present:\n") # print the categorical
543   ↳ shortlist that exists in the dataset
544 # print(intersect(feat_categorical_cand, feat_present))
545
546 ## -----
547 ## 8.2 Create FICO average and refresh shortlists
548 ## -----

```

```

539 lc_model_features <- lc_model_base
540   ↪ before engineering additional features
541
542 if (all(c("fico_range_low", "fico_range_high") %in% names(lc_model_features))) { # proceed only
543   ↪ if both FICO bounds are available
544   lc_model_features <- lc_model_features %>%
545     ↪ adding a single FICO summary measure
546   dplyr::mutate(
547     fico_avg = 0.5 * (fico_range_low + fico_range_high)           # compute the midpoint of the
548     ↪ reported FICO band
549   )
550   cat("\nCreated fico_avg from fico_range_low and fico_range_high.\n")# log feature creation for
551   ↪ reproducibility
552
553   feat_numeric_cand <- setdiff(                                     # remove the raw FICO bounds
554     ↪ from the numeric candidate list
555     feat_numeric_cand,
556     c("fico_range_low", "fico_range_high")
557   )
558   feat_numeric_cand <- union(feat_numeric_cand, "fico_avg")          # add fico_avg to the numeric
559   ↪ candidate list
560 }
561
562 feat_numeric_present <- intersect(feat_numeric_cand, names(lc_model_features))      # keep only
563   ↪ numeric candidates present after the FICO update
564 feat_categorical_present <- intersect(feat_categorical_cand, names(lc_model_features))# keep only
565   ↪ categorical candidates present after the FICO update
566
567 # cat("\nNumeric candidates present after FICO adjustment:\n")                  # list numeric candidates
568   ↪ that remain after the FICO transformation
569 # print(feat_numeric_present)
570 # cat("\nCategorical candidates present:\n")                                    # list categorical
571   ↪ candidates that remain after the FICO transformation
572 # print(feat_categorical_present)
573
574 ## -----
575 ## 8.3 Tidy categorical variables (factor + lump rare purpose levels)
576 ## -----
577
578 purpose_min_count <- round(0.02 * nrow(lc_model_features))                      # set the minimum count for
579   ↪ keeping a named purpose level
580
581 lc_model_features <- lc_model_features %>%                                     # coerce shortlisted
582   ↪ categorical variables to factor for modeling
583   dplyr::mutate(
584     dplyr::across(
585       dplyr::all_of(feat_categorical_present),
586       ~ as.factor(.)
587     )
588   )
589
590 if ("purpose" %in% feat_categorical_present) { # lump rare purpose levels
591   ↪ only if the purpose variable is present
592   lc_model_features <- lc_model_features %>%
593     ↪ purpose categories into an other group
594   dplyr::mutate(
595     purpose = forcats::fct_lump_min(
596       purpose,
597       min        = purpose_min_count,
598       ↪ size threshold
599     )
600   )
601
602   # collapse low-frequency
603   # enforce the minimum cell

```

```

584     other_level = "other"                                # name for the pooled
585     ↪ rare-category level
586   )
587 }
588
589 ## -----
590 ## 8.4 Drop high-missing or near-constant features within shortlist
591 ## -----
592
593 vars_candidate <- intersect(                           # restrict predictors to
594   ↪ shortlisted variables that exist in the dataset
595   c(feat_numeric_present, feat_categorical_present),
596   names(lc_model_features)
597 )
598
599 tmp_candidate <- lc_model_features %>%                # create a temporary
600   ↪ candidate-only table for diagnostics
601   dplyr::select(dplyr::all_of(vars_candidate))
602
603 na_summary_short <- tmp_candidate %>%                  # compute missingness rates
604   ↪ for each candidate predictor
605   dplyr::summarise(dplyr::across(
606     dplyr::everything(),
607     ~ mean(is.na(.))
608   )) %>%
609   tidyr::pivot_longer(
610     dplyr::everything(),
611     names_to = "variable",
612     values_to = "prop_na"
613   )
614
615 high_na_vars <- na_summary_short %>%                  # identify predictors with
616   ↪ high missingness to exclude from modeling           # drop predictors with more
617   dplyr::filter(prop_na > 0.40) %>%
618     ↪ than 40 percent missing values
619   dplyr::pull(variable)
620
621 # cat("\nWithin shortlist: variables with more than 40 percent missingness:\n") # print variables
622 # dropped due to high missingness
623 # print(high_na_vars)
624
625 num_sd_short <- tmp_candidate %>%                      # compute standard deviations
626   ↪ for numeric predictors to detect near-constants
627   dplyr::select(where(is.numeric)) %>%
628   dplyr::summarise(dplyr::across(
629     dplyr::everything(),
630     sd,
631     na.rm = TRUE
632   )) %>%
633   tidyr::pivot_longer(
634     dplyr::everything(),
635     names_to = "variable",
636     values_to = "sd"
637   )
638
639 near_const_num <- num_sd_short %>%                    # identify numeric predictors
640   ↪ with negligible variation
641   dplyr::filter(sd < 1e-6) %>%
642     ↪ standard deviation as near-constant
643   dplyr::pull(variable)                                 # treat extremely small

```

```

635
636 # cat("\nWithin shortlist: near-constant numeric vars:\n")
637   ↪ dropped due to near-constant values
638 # print(near_const_num)
639
640 fac_maxprop_short <- tmp_candidate %>%
641   ↪ category share for each factor predictor
642   dplyr::select(where(is.factor)) %>%
643   dplyr::summarise(
644     dplyr::across(
645       dplyr::everything(),
646       ~ {
647         tab <- table(., useNA = "no")
648           ↪ levels excluding missing values
649         max(prop.table(tab))
650           ↪ most frequent level
651       }
652     )
653   ) %>%
654   tidyr::pivot_longer(
655     dplyr::everything(),
656     names_to = "variable",
657     values_to = "max_prop"
658   )
659
660 near_const_fac <- fac_maxprop_short %>%
661   ↪ dominated by a single category
662   dplyr::filter(max_prop > 0.99) %>%
663     ↪ than 99 percent mass in one level
664   dplyr::pull(variable)
665
666 # cat("\nWithin shortlist: near-constant factor vars:\n")
667   ↪ dropped due to near-constant levels
668 # print(near_const_fac)
669
670 vars_drop_short <- unique(c(
671   ↪ flagged for removal within the shortlist
672   high_na_vars,
673   near_const_num,
674   near_const_fac
675 ))
676
677 vars_predictors_final <- setdiff(vars_candidate, vars_drop_short)
678   ↪ set after removing problematic variables
679
680 # cat("\nFinal predictor set size:\n")
681   ↪ predictors kept for modeling
682 # print(length(vars_predictors_final))
683 # print(vars_predictors_final)
684
685 lc_model_orig <- lc_model_features %>%
686   ↪ modeling dataset with outcome and predictors
687   dplyr::select(
688     default,
689       ↪ outcome
690     term_60,
691       ↪ indicator for later analyses
692     issue_year,
693       ↪ stratification and controls
694     dplyr::all_of(vars_predictors_final)
695       ↪ origination-time predictors
696
697   # construct the final
698   # include the binary default
699   # include the 60-month term
700   # include issuance year for
701   # include the final set of

```



```

732 "Create fico_avg; convert to factors; lump rare purpose; handle high-missing/near-constant
733   ↪ vars",
734   "Keep: default, term_60, issue_year + final predictor shortlist"
735 )
736 ) %>%
737 dplyr::mutate(                                     # add derived columns that
738   ↪ quantify how much was removed at each step
739   rows_removed = dplyr::if_else(                  # compute rows removed
740     ↪ relative to the previous step
741     is.na(dplyr::lag(n_rows)), NA_integer_,
742     pmax(0L, dplyr::lag(n_rows) - n_rows)
743   ),
744   cols_dropped = dplyr::if_else(                 # compute columns dropped
745     ↪ relative to the previous step
746     is.na(dplyr::lag(n_cols)), NA_integer_,
747     pmax(0L, dplyr::lag(n_cols) - n_cols)
748   ),
749   cols_added = dplyr::if_else(                  # compute columns added
750     ↪ relative to the previous step
751     is.na(dplyr::lag(n_cols)), NA_integer_,
752     pmax(0L, n_cols - dplyr::lag(n_cols))
753   )
754 )
755
756 data_audit_keep <- data_audit %>%                # keep only the columns you
757   ↪ want to display in the audit table
758   dplyr::select(step, action, n_rows, rows_removed, n_cols, cols_dropped, cols_added, notes)
759
760 data_audit_pretty <- data_audit_keep %>%          # create a console-friendly
761   ↪ version with comma formatting
762   dplyr::mutate(
763     dplyr::across(
764       c(n_rows, rows_removed, n_cols, cols_dropped, cols_added),
765       ~ ifelse(is.na(.), NA, scales::comma(.))
766     )
767   )
768
769 print(data_audit_pretty)                           # print the audit table for a
770   ↪ quick visual check in the console
771
772 data_audit_tex <- data_audit_keep %>%            # convert the audit table
773   ↪ into a LaTeX-friendly data frame
774   dplyr::mutate(
775     dplyr::across(c(step, n_rows, rows_removed, n_cols, cols_dropped, cols_added), as.integer)
776   ) %>%
777   as.data.frame()
778
779 colnames(data_audit_tex) <- c(                  # set display column names
780   ↪ for the exported LaTeX table
781   "Step", "Action", "Rows", "Rows removed", "Cols", "Cols dropped", "Cols added", "Notes"
782 )
783
784 stargazer::stargazer(                            # export the audit trail as a
785   ↪ LaTeX table for the paper
786   data_audit_tex,
787   type = "latex",
788   summary = FALSE,
789   rownames = FALSE,
790   digits = 0,
791   no.space = TRUE,
792   table.placement = "!htbp",

```

```

782 title = "Data preparation audit trail",
783 label = "tab:data_audit",
784 out = "table_data_audit.tex"
785 )
786
787 ## -----
788 ## 9.2 Default rate by grade and term length
789 ## -----
790
791 plot_def_by_grade_term <- lc_model_orig %>%
792   dplyr::group_by(grade, term_60) %>%
793   dplyr::summarise(
794     def_rate = mean(default, na.rm = TRUE),
795     n         = dplyr::n(),
796     .groups   = "drop"
797   ) %>%
798   dplyr::mutate(
799     term_60 = factor(
800       levels = c(0, 1),
801       labels = c("36 months", "60 months")
802     )
803   ) %>%
804   ggplot(aes(x = grade, y = def_rate, group = term_60, colour = term_60)) + # plot default rate
805   geom_line(linewidth = 0.7) + # connect grade points to
806   geom_point(aes(size = n), alpha = 0.65) + # size points by number of
807   scale_y_continuous(labels = scales::percent_format(accuracy = 1)) + # express default rates as
808   scale_colour_manual(values = c("36 months" = "#0072B2", "60 months" = "#D55E00")) + # set
809   labs(
810     x      = "LendingClub grade",
811     y      = "Default rate",
812     colour = "Term",
813     size   = "Number of loans"
814   ) +
815   theme_academic() + # apply the project theme
816   theme(legend.position = "bottom") # place legend below the plot
817   # for readability
818
819 ## -----
820 ## 9.3 Interest rate versus default
821 ## -----
822
823 int_rate_thr_75 <- stats::quantile(lc_model_orig$int_rate, 0.75, na.rm = TRUE) # compute the
824   # seventy-fifth percentile interest-rate threshold
825
826 df_int_default <- lc_model_orig %>%
827   dplyr::mutate(int_rate_round = round(int_rate, 1)) %>%
828   dplyr::group_by(int_rate_round) %>%
829   # group by the interest-rate
830   # bin

```

```

828 dplyr::summarise(
829   def_rate = mean(default, na.rm = TRUE),
830   ↪   interest-rate bin
831   n        = dplyr::n(),
832   ↪   very small bins
833   .groups  = "drop"
834 ) %>%
835   dplyr::filter(n > 200)
836   ↪   observations to avoid noisy rates
837
838 y_max <- max(df_int_default$def_rate, na.rm = TRUE)
839   ↪   rate for placing the annotation label
840
841 plot_int_default <- df_int_default %>%
842   ↪   function of interest rate
843   ggplot(aes(x = int_rate_round, y = def_rate)) +
844     geom_line(color = "#0072B2") +
845     ↪   across interest-rate bins
846     geom_point(size = 1.2, color = "#0072B2") +
847     ↪   bin-level estimates
848     geom_vline(
849       ↪   percentile threshold used for the high-interest treatment
850       xintercept = int_rate_thr_75,
851       linetype   = 2,
852       color      = "#D55E00"
853     ) +
854     annotate(
855       ↪   directly on the plot for interpretation
856       "text",
857       x        = as.numeric(int_rate_thr_75),
858       y        = y_max,
859       label    = paste0("75th pct: ", round(int_rate_thr_75, 2), "%"),
860       hjust    = -0.05,
861       vjust    = 1.2,
862       size     = 3,
863       color    = "#D55E00"
864     ) +
865     scale_y_continuous(labels = scales::percent_format(accuracy = 1)) + # express default rates as
866     ↪   percentages
867     labs(
868       x = "Interest rate (%)",
869       y = "Default rate"
870     ) +
871     theme_academic()                                     # apply the project theme
872
873 plot_def_grade_int_panel <- plot_def_by_grade_term + plot_int_default +  # combine the grade-term
874   ↪   plot and the interest-default plot
875   patchwork::plot_layout(ncol = 2) +                      # place the two plots
876   ↪   side-by-side
877   patchwork::plot_annotation(                            # add panel tags for
878     ↪   referencing in the paper
879     tag_levels = "a",
880     tag_prefix = "(",
881     tag_suffix = ")"
882   )
883
884 print(plot_def_grade_int_panel)                         # display the combined panel
885   ↪   in the plotting device
886
887 ggplot2::ggsave(                                       # save the combined panel at
888   ↪   high resolution

```

```

874 filename = "plot_def_grade_int_panel.png",
875 plot      = plot_def_grade_int_panel,
876 width     = 12, height = 6, units = "in",
877 dpi       = 1000
878 )
879
880 ## -----
881 ## 9.4 Risk surface: grade by FICO by term (heatmap)
882 ## -----
883
884 if (!"fico_bin" %in% names(lc_model_orig)) {                                # create a FICO bin only if
885   → it does not already exist
886   lc_model_orig <- lc_model_orig %>%
887     dplyr::mutate(
888       fico_bin = cut(                                         # discretize fico_avg into
889         fico_avg,                                         standard bands for a readable heatmap
890         breaks = c(600, 640, 660, 680, 700, 720, 740, 760, 800),
891         include.lowest = TRUE,
892         right = FALSE
893     )
894   )
895 }
896
897 plot_heat_grade_fico <- lc_model_orig %>%                                # compute cell-level default
898   → rates and plot them as a heatmap
899   dplyr::filter(!is.na(fico_bin)) %>%                                # drop observations without
899   → a valid FICO bin
900   dplyr::group_by(grade, fico_bin, term_60) %>%                            # group by grade, FICO band,
900   → and term group
901   dplyr::summarise(                                                       # cell default rate
902     def_rate = mean(default, na.rm = TRUE),                                # cell size kept for
903     n        = dplyr::n(),                                                 diagnostics if needed
904     .groups = "drop"                                                       # build readable facet
905   ) %>%
906   dplyr::mutate(                                                       # map grade and FICO band to
907     term_60_lab = factor(                                             heatmap axes and default rate to fill
908       → labels for term groups
909       term_60,                                         levels = c(0, 1),
910       labels = c("36-month loans", "60-month loans")           # draw heatmap tiles with
911   ) %>%                                                               # print default rates inside
912   ggplot(aes(x = grade, y = fico_bin, fill = def_rate)) +                # show separate heatmaps for
913     geom_tile(color = "white") +                                         # map default rate to a
914     geom_text(                                                       # continuous color scale
915       aes(label = scales::percent(def_rate, accuracy = 1)),           name = "Default rate",
916       size = 3                                                       labels = scales::percent_format(accuracy = 1),
917     ) +                                                               low = "#FFF6DB",
918     facet_wrap(~ term_60_lab) +                                         high = "#7F1204"
919     scale_fill_gradient(                                             
920       →
921       name = "Default rate",
922       labels = scales::percent_format(accuracy = 1),
923       low = "#FFF6DB",
924       high = "#7F1204"
925

```

```

923 ) +
924   labs(
925     x = "Credit grade",
926     y = "FICO band (average of low and high)"
927   ) +
928   theme(
929     axis.text.x = element_text(angle = 45, hjust = 1),           # rotate grade labels for
930     ↵ readability
931     panel.grid = element_blank()                                # remove grid lines to keep
932     ↵ the heatmap clean
933 )
934
935 print(plot_heat_grade_fico)                                     # display the heatmap in the
936   ↵ plotting device
937
938 ggplot2::ggsave(                                              # save the heatmap to disk
939   ↵ at high resolution
940   filename = "plot_heat_grade_fico.png",
941   plot      = plot_heat_grade_fico,
942   width     = 12, height = 6, units = "in",
943   dpi       = 1000
944 )
945
946 ## =====
947 ## 10. ML dataset: split, impute, design matrices
948 ## =====
949
950 ml_data <- lc_model_orig %>%                                    # start from the
951   ↵ final modeling dataset
952   dplyr::mutate(
953     default    = as.integer(default),                            # ensure the
954     ↵ outcome is stored as 0 or 1 integer
955     term_60    = as.integer(term_60),                            # ensure the term
956     ↵ indicator is stored as 0 or 1 integer
957     issue_year = factor(issue_year),                           # treat issuance
958     ↵ year as a categorical control
959   ) %>%
960   dplyr::mutate(
961     dplyr::across(                                             # convert
962       ↵ selected categorical predictors to factors if present
963       .cols = intersect(
964         c("term", "grade", "sub_grade", "emp_length", "home_ownership",
965           "verification_status", "purpose", "application_type", "initial_list_status"),
966         names(.)
967       ),
968       .fns  = ~ as.factor(.x)                                   # coerce to
969       ↵ factor to enable stable one-hot encoding later
970     )
971   )
972
973 stopifnot(all(ml_data$default %in% c(0L, 1L)))                  # verify that the
974   ↵ outcome contains only 0 and 1
975
976 ## -----
977 ## 10.2 Stratified train/test split (70/30 within default)
978 ## -----

```

```

973 train_frac <- 0.70
974   ↳ observations assigned to the training set
974 set.seed(base_seed + 3)
975   ↳ the split is reproducible
975
976 ml_data <- ml_data %>%
977   ↳ identifier so sampled rows can be tracked
978     dplyr::mutate(row_id = dplyr::row_number())
978
979 train_ids <- ml_data %>%
980   ↳ rows within each outcome class
980     dplyr::group_by(default) %>%
981       ↳ default outcome to preserve class balance
981     dplyr::slice_sample(prop = train_frac) %>%
982       ↳ training fraction from each class
982     dplyr::ungroup() %>%
983     dplyr::pull(row_id)
983   ↳ selected training row identifiers
984
984 ml_data_split <- ml_data %>%
985   ↳ indicator and remove the helper id
986     dplyr::mutate(is_train = row_id %in% train_ids) %>%
987       ↳ each row is in the training split
987     dplyr::select(-row_id)
987   ↳ identifier after the split is defined
988
989 train_ml <- ml_data_split %>% dplyr::filter(is_train)
990   ↳ used for fitting and preprocessing
990 test_ml <- ml_data_split %>% dplyr::filter(!is_train)
990   ↳ held out for final evaluation
991
992 stopifnot(abs(mean(train_ml$default) - mean(test_ml$default)) < 0.01)
992   ↳ confirm that outcome prevalence is similar across splits
993
994 ## -----
995 ## 10.3 Median and Unknown imputation (fit on train, apply to test)
996 ## -----
997
998 num_vars <- train_ml %>%
999   ↳ numeric predictors using training data only
1000     dplyr::select(-default, -is_train) %>%
1001       ↳ and split flag from predictor lists
1000     dplyr::select(where(is.numeric)) %>%
1001       names()
1002
1003 fac_vars <- train_ml %>%
1004   ↳ predictors using training data only
1004     dplyr::select(-default, -is_train) %>%
1005       ↳ and split flag from predictor lists
1005     dplyr::select(where(is.factor)) %>%
1006       names()
1007
1008 num_impute_vals <- train_ml %>%
1009   ↳ training medians for numeric imputation
1009     dplyr::summarise(dplyr::across(dplyr::all_of(num_vars), ~ median(.x, na.rm = TRUE))) # store one
1010       ↳ median per numeric variable
1011
1011 train_ml_imp <- train_ml %>%
1012   ↳ imputation rules to the training set
1012     dplyr::mutate(

```

```

1013 dplyr::across(dplyr::all_of(fac_vars),                                # replace missing
1014   ~ forcats::fct_na_value_to_level(.x, level = "Unknown")),
1015 dplyr::across(dplyr::all_of(num_vars),                                     # replace missing
1016   ~ numeric_values_using_training_medians
1017   ~ replace(.x, is.na(.x), num_impute_vals[[dplyr::cur_column()]]))
1018 )
1019 test_ml_imp <- test_ml %>%
1020   → imputation rules to the test set
1021 dplyr::mutate(
1022   dplyr::across(dplyr::all_of(fac_vars),                                # replace missing
1023     ~ factor values with the same explicit level
1024     ~ forcats::fct_na_value_to_level(.x, level = "Unknown")),
1025   dplyr::across(dplyr::all_of(num_vars),                                     # replace missing
1026     ~ numeric_values_using_training_medians
1027     ~ replace(.x, is.na(.x), num_impute_vals[[dplyr::cur_column()]]))
1028 )
1029 for (v in fac_vars) {
1030   → factor levels with training levels for stable one-hot encoding
1031   lvls_train <- levels(train_ml_imp[[v]])                                # align test
1032   → training levels for this factor
1033   x_chr      <- as.character(test_ml_imp[[v]])                            # store the
1034   → values to character for safe recoding
1035   x_chr[is.na(x_chr) | !(x_chr %in% lvls_train)] <- "Unknown"          # convert test
1036   → unseen test levels to Unknown
1037   test_ml_imp[[v]] <- factor(x_chr, levels = lvls_train)                 # map missing and
1038   → factor using the training level set
1039 }
1040 stopifnot(!anyNA(train_ml_imp %>% dplyr::select(-default, -is_train)))    # rebuild the
1041 stopifnot(!anyNA(test_ml_imp  %>% dplyr::select(-default, -is_train)))    # confirm there
1042 ## -----
1043 ## 10.4 Design matrices (glmnet-ready one-hot encoding)
1044 ## -----
1045 mm_form <- default ~ . - is_train                                         # specify the
1046   → model formula and drop the split flag
1047 mf_train <- stats::model.frame(mm_form, data = train_ml_imp)               # build the
1048   → training model frame and store factor level information
1049 x_train <- stats::model.matrix(mm_form, data = mf_train)                   # create the
1050   → training design matrix with one-hot encoded factors
1051 mf_test  <- stats::model.frame(                                           # build the test
1052   → model frame using the training factor levels
1053   mm_form,
1054   data = test_ml_imp,
1055   xlev = .getXlevels(stats::terms(mf_train), mf_train)
1056 )
1057 x_test   <- stats::model.matrix(mm_form, data = mf_test)                  # create the
1058   → test design matrix with the same columns as training
1059 y_train <- as.numeric(train_ml_imp$default)                                 # extract
1060   → training labels as numeric 0 or 1
1061 y_test  <- as.numeric(test_ml_imp$default)                                  # extract test
1062   → labels as numeric 0 or 1

```

```

1055
1056 stopifnot(ncol(x_train) == ncol(x_test))                                # verify that
   ↳ train and test have identical feature columns
1057
1058 ## -----
1059 ## 10.5 Metrics helper
1060 ## -----
1061
1062 if (!exists("metrics_binary")) {                                         # define the
   ↳ helper only once to avoid overwriting in later runs
     if (!requireNamespace("pROC", quietly = TRUE)) {                         # check that the
       ↳ ROC and AUC package is available
         install.packages("pROC")                                              # install the
       ↳ package if it is missing
     }
1063 library(pROC)                                                       # attach pROC
   ↳ for ROC and AUC computation
1064
1065 metrics_binary <- function(y_true, p_hat, threshold = 0.5) {           # compute
   ↳ standard classification metrics from labels and probabilities
     stopifnot(length(y_true) == length(p_hat))                               # confirm that
       ↳ labels and predictions have the same length
     y_true <- as.integer(y_true)                                             # store labels
       ↳ as 0 or 1 integers
     p_hat <- as.numeric(p_hat)                                              # store
       ↳ predicted probabilities as numeric values
1066
1067 p_hat <- pmin(pmax(p_hat, 1e-6), 1 - 1e-6)                           # clamp
   ↳ probabilities away from zero and one for stability
1068
1069 pred_class <- ifelse(p_hat >= threshold, 1L, 0L)                      # convert
   ↳ probabilities to class predictions using the threshold
1070
1071 accuracy <- mean(pred_class == y_true)                                    # compute
   ↳ accuracy as the share of correctly classified observations
1072 error      <- 1 - accuracy                                              # compute
   ↳ misclassification error as one minus accuracy
1073
1074 tp <- sum(pred_class == 1L & y_true == 1L)                             # count true
   ↳ positives at the chosen threshold
1075 fp <- sum(pred_class == 1L & y_true == 0L)                             # count false
   ↳ positives at the chosen threshold
1076 tn <- sum(pred_class == 0L & y_true == 0L)                            # count true
   ↳ negatives at the chosen threshold
1077 fn <- sum(pred_class == 0L & y_true == 1L)                            # count false
   ↳ negatives at the chosen threshold
1078
1079 roc_obj <- pROC::roc(response = y_true, predictor = p_hat, quiet = TRUE) # compute the
   ↳ ROC curve object from labels and probabilities
1080 auc_val <- as.numeric(pROC::auc(roc_obj))                                # extract the
   ↳ scalar AUC value
1081
1082 brier <- mean((p_hat - y_true)^2)                                         # compute the
   ↳ Brier score as mean squared probability error
1083
1084 dplyr::tibble(                                                               # return a
   ↳ one-row tibble for easy table binding
     error      = error,
     accuracy   = accuracy,
     auc        = auc_val,
     brier      = brier,
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094

```

```

1095     tp = tp, fp = fp, tn = tn, fn = fn
1096   )
1097 }
1098 }
1099
1100 ## =====
1101 ## 11. Regularised regression: logit & elastic net (block (a))
1102 ## =====
1103
1104 ## -----
1105 ## 11.1 Unpenalised logistic regression (GLM)
1106 ## -----
1107
1108 logit_formula <- default ~ . - is_train                                # define the model
1109   # formula and exclude the split indicator
1110
1111 logit_model <- glm(                                                 # fit a baseline
1112   formula = logit_formula,                                              # use the shared
1113   family = binomial(link = "logit"),                                         # specify a Bernoulli
1114   data = train_ml_imp,                                                    # use the imputed
1115   # training dataset
1116 )
1117
1118 stopifnot(isTRUE(logit_model$converged))                                # verify that the
1119   # optimization converged successfully
1120
1121 logit_p_test <- predict(                                                 # generate predicted
1122   # default probabilities on the test set
1123   logit_model,                                                       # fitted logistic
1124   newdata = test_ml_imp,                                              # imputed test
1125   # dataset
1126   type = "response"                                                 # return
1127   # probabilities rather than linear predictors
1128 )
1129
1130
1131 logit_metrics <- metrics_binary(y_test, logit_p_test) %>%                # compute test-set
1132   # metrics for the logistic regression
1133   dplyr::mutate(
1134     model = "logit",                                                 # label used when
1135     # combining results across models
1136     alpha = NA_real_,                                              # keep placeholder
1137     # fields for consistent output structure
1138     lambda = NA_real_,                                             # keep placeholder
1139     # fields for consistent output structure
1140   )
1141
1142 ## -----
1143 ## 11.2 Elastic-net logistic regression (glmnet; tune alpha and lambda by CV)
1144 ## -----
1145
1146 alpha_grid <- c(0, 0.25, 0.5, 0.75, 1)                                    # define the alpha
1147   # grid from ridge to lasso
1148
1149 set.seed(base_seed + 4)                                                 # set a seed so
1150   # cross-validation folds are reproducible

```

```

1139 K <- 5                                     # set the number of
1140   ↪ cross-validation folds
1141 foldid <- sample(rep(seq_len(K), length.out = nrow(x_train)))
1142   ↪ training row to a fold using a fixed fold id vector
1143 glmnet_results_list <- vector("list", length(alpha_grid))      # preallocate a list
1144   ↪ to store performance metrics by alpha
1145 glmnet_fits       <- vector("list", length(alpha_grid))      # preallocate a list
1146   ↪ to store fitted cv.glmnet objects by alpha
1147
1148 for (i in seq_along(alpha_grid)) {           # iterate over all
1149   ↪ alpha values in the grid
1150     a <- alpha_grid[i]                         # set the current
1151     ↪ alpha value
1152     cat("\n[glmnet] Fitting cv.glmnet with alpha =", a, "...\\n")
1153   ↪ long runs
1154
1155     cv_fit <- glmnet::cv.glmnet(               # fit an elastic-net
1156       ↪ logistic regression with cross-validation
1157       x             = x_train,                  # training design
1158       ↪ matrix
1159       y             = y_train,                  # training labels
1160       ↪ family        = "binomial",            # logistic regression
1161       ↪ loss
1162       alpha         = a,                      # elastic-net mixing
1163       ↪ parameter
1164       nfolds        = K,                      # number of folds
1165       ↪ used in cross-validation
1166       foldid        = foldid,                 # reuse identical
1167       ↪ folds across alpha values
1168       standardize  = TRUE,                   # standardize
1169       ↪ predictors within glmnet
1170       type.measure  = "deviance",            # select lambda by
1171       ↪ cross-validated deviance
1172     )
1173
1174     idx_min    <- which(cv_fit$lambda == cv_fit$lambda.min)[1]      # locate the index of
1175       ↪ the selected lambda within the lambda path
1176     cv_err_min <- cv_fit$cvm[idx_min]                                # store the
1177       ↪ cross-validated deviance at the selected lambda
1178
1179     p_hat_test <- predict(                                         # predict test-set
1180       ↪ default probabilities at lambda.min
1181       cv_fit,                                              # fitted
1182       ↪ cross-validation object
1183       newx = x_test,                                         # test design matrix
1184       s     = "lambda.min",                                    # choose the lambda
1185       ↪ that minimizes CV deviance
1186       type = "response",                                     # return
1187       ↪ probabilities
1188     )[, 1]
1189
1190     glmnet_results_list[[i]] <- metrics_binary(y_test, p_hat_test) %>%
1191       ↪ on the common test set
1192     dplyr::mutate(
1193       model    = "glmnet",                                     # model identifier
1194       ↪ used for downstream tables and plots
1195       alpha    = a,                                         # alpha value used in
1196       ↪ this fit
1197       lambda   = cv_fit$lambda.min,                          # selected lambda
1198       ↪ value

```

```

1175     cv_error = cv_err_min
1176     ↪ deviance used for selecting among alpha values
1177   )
1178
1179   glmnet_fits[[i]] <- cv_fit
1180   ↪ object so the best model can be retrieved later
1181 }
1182
1183 glmnet_metrics_all <- dplyr::bind_rows(glmnet_results_list)
1184   ↪ across all alpha values into one table
1185
1186 # cat("\n[glmnet] Test-set performance (all alphas):\n")
1187   ↪ alpha grid results for manual inspection
1188 # print(glmnet_metrics_all)
1189
1190 glmnet_best_metrics <- glmnet_metrics_all %>%
1191   ↪ elastic-net model across alpha values
1192   dplyr::arrange(cv_error, dplyr::desc(auc)) %>%
1193   ↪ CV deviance and break ties by higher test AUC
1194   dplyr::slice(1)
1195
1196 alpha_best <- glmnet_best_metrics$alpha
1197   ↪ alpha value
1198 lambda_best <- glmnet_best_metrics$lambda
1199   ↪ corresponding selected lambda value
1200
1201 best_idx <- which(alpha_grid == alpha_best)[1]
1202   ↪ of the best alpha in the stored fit list
1203 glmnet_best_fit <- glmnet_fits[[best_idx]]
1204   ↪ cv.glmnet object for the best alpha
1205
1206 glmnet_best_p_test <- predict(
1207   ↪ test-set probabilities for the selected elastic-net fit
1208   glmnet_best_fit,
1209   ↪ fit object
1210   newx = x_test,
1211   s    = lambda_best,
1212   ↪ value
1213   type = "response"
1214   ↪ probabilities
1215 )[, 1]
1216
1217 ## -----
1218 ## 11.3 Regression-family summary table
1219 ## -----
1220
1221 regression_family_metrics <- dplyr::bind_rows(
1222   ↪ regression and elastic-net results
1223   logit_metrics,
1224   glmnet_best_metrics
1225 ) %>%
1226   dplyr::select(model, alpha, lambda, error, accuracy, auc, brier,
1227   ↪ columns needed for the summary table
1228     tp, fp, tn, fn, cv_error)
1229
1230 cat("\nRegression-family models: summary metrics:\n")
1231   ↪ the regression-family metrics table
1232 print(regression_family_metrics)
1233   ↪ regression-family metrics table to the console
1234
1235 ## =====

```

```

1218 ## 12. Tree-based models: CART, Bagging (custom + benchmark), Random Forest
1219 ## =====
1220
1221 ## -----
1222 ## 12.0 Data objects (same predictors; no new split)
1223 ## -----
1224
1225 tree_formula <- default ~ . - is_train                                # use the same
1226   ↳ predictor set as in the regression models
1227
1228 train_tree <- train_ml_imp                                              # training data after
1229   ↳ imputation
1230 test_tree  <- test_ml_imp                                                 # test data after
1231   ↳ imputation
1232
1233 train_tree$default <- factor(train_tree$default, levels = c(0, 1))      # convert outcome to
1234   ↳ factor for classification trees
1235 test_tree$default  <- factor(test_tree$default,  levels = c(0, 1))       # keep test levels
1236   ↳ identical to the training levels
1237
1238 y_train_vec <- as.integer(as.character(train_tree$default))             # store training
1239   ↳ labels as integer zero or one for metric functions
1240 y_test_vec  <- as.integer(as.character(test_tree$default))                # store test labels
1241   ↳ as integer zero or one for metric functions
1242
1243 get_prob1 <- function(prob_obj) {                                         # extract the
1244   ↳ probability for class one from different predict outputs
1245     if (is.null(dim(prob_obj))) return(as.numeric(prob_obj))            # return directly when
1246       ↳ predict already gives a vector
1247     cn <- colnames(prob_obj)                                             # read the column
1248       ↳ names that represent class labels
1249     if (!is.null(cn) && "1" %in% cn) return(prob_obj[, "1"])           # take the column
1250       ↳ labeled one when it exists
1251     if (!is.null(cn) && "0" %in% cn) return(rep(0, nrow(prob_obj)))      # return zeros when a
1252       ↳ stump predicts only class zero
1253     prob_obj[, ncol(prob_obj)]                                           # fall back to the
1254       ↳ last column if labels are missing
1255 }
1256
1257 ## -----
1258 ## 12.1 CART tree (rpart) with CV cp and stump safeguard
1259 ## -----
1260
1261 cat("\n[CART] Fitting CART tree...\n")                                     # print a progress
1262   ↳ message
1263
1264 set.seed(base_seed + 5)                                                       # set a seed so the
1265   ↳ CART fit is reproducible
1266
1267 cart_grow <- rpart(                                                       # first grow a
1268   ↳ reasonably flexible tree
1269   formula = tree_formula,                                                 # use the shared tree
1270     ↳ formula
1271   data    = train_tree,                                                    # fit on the training
1272     ↳ data
1273   method  = "class",                                                       # fit a classification
1274     ↳ tree
1275   control = rpart.control(                                                 # allow growth so
1276     ↳ cp      = 0.0005,
1277       ↳ pruning can be meaningful

```

```

1258     minbucket = 200,                                     # enforce a minimum
1259     ↪   leaf size for stability
1260     maxdepth  = 10                                      # cap depth to avoid
1261     ↪   overly complex trees
1262   )
1263 
1263 ctab   <- cart_grow$cptable                                # extract the cross
1264   ↪   validation table produced by rpart
1264 cp_min <- ctab[which.min(ctab[, "xerror"]), "CP"]          # choose the
1264   ↪   complexity parameter with minimum cross validated error
1265 
1266 cart_final <- prune(cart_grow, cp = cp_min)                # prune the grown
1266   ↪   tree using the selected complexity parameter
1267 
1268 if (!is.null(cart_final$cptable) && cart_final$cptable[1, "nsplit"] == 0) { # detect the case
1268   ↪   where pruning removes all splits
1269   cat("[CART] Warning: pruned tree became a stump. Falling back to grown tree.\n") # report the
1269   ↪   fallback choice
1270   tree_fit    <- cart_grow                                    # use the grown tree
1270   ↪   as a fallback model
1271   cart_cp_used <- cp_min                                     # keep the selected
1271   ↪   complexity parameter for reporting
1272 } else {
1273   tree_fit    <- cart_final                                 # use the pruned tree
1273   ↪   when it still contains splits
1274   cart_cp_used <- cp_min                                     # store the
1274   ↪   complexity parameter used for pruning
1275 }
1276 
1277 tree_prob_test <- get_probi(predict(tree_fit, newdata = test_tree, type = "prob")) # compute test
1277   ↪   probabilities and extract the class one column
1278 
1279 tree_metrics <- metrics_binary(y_test_vec, tree_prob_test) %>%                      # evaluate CART
1279   ↪   performance on the common test set
1280   dplyr::mutate(
1281     model  = "cart_tree",                                     # model identifier
1281     ↪   used in later tables and plots
1282     alpha   = NA_real_,                                     # placeholder to
1282     ↪   match the common metric table structure
1283     lambda = cart_cp_used,                                # store the chosen
1283     ↪   complexity parameter for reference
1284   )
1285 
1286 cat("\n[CART] Test-set performance:\n")                         # print a header for
1286   ↪   the CART metrics
1287 print(tree_metrics)                                            # print CART metrics
1288 
1289 ## -----
1290 ## 12.2 Custom bagging with rpart and OOB tracking
1291 ## -----
1292 
1293 bagging_rpart_oob_trace <- function(train_df, test_df, formula, y_true,
1294   B = 80,
1295   minbucket = 50,
1296   maxdepth  = 20,
1297   cp        = 0.0005) {
1298 
1299   n_train <- nrow(train_df)                                  # number of training
1299   ↪   observations

```

```

1300 n_test <- nrow(test_df)                                # number of test
1301   ↳ observations
1302 test_sum <- numeric(n_test)                            # running sum of
1303   ↳ predicted probabilities on the test set
1304 oob_sum <- numeric(n_train)                            # running sum of out
1305   ↳ of bag probabilities on the training set
1306 oob_count <- integer(n_train)                          # number of out of
1307   ↳ bag predictions per training observation
1308
1309 err_path <- numeric(B)                                # store out of bag
1310   ↳ misclassification error after each tree
1311
1312 for (b in seq_len(B)) {                                 # loop over the
1313   ↳ number of bootstrap trees
1314     if (b %% 10 == 0) cat("[Bagging] Tree", b, "of", B, "...\\n")    # print progress
1315   ↳ every ten trees
1316
1317 idx_b <- sample.int(n = n_train, size = n_train, replace = TRUE)      # draw a bootstrap
1318   ↳ sample of training indices
1319
1320 inbag_flag <- logical(n_train)                         # initialize in-bag
1321   ↳ indicator for this tree
1322 inbag_flag[idx_b] <- TRUE                             # mark in-bag rows
1323
1324 oob <- which(!inbag_flag)                            # identify out of bag
1325   ↳ rows for this tree
1326
1327 train_b <- train_df[idx_b, , drop = FALSE]           # create the
1328   ↳ bootstrap training set for this tree
1329
1330 fit_b <- rpart(                                       # fit a base CART
1331   ↳ learner on the bootstrap sample
1332     formula = formula,
1333     data = train_b,
1334     method = "class",
1335     control = rpart.control()                           # skip internal cross
1336   ↳ validation to speed up bagging
1337     cp = cp, minbucket = minbucket, maxdepth = maxdepth, xval = 0
1338   )
1339
1340 p_test <- get_prob1(predict(fit_b, newdata = test_df, type = "prob"))  # compute predicted
1341   ↳ probabilities on the test set
1342 test_sum <- test_sum + p_test                            # accumulate test
1343   ↳ probabilities for bagged averaging
1344
1345 if (length(oob) > 0) {                                  # update out of bag
1346   ↳ predictions when out of bag rows exist
1347     p_oob <- get_prob1(
1348       predict(fit_b, newdata = train_df[oob, , drop = FALSE], type = "prob")
1349     )
1350     oob_sum[oob] <- oob_sum[oob] + p_oob                # accumulate out of
1351     ↳ bag probability sums for those rows
1352     oob_count[oob] <- oob_count[oob] + 1L                 # increment the
1353     ↳ number of out of bag votes for those rows
1354   }
1355
1356 has_oob <- oob_count > 0L                            # identify training
1357   ↳ rows with at least one out of bag prediction so far
1358 if (any(has_oob)) {

```

```

1341     curr_probs <- oob_sum[has_oob] / oob_count[has_oob]           # compute current
1342     ↪  average out of bag probabilities
1343     curr_preds <- ifelse(curr_probs > 0.5, 1, 0)                 # convert
1344     ↪  probabilities to classes using threshold one half
1345     err_path[b] <- mean(curr_preds != y_true[has_oob])            # compute out of bag
1346     ↪  misclassification error at this step
1347   } else {
1348     err_path[b] <- NA_real_                                         # keep missing when
1349     ↪  no out of bag predictions exist yet
1350   }
1351 }
1352
1353 test_prob <- test_sum / B                                         # compute the final
1354   ↪  bagged test probabilities
1355
1356 oob_prob <- rep(NA_real_, n_train)                                    # initialize the out
1357   ↪  of bag probability vector
1358 has_oob <- oob_count > 0L                                         # identify
1359   ↪  observations with at least one out of bag prediction
1360 oob_prob[has_oob] <- oob_sum[has_oob] / oob_count[has_oob]          # compute final out
1361   ↪  of bag probabilities by averaging votes
1362
1363 oob_idx <- which(has_oob)                                           # store indices of
1364   ↪  observations with out of bag probabilities
1365
1366 list(test_prob = test_prob, oob_prob = oob_prob, oob_idx = oob_idx, err_path = err_path) # return test and out of bag outputs for reporting
1367 }
1368
1369 set.seed(base_seed + 6)                                              # set a seed so
1370   ↪  bagging is reproducible
1371
1372 cat("\n[Bagging] Fitting custom rpart bagging with OOB...\n")      # print a progress
1373   ↪  message
1374
1375 bag_res <- bagging_rpart_oob_trace(                                # training data used
1376   train_df = train_tree,                                               # training data used
1377   ↪  to fit each bootstrap tree
1378   test_df = test_tree,                                                 # test data used to
1379   ↪  compute bagged test predictions
1380   formula = tree_formula,                                             # formula shared
1381   ↪  across tree models
1382   y_true = y_train_vec,                                                # training truth used
1383   ↪  for out of bag error tracking
1384   B = 80,                                                               # number of bootstrap
1385   ↪  trees
1386   minbucket = 50,                                                       # minimum leaf size
1387   ↪  for each base tree
1388   maxdepth = 20,                                                       # maximum depth for
1389   ↪  each base tree
1390   cp = 0.0005,                                                          # complexity parameter
1391   ↪  controlling base tree growth
1392 )
1393
1394 bagging_prob_test <- bag_res$test_prob                               # store bagged test
1395   ↪  probabilities
1396
1397 bagging_test_metrics <- metrics_binary(y_test_vec, bagging_prob_test) %>% # compute test-set
1398   ↪  metrics for custom bagging
1399   dplyr::mutate(model = "bagging_rpart_test", alpha = NA_real_, lambda = NA_real_)
1400

```

```

1380 cat("\n[Bagging custom] Test-set performance:\n")                                # print a header for
1381   ↪ custom bagging test metrics
1381 print(bagging_test_metrics)
1382   ↪ test metrics
1382
1383 y_oob <- y_train_vec[bag_res$oob_idx]                                              # store out of bag
1384   ↪ truth labels for observations with OOB predictions
1384
1385 bagging_oob_metrics <- metrics_binary(                                               # compute metrics on
1386   ↪ the out of bag predictions only
1386     y_true = y_oob,
1387     p_hat  = bag_res$oob_prob[bag_res$oob_idx]
1388   ) %>%
1389   dplyr::mutate(model = "bagging_rpart_oob", alpha = NA_real_, lambda = NA_real_)
1390
1391 cat("\n[Bagging custom] OOB performance:\n")                                         # print a header for
1392   ↪ custom bagging out of bag metrics
1392 print(bagging_oob_metrics)                                                       # print custom bagging
1393   ↪ out of bag metrics
1393
1394 ## -----
1395 ## 12.3 Benchmark bagging (ipred::bagging) and its OOB error
1396 ## -----
1397
1398 set.seed(base_seed + 7)                                                               # set a seed so
1399   ↪ benchmark bagging is reproducible
1400
1400 cat("\n[Bagging] Fitting ipred::bagging (benchmark)... \n")                         # print a progress
1401   ↪ message
1401
1402 bag_std <- ipred::bagging(                                                        # formula shared with
1403   formula = tree_formula,                                                 # training data used
1404   ↪ other tree models
1404   data   = train_tree,                                                 # number of bootstrap
1405   ↪ for fitting
1405   nbagg  = 80,                                                       # request out of bag
1405   ↪ trees
1406   coob   = TRUE,                                                       # speed up by
1406   ↪ error reporting
1407   control = rpart.control(xval = 0, cp = 0.0, minbucket = 50)           # skipping internal CV inside each tree
1407   )
1408
1409 if (!is.null(bag_std$err)) {                                                       # ipred stores an out
1410   ↪ of bag error path in the err slot
1411   oob_err_ipred <- tail(bag_std$err, 1)                                     # store the final out
1411   ↪ of bag misclassification error
1412   cat("\n[ipred::bagging] Final OOB misclassification error:\n")          # print a header for
1412   ↪ the final out of bag error
1413   print(oob_err_ipred)                                                       # print the final out
1413   ↪ of bag error
1414 }
1415
1416 bag_std_prob_test <- get_prob1(predict(bag_std, newdata = test_tree, type = "prob")) # compute
1416   ↪ benchmark bagging probabilities on the test set
1417
1418 bagging_ipred_test_metrics <- metrics_binary(y_test_vec, bag_std_prob_test) %>%      # compute test
1418   ↪ metrics for benchmark bagging
1419   dplyr::mutate(model = "bagging_ipred_test", alpha = NA_real_, lambda = NA_real_)
```

```

1421 cat("\n[ipred::bagging] Test-set performance:\n")                                # print a header for
1422   ↪ benchmark bagging test metrics
1423 print(bagging_ipred_test_metrics)                                                 # print benchmark
1424   ↪ bagging test metrics
1425
1426 ## -----
1427 ## 12.4 Random forest (ranger)
1428 ## -----
1429
1430 set.seed(base_seed + 8)                                                       # set a seed so the
1431   ↪ random forest fit is reproducible
1432
1433 cat("\n[RF] Fitting random forest (ranger)...\n")                            # print a progress
1434   ↪ message
1435
1436 rf_best_fit <- ranger(                                                   
1437   formula          = tree_formula,                                         # formula shared
1438   ↪ across models
1439   data             = train_tree,                                           # training data used
1440   ↪ for fitting
1441   num.trees        = 300,                                                 # number of trees in
1442   ↪ the forest
1443   mtry            = floor(sqrt(ncol(train_tree) - 1)),                 # default heuristic
1444   ↪ for candidate variables per split
1445   min.node.size   = 100,                                                 # minimum terminal
1446   ↪ node size for regularization
1447   probability     = TRUE,                                                # request class
1448   ↪ probability predictions
1449   importance      = "impurity",                                         # compute
1450   ↪ impurity-based variable importance
1451 )
1452
1453 rf_pred_test <- predict(rf_best_fit, data = test_tree)$predictions[, "1"]    # extract predicted
1454   ↪ probabilities for class one on the test set
1455
1456 rf_best_metrics <- metrics_binary(y_test_vec, rf_pred_test) %>%           # compute test metrics
1457   ↪ for the random forest
1458   dplyr::mutate(model = "random_forest_best", alpha = NA_real_, lambda = NA_real_)
1459
1460 cat("\n[RF] Test-set performance:\n")                                         # print a header for
1461   ↪ random forest test metrics
1462 print(rf_best_metrics)                                                       # print random forest
1463   ↪ test metrics
1464
1465 ## -----
1466 ## 12.5 Collect tree-family test-set metrics (same test set)
1467 ## -----
1468
1469 tree_family_metrics <- dplyr::bind_rows(                                     # combine test-set
1470   ↪ metrics across all tree-based models
1471   tree_metrics,
1472   bagging_test_metrics,
1473   bagging_ipred_test_metrics,
1474   rf_best_metrics
1475 )
1476
1477 cat("\nTree-based models: test-set performance summary (same test set):\n") # print a header for
1478   ↪ the combined tree metrics
1479 print(tree_family_metrics)                                                   # print the combined
1480   ↪ tree metrics

```

```

1464 cat("\nCustom bagging OOB metric (for discussion):\n") # print a header for
1465   ↪ the out of bag metrics
1466 print(bagging_oob_metrics) # print the out of bag
1467   ↪ metrics for custom bagging
1468
1469 ## =====
1470 ## 12.6 OOB reporting outputs (tables plus optional OOB curve plot)
1471 ## =====
1472
1473 ## -----
1474 ## 12.6.1 OOB summary table (final + best + extra metrics where possible)
1475 ## -----
1476
1477 stopifnot(exists("bag_res"), exists("bag_std"), exists("rf_best_fit"),
1478   ↪ exists("bagging_oob_metrics")) # ensure required objects exist before reporting
1479
1480 custom_err_path <- bag_res$err_path # extract the custom
1481   ↪ bagging out of bag error path
1482 custom_err_path <- custom_err_path[!is.na(custom_err_path)] # remove missing
1483   ↪ entries from early iterations if present
1484
1485 summarise_oob_path <- function(err_path) { # summarise an out of
1486   ↪ bag error path into a few reporting statistics
1487     tibble::tibble(
1488       n_trees_final = length(err_path), # total number of
1489         ↪ trees reported in the error path
1490       oob_error_final = as.numeric(tail(err_path, 1)), # final out of bag
1491         ↪ misclassification error
1492       oob_error_min = as.numeric(min(err_path)), # minimum out of bag
1493         ↪ misclassification error achieved
1494       n_trees_at_min = as.integer(which.min(err_path)) # tree index where the
1495         ↪ minimum out of bag error occurs
1496     )
1497   }
1498
1499 oob_custom_sum <- summarise_oob_path(custom_err_path) # compute the custom
1500   ↪ bagging out of bag summary statistics
1501
1502 oob_ipred_sum <- tibble::tibble( # create a comparable
1503   ↪ summary row for ipred bagging
1504     n_trees_final = 80L, # number of trees used
1505       ↪ in ipred bagging
1506     oob_error_final = as.numeric(oob_err_ipred), # final out of bag
1507       ↪ error from ipred
1508     oob_error_min = NA_real_, # not available
1509       ↪ because a full path is not used here
1510     n_trees_at_min = NA_integer_ # not available
1511       ↪ because a full path is not used here
1512   )
1513
1514 oob_rf_sum <- tibble::tibble( # create a comparable
1515   ↪ summary row for random forest
1516     n_trees_final = as.integer(rf_best_fit$num.trees), # number of trees used
1517       ↪ in the forest
1518     oob_error_final = as.numeric(rf_best_fit$prediction.error), # built-in out of bag
1519       ↪ misclassification error from ranger
1520     oob_error_min = NA_real_, # not available
1521       ↪ because a full path is not stored by default
1522     n_trees_at_min = NA_integer_ # not available
1523       ↪ because a full path is not stored by default
1524   )

```

```

1504
1505 oob_report_tbl <- dplyr::bind_rows(
1506   ↪ summaries into one reporting table
1507   dplyr::mutate(oob_custom_sum, method = "Bagging custom",           extra_note = "Full OOB path and
1508   ↪ OOB probabilities"),
1509   dplyr::mutate(oob_ipred_sum,  method = "Bagging ipred",            extra_note = "Only final OOB
1510   ↪ misclassification available"),
1511   dplyr::mutate(oob_rf_sum,     method = "Random forest ranger",      extra_note = "Built-in OOB error
1512   ↪ as a single value")
1513 ) %>%
1514   dplyr::select(method, n_trees_final, oob_error_final, oob_error_min, n_trees_at_min,
1515   ↪ extra_note) %>%
1516   dplyr::mutate(
1517     oob_auc    = dplyr::if_else(method == "Bagging custom", bagging_oob_metrics$auc,    NA_real_),
1518     ↪ # add OOB AUC only where OOB probabilities exist
1519     oob_brier = dplyr::if_else(method == "Bagging custom", bagging_oob_metrics$brier, NA_real_)
1520     ↪ # add OOB Brier only where OOB probabilities exist
1521   )
1522
1523 oob_report_tex <- oob_report_tbl %>%                                # convert the OOB
1524   ↪ report into a LaTeX-ready data frame
1525   dplyr::transmute(
1526     Method      = method,                                              # method label shown
1527     ↪ in the paper table
1528     `Number of Trees` = as.integer(n_trees_final),                      # number of trees
1529     ↪ used by each method
1530     `OOB error (misclassification)` = as.numeric(oob_error_final),       # out of bag
1531     ↪ misclassification error
1532     `OOB AUC`    = as.numeric(oob_auc),                                     # OOB AUC when
1533     ↪ available
1534     `OOB Brier`  = as.numeric(oob_brier),                                    # OOB Brier when
1535     ↪ available
1536   ) %>%
1537   as.data.frame()
1538
1539 stargazer::stargazer(                                              # export the OOB
1540   ↪ report table to LaTeX
1541   oob_report_tex,
1542   type = "latex",
1543   summary = FALSE,
1544   rownames = FALSE,
1545   digits = 3,
1546   no.space = TRUE,
1547   table.placement = "!htbp",
1548   title = "Out-of-bag performance summary",
1549   label = "tab:oob_report",
1550   out = "table_oob_report.tex"
1551 )
1552
1553 ## -----
1554 ## 12.6.2 OOB curve plot
1555 ## -----
1556
1557 ipred_oob_final <- as.numeric(oob_err_ipred)                         # store the final
1558   ↪ ipred out of bag misclassification error
1559
1560 oob_plot_df <- tibble::tibble(                                           # build a plotting
1561   ↪ data frame for the custom bagging error path
1562   B        = seq_along(custom_err_path),                                    # tree index from one
1563   ↪ to the final number of trees

```

```

1547   oob_err = custom_err_path                                # out of bag
1548   ↳ misclassification error at each tree index
1549 )
1550
1550 B_star    <- oob_custom_sum$n_trees_at_min                # tree index where the
1551   ↳ custom bagging OOB error is minimized
1551 oob_star <- oob_custom_sum$oob_error_min                 # minimum custom
1551   ↳ bagging OOB error value
1552
1553 y_lo <- min(oob_plot_df$oob_err, ipred_oob_final) - 0.0008 # lower plot limit
1553   ↳ with a small margin
1554 y_hi <- max(oob_plot_df$oob_err, ipred_oob_final) + 0.0008  # upper plot limit
1554   ↳ with a small margin
1555
1556 plot_oob_bagging <- ggplot2::ggplot(oob_plot_df, ggplot2::aes(x = B, y = oob_err)) +
1557   ggplot2::geom_line(linewidth = 1.0, colour = "#0072B2") +          # plot the custom
1557   ↳ bagging OOB error path
1558   ggplot2::geom_point()                                              # mark the minimum
1558   ↳ point on the custom OOB curve
1559   data = tibble::tibble(B = B_star, oob_err = oob_star),
1560     ggplot2::aes(x = B, y = oob_err),
1561     size = 2.0, colour = "#0072B2"
1562 )
1563 ggplot2::geom_hline()                                         # add a horizontal
1563   ↳ reference line for the ipred final OOB error
1564   yintercept = ipred_oob_final,
1565   linetype = 2, linewidth = 0.9, colour = "#D30E00"
1566 )
1567 ggplot2::annotate()                                         # label the ipred
1567   ↳ reference line directly on the plot
1568   "text",
1569   x = max(oob_plot_df$B),
1570   y = ipred_oob_final,
1571   label = paste0("ipred final OOB = ", scales::percent(ipred_oob_final, accuracy = 0.1)),
1572   hjust = 1.05, vjust = -0.7, size = 4.2, colour = "#D30E00"
1573 )
1574 ggplot2::annotate()                                         # label the minimum
1574   ↳ point on the custom OOB curve
1575   "text",
1576   x = B_star,
1577   y = oob_star,
1578   label = paste0("min at B = ", B_star, ": ", scales::percent(oob_star, accuracy = 0.1)),
1579   hjust = -0.05, vjust = 1.2, size = 4.2, colour = "#0072B2"
1580 )
1581 ggplot2::scale_y_continuous()
1582   limits = c(y_lo, y_hi),
1583   labels = scales::percent_format(accuracy = 0.1)
1584 )
1585 ggplot2::labs(x = "Number of trees", y = "OOB misclassification error") + # label axes for
1585   ↳ interpretation
1586 theme_academic() +                                         # apply the project
1586   ↳ theme
1587 theme(
1588   axis.title.x = ggplot2::element_text(size = 14),           # slightly larger axis
1588   ↳ titles for paper readability
1589   axis.title.y = ggplot2::element_text(size = 14),
1590   axis.text.x = ggplot2::element_text(size = 12),
1591   axis.text.y = ggplot2::element_text(size = 12)
1592 )
1593

```

```

1594 print(plot_oob_bagging)                                # display the OOB curve
1595   ↵ plot
1596 ggplot2::ggsave(                                     # save the OOB curve
1597   ↵ plot at high resolution
1598   filename = "plot_oob_bagging.png",
1599   plot      = plot_oob_bagging,
1600   width     = 12, height = 4, units = "in",
1600   dpi       = 1000
1601 )
1602
1603 ## =====
1604 ## 13. Model comparison
1605 ## =====
1606
1607 ## -----
1608 ## 13.1 Sanity checks: required objects from Parts 11 and 12
1609 ## -----
1610
1611 need_objs <- c(                                         # objects that must
1612   ↵ exist before running Part 13                      # common test labels
1612   "y_test", "metrics_binary",
1612   ↵ and metrics helper
1613   "logit_p_test", "glmnet_best_p_test",                # Part 11 test
1613   ↵ probabilities
1614   "tree_prob_test", "bagging_prob_test", "bag_std_prob_test", "rf_pred_test", # Part 12 test
1614   ↵ probabilities
1615   "bagging_oob_metrics", "rf_best_fit"                 # Part 12 out of bag
1615   ↵ info and fitted random forest object
1616 )
1617
1618 missing <- need_objs[!vapply(need_objs, exists, logical(1L))] # identify which
1618   ↵ required objects are missing
1619 if (length(missing) > 0) {                                # stop early with a
1619   ↵ clear message when prerequisites are not met
1620   stop("Part 13 missing objects (run Parts 11 and 12 first):\n ", #,
1621     paste(missing, collapse = ", "))
1622 }
1623
1624 y_test_use <- as.integer(y_test)                           # ensure test labels
1624   ↵ are stored as 0 or 1 integers
1625
1626 ## -----
1627 ## 13.2 Collect predictions in one place (same test set)
1628 ## -----
1629
1630 pred_tbl <- tibble::tibble(                                 # build a single
1630   ↵ table with true labels and model probabilities
1631   y          = y_test_use,                                  # true test label
1631   ↵ stored as 0 or 1
1632   p_logit    = as.numeric(logit_p_test),                  # logit predicted
1632   ↵ probability of default equals one
1633   p_enet     = as.numeric(glmnet_best_p_test),            # elastic net
1633   ↵ predicted probability of default equals one
1634   p_cart     = as.numeric(tree_prob_test),                # CART predicted
1634   ↵ probability of default equals one
1635   p_bag_custom = as.numeric(bagging_prob_test),           # custom bagging
1635   ↵ predicted probability of default equals one
1636   p_bag_ipred = as.numeric(bag_std_prob_test),             # ipred bagging
1636   ↵ predicted probability of default equals one

```

```

1637 p_rf      = as.numeric(rf_pred_test)                      # random forest
1638   ↪ predicted probability of default equals one
1639 )
1640 pred_long <- pred_tbl %>%
1641   ↪ format for plotting and grouped calculations
1642   tidyverse::pivot_longer(
1643     cols      = dplyr::starts_with("p_"),
1644       ↪ probability columns
1645     names_to  = "model",
1646     values_to = "p_hat"
1647   ) %>%
1648   dplyr::mutate(
1649     model = dplyr::recode(
1650       ↪ names to paper-ready model names
1651       model,
1652       p_logit    = "Logit",
1653       p_enet     = "Elastic net",
1654       p_cart     = "CART",
1655       p_bag_custom = "Bagging (custom)",
1656       p_bag_ipred = "Bagging (ipred)",
1657       p_rf       = "Random forest"
1658     ),
1659     model = factor(
1660       ↪ for legends and tables
1661       model,
1662       levels = c("Logit", "Elastic net", "CART",
1663                  "Bagging (custom)", "Bagging (ipred)", "Random forest")
1664     )
1665   )
1666   ## -----
1667   ## 13.3 Main comparison table (computed from stored predictions)
1668   ## -----
1669 metric_one <- function(y, p) {                                # wrapper so metric
1670   ↪ computation reads cleanly below
1671   metrics_binary(y_true = y, p_hat = p)                         # apply the same
1672   ↪ metric definition to every model
1673 }
1674
1675 model_metrics_test <- tibble::tibble(                          # map model names to
1676   ↪ the probability columns in pred_tbl
1677   model = levels(pred_long$model),                               # model names in the
1678   ↪ fixed order
1679   p_col = c("p_logit", "p_enet", "p_cart", "p_bag_custom", "p_bag_ipred", "p_rf")
1680 ) %>%
1681   dplyr::mutate(
1682     met = purrr::map(p_col, ~ metric_one(pred_tbl$y, pred_tbl[[.x]]))  # compute metrics on
1683     ↪ the same holdout test set
1684   ) %>%
1685   tidyverse::unnest(met) %>%                                    # unpack the returned
1686   ↪ metric tibble columns
1687   dplyr::select(model, error, accuracy, auc, brier, tp, fp, tn, fn) %>% # keep only reporting
1688   ↪ columns
1689   dplyr::arrange(dplyr::desc(auc))                                # rank primarily by
1690   ↪ AUC
1691
1692 cat("\n==== Test-set performance (all models, same test set) ====\n") # print a header for
1693   ↪ the comparison table

```

```

1683 print(model_metrics_test)                                # print the model
1684   ↪ comparison table
1685
1685 model_metrics_tex <- model_metrics_test %>%
1686   ↪ LaTeX-export version with cleaned model labels      # create a
1686 dplyr::mutate(
1687   Model = dplyr::case_when(
1688     ↪ naming for the paper table
1688     model %in% c("logit", "Logit", "Logistic") ~ "Logit",
1689     model %in% c("enet", "ENet", "ElasticNet") ~ "Elastic Net",
1690     model %in% c("cart", "CART") ~ "CART",
1691     model %in% c("bag_custom", "Bag (custom)") ~ "Bagging (custom)",
1692     model %in% c("bag_ipred", "Bag (ipred)") ~ "Bagging (ipred)",
1693     model %in% c("rf", "RF", "RandomForest") ~ "Random Forest",
1694     TRUE ~ as.character(model)
1695   )
1696 ) %>%
1697 dplyr::select(
1698   Model,
1699   Error = error,
1700   Accuracy = accuracy,
1701   AUC = auc,
1702   Brier = brier,
1703   TP = tp, FP = fp, TN = tn, FN = fn
1704 ) %>%
1705 as.data.frame()
1706
1707 stargazer::stargazer(
1708   model_metrics_tex,                                         # LaTeX-ready data
1709   ↪ frame with performance metrics
1710   type = "latex",
1711   summary = FALSE,
1711   rownames = FALSE,
1712   digits = 3,                                              # format metrics with a
1713   ↪ consistent number of decimals
1714   no.space = TRUE,
1714   table.placement = "!htbp",
1715   title = "Test-set performance comparison (common test set)",
1716   label = "tab:model_metrics_test",
1717   out = "table_model_metrics_test.tex"
1718 )
1719
1720 ## -----
1721 ## 13.4 Plot A and B: ROC and PR curves (two-panel)
1722 ## -----
1723
1724 roc_df <- pred_long %>%                                # build ROC curve
1725   ↪ points for each model
1725 dplyr::group_by(model) %>%
1726 dplyr::group_modify(~{
1727   r <- pROC::roc(response = .x$y, predictor = .x$p_hat, quiet = TRUE)    # compute the ROC
1728   ↪ curve from probabilities and true labels
1728   tibble::tibble(
1729     fpr = 1 - r$specificities,                                         # false positive rate
1729     ↪ across thresholds
1730     tpr = r$sensitivities,                                            # true positive rate
1730     ↪ across thresholds
1731   )
1732 }) %>%
1733 dplyr::ungroup()
1734

```

```

1735 plot_roc <- ggplot2::ggplot(roc_df, ggplot2::aes(x = fpr, y = tpr, colour = model)) +
1736   ggplot2::geom_line(linewidth = 0.9) +                                     # plot ROC curves for
1737   ggplot2::geom_abline(slope = 1, intercept = 0, linetype = 2, colour = "grey60") + # add random
1738   ggplot2::scale_colour_manual(values = c(                                # apply a fixed
1739     "Logit"           = "#0072B2",
1740     "Elastic net"    = "#D55E00",
1741     "CART"            = "#009E73",
1742     "Bagging (custom)" = "#CC79A7",
1743     "Bagging (ipred)" = "#E69F00",
1744     "Random forest"  = "#56B4E9"
1745   )) +
1746   ggplot2::labs(x = "False positive rate", y = "True positive rate", colour = NULL) +
1747   theme_academic() +
1748   ggplot2::theme(legend.position = "bottom")                                # place legend below
1749   ↪ the plot for paper layout
1750
1750 pr_curve_df <- function(y, p, model_name, n_points = 250L) {                # construct a
1751   ↪ precision recall curve by sweeping thresholds
1752   thr <- stats::quantile(p, probs = seq(0, 1, length.out = n_points), na.rm = TRUE) # build a grid
1753   ↪ of thresholds based on probability quantiles
1754   thr <- sort(unique(as.numeric(thr))), decreasing = TRUE)                      # enforce unique
1755   ↪ thresholds from high to low
1756   out <- lapply(thr, function(t) {
1757     pred <- as.integer(p >= t)                                                 # classify as one when
1758     ↪ probability is at least the threshold
1759     tp <- sum(pred == 1L & y == 1L)                                         # true positives at
1760     ↪ this threshold
1761     fp <- sum(pred == 1L & y == 0L)                                         # false positives at
1762     ↪ this threshold
1763     fn <- sum(pred == 0L & y == 1L)                                         # false negatives at
1764     ↪ this threshold
1765     precision <- ifelse(tp + fp == 0L, NA_real_, tp / (tp + fp))          # precision defined as
1766     ↪ true positives divided by predicted positives
1767     recall <- ifelse(tp + fn == 0L, NA_real_, tp / (tp + fn))              # recall defined as
1768     ↪ true positives divided by actual positives
1769     c(precision = precision, recall = recall)
1770   })
1771   mat <- do.call(rbind, out)                                              # stack the threshold
1772   ↪ results into a matrix
1773   tibble::tibble(
1774     recall    = mat[, "recall"],                                             # recall values across
1775     ↪ thresholds
1776     precision = mat[, "precision"],                                         # precision values
1777     ↪ across thresholds
1778     model     = model_name,                                                 # model name used in
1779     ↪ the legend
1780   ) %>%
1781   dplyr::filter(!is.na(precision), !is.na(recall))                         # drop undefined
1782   ↪ points for stable plotting
1783 }
1784
1784 pr_df <- dplyr::bind_rows(                                              # build PR curve
1785   ↪ points for each model on the same test set
1786   pr_curve_df(pred_tbl$y, pred_tbl$p_logit,      "Logit"),
1787   pr_curve_df(pred_tbl$y, pred_tbl$p_enet,       "Elastic net"),
1788   pr_curve_df(pred_tbl$y, pred_tbl$p_cart,        "CART"),
1789   pr_curve_df(pred_tbl$y, pred_tbl$p_bag_custom, "Bagging (custom)"),
1790   pr_curve_df(pred_tbl$y, pred_tbl$p_bag_ipred,  "Bagging (ipred)"),

```

```

1777 pr_curve_df(pred_tbl$y, pred_tbl$p_rf,           "Random forest")
1778 ) %>%
1779   dplyr::mutate(model = factor(model, levels = levels(pred_long$model)))    # enforce the same
1780   ↪ legend ordering as pred_long
1781
1782 plot_pr <- ggplot2::ggplot(pr_df, ggplot2::aes(x = recall, y = precision, colour = model)) +
1783   ggplot2::geom_line(linewidth = 0.9) +                                     # plot PR curves for
1784   ↪ each model
1785   ggplot2::scale_colour_manual(values = c(                                # reuse the same
1786     ↪ palette as the ROC plot
1787     "Logit"              = "#0072B2",
1788     "Elastic net"        = "#D55E00",
1789     "CART"               = "#009E73",
1790     "Bagging (custom)"  = "#CC79A7",
1791     "Bagging (ipred)"   = "#E69F00",
1792     "Random forest"     = "#56B4E9"
1793   )) +
1794   ggplot2::labs(x = "Recall", y = "Precision", colour = NULL) +
1795   theme_academic() +
1796   ggplot2::theme(legend.position = "bottom")                                # keep the legend
1797   ↪ consistent with the ROC plot
1798
1799 plot_roc_pr <- (plot_roc | plot_pr) +
1800   patchwork::plot_annotation(
1801     tag_levels = "a",                                                 # label panels for
1802     ↪ referencing in the text
1803     tag_prefix = "(",
1804     tag_suffix = ")"
1805
1806 print(plot_roc_pr)                                              # print the combined
1807   ↪ ROC and PR panel
1808
1809 ggsave(
1810   filename = "plot_roc_pr.png",                                         # save the combined
1811   ↪ ROC and PR panel as a high-resolution figure
1812   plot      = plot_roc_pr,
1813   width     = 10, height = 6, units = "in",
1814   dpi       = 1000
1815 )
1816
1817 ## -----
1818 ## 13.5 Plot C: Cumulative gains and lift
1819 ## -----
1820
1821 lift_df <- pred_long %>%                                         # compute gains curves
1822   ↪ by sorting predicted risk within each model
1823   dplyr::group_by(model) %>%
1824     dplyr::arrange(dplyr::desc(p_hat), .by_group = TRUE) %>%          # order borrowers from
1825     ↪ highest predicted risk to lowest
1826     dplyr::mutate(
1827       frac_population      = dplyr::row_number() / dplyr::n(),            # share of the
1828       ↪ population screened at each cutoff
1829       cum_defaults         = cumsum(y),                                     # cumulative number of
1830       ↪ defaults captured as screening expands
1831       total_defaults       = sum(y),                                       # total defaults in
1832       ↪ the test set for this model group
1833       frac_defaults_captured = cum_defaults / total_defaults             # share of defaults
1834       ↪ captured at each cutoff
1835     ) %>%
1836     dplyr::ungroup()

```

```

1825
1826 plot_lift <- ggplot2::ggplot(lift_df, ggplot2::aes(x = frac_population, y =
1827   ↳ frac_defaults_captured, colour = model)) +
1828     ggplot2::geom_line(linewidth = 0.9) +                                # plot gains curves
1829     ↳ for each model
1830     ggplot2::geom_abline(slope = 1, intercept = 0, linetype = 2, colour = "grey60") + # add random
1831     ↳ screening baseline
1832     ggplot2::scale_colour_manual(values = c(                                # reuse the same
1833       ↳ palette for model consistency
1834         "Logit"           = "#0072B2",
1835         "Elastic net"      = "#D55E00",
1836         "CART"             = "#009E73",
1837         "Bagging (custom)" = "#CC79A7",
1838         "Bagging (ipred)"  = "#E69F00",
1839         "Random forest"    = "#56B4E9"
1840     )) +
1841       ggplot2::scale_x_continuous(labels = scales::percent_format(accuracy = 1)) +
1842       ggplot2::scale_y_continuous(labels = scales::percent_format(accuracy = 1)) +
1843       ggplot2::labs(
1844         x = "Share of borrowers screened (highest risk first)",
1845         y = "Share of defaults captured",
1846         colour = NULL
1847     ) +
1848       theme_academic() +
1849       ggplot2::theme(legend.position = "bottom")                                # keep legend
1850     ↳ placement consistent across figures
1851
1852 print(plot_lift)                                                       # print the lift plot
1853
1854 ggpplot2::ggsave(
1855   filename = "plot_lift.png",                                            # save the lift plot
1856   ↳ as a high-resolution figure
1857   plot      = plot_lift,
1858   width     = 12, height = 6, units = "in",
1859   dpi       = 1000
1860 )
1861
1862 ## -----
1863 ## 13.6 Plot D: Predicted risk distributions by class (faceted)
1864 ## -----
1865
1866 plot_risk_dist <- ggplot2::ggplot(
1867   pred_long %>% dplyr::mutate(y_fac = factor(y, levels = c(0, 1), labels = c("No default",
1868     ↳ "Default"))), # label true outcomes for plotting
1869   ggplot2::aes(x = p_hat, fill = y_fac)
1870   ↳ # map predicted risk and outcome group
1871 ) +
1872   ggplot2::geom_density(alpha = 0.35, adjust = 1.1) +                      # plot overlapping
1873   ↳ densities to compare separation by outcome
1874   ggplot2::facet_wrap(~ model, ncol = 3, scales = "free_y") +                # facet by model to
1875   ↳ compare distributions side-by-side
1876   ggplot2::scale_x_continuous(labels = scales::percent_format(accuracy = 1)) +
1877   ggplot2::labs(x = "Predicted default probability", y = "Density", fill = NULL) +
1878   theme_academic() +
1879   ggplot2::theme(legend.position = "bottom")                                    # place legend below
1880   ↳ to match other figures
1881
1882 print(plot_risk_dist)                                                     # print the risk
1883 ↳ distribution figure
1884
1885 ggpplot2::ggsave(

```

```

1874 filename = "plot_risk_dist.png",                                     # save the risk
1875   ↪ distribution figure as a high-resolution plot
1876 plot      = plot_risk_dist,
1877 width     = 10, height = 6, units = "in",
1878 dpi       = 1000
1879 )
1880 ## =====
1881 ## 14. Double machine learning extension: Effect of 60-month term on default
1882 ## =====
1883 ## -----
1884 ## 14.1 DML dataset (train only; imputed; numeric outcome and treatment)
1885 ## -----
1886
1887 dml_df <- train_ml_imp %>%
1888   ↪ imputed training sample                                         # start from the
1889   dplyr::select(-is_train) %>%
1890   ↪ since all rows here are training rows                           # drop the split flag
1891   dplyr::mutate(
1892     default = as.numeric(default),                                    # store the outcome
1893     ↪ as numeric zero or one for DML algebra
1894     term_60 = as.numeric(term_60),                                     # store the treatment
1895     ↪ as numeric zero or one for DML algebra
1896   )
1897
1898 n_dml <- nrow(dml_df)                                              # store the DML
1899   ↪ training sample size
1900 y      <- dml_df$default                                         # extract the outcome
1901   ↪ vector
1902 d      <- dml_df$term_60                                           # extract the
1903   ↪ treatment vector
1904
1905 ## -----
1906 ## 14.2 Cross-fitting setup (precompute X once; fixed lambdas for speed)
1907 ## -----
1908
1909 K <- 5                                                               # set the number of
1910   ↪ folds used for cross-fitting
1911 set.seed(base_seed + 9)                                              # set a seed so fold
1912   ↪ assignment is reproducible
1913
1914 fold_id <- sample(rep(seq_len(K), length.out = n_dml))             # assign roughly
1915   ↪ balanced fold labels
1916
1917 X_full <- model.matrix(                                             # build the one-hot
1918   ↪ encoded covariate matrix and exclude outcome and treatment
1919   ~ . - default - term_60,
1920   data = dml_df
1921 )
1922
1923 alpha_dml <- if (exists("alpha_best")) alpha_best else 0.75          # reuse the selected
1924   ↪ elastic-net mixing value when available
1925
1926 set.seed(base_seed + 10)                                              # set a seed so the
1927   ↪ global lambda selection is reproducible
1928
1929 cv_g_full <- glmnet::cv.glmnet(                                         # select one lambda
1930   ↪ for the outcome model g(X) on the full training sample
1931   x           = X_full,
1932   y           = y,

```

```

1920 family      = "binomial",
1921 alpha       = alpha_dml,
1922 nfolds      = 5,
1923 standardize = TRUE,
1924 type.measure = "deviance"
1925 )
1926
1927 cv_m_full <- glmnet::cv.glmnet(
1928   ↪ for the treatment model  $m(X)$  on the full training sample
1929   x          = X_full,
1930   y          = d,
1931   family     = "binomial",
1932   alpha      = alpha_dml,
1933   nfolds     = 5,
1934   standardize = TRUE,
1935   type.measure = "deviance"
1936 )
1937 lambda_g <- cv_g_full$lambda.min
1938   ↪ lambda for  $g(X)$  equals expected outcome given  $X$ 
1939 lambda_m <- cv_m_full$lambda.min
1940   ↪ lambda for  $m(X)$  equals expected treatment given  $X$ 
1941
1942 m_hat <- numeric(n_dml)                                # allocate
1943   ↪ out-of-fold predicted propensities
1944 g_hat <- numeric(n_dml)                                # allocate
1945   ↪ out-of-fold predicted default risks
1946
1947 ## -----
1948 ## 14.3 Cross-fitted nuisance predictions:  $g_{\text{hat}}$  and  $m_{\text{hat}}$ 
1949 ## -----
1950
1951 for (k in seq_len(K)) {                                    # loop over folds for
1952   ↪ cross-fitting
1953   cat("[DML] Fold", k, "of", K, "...\\n")               # print progress so
1954   ↪ long runs are visible
1955
1956   idx_tr <- which(fold_id != k)                          # training indices
1957   ↪ for this fold
1958   idx_va <- which(fold_id == k)                          # validation indices
1959   ↪ for this fold
1960
1961   fit_m <- glmnet::glmnet(                                 # fit the treatment
1962     ↪ model on the training folds only
1963     x          = X_full[idx_tr, , drop = FALSE],
1964     y          = d[idx_tr],
1965     family     = "binomial",
1966     alpha      = alpha_dml,
1967     lambda     = lambda_m,
1968     standardize = TRUE
1969   )
1970
1971   m_hat[idx_va] <- predict(                               # predict
1972     ↪ propensities on the held-out fold
1973     fit_m,
1974     newx = X_full[idx_va, , drop = FALSE],
1975     s      = lambda_m,
1976     type   = "response"
1977   )[, 1]
1978

```

```

1969 fit_g <- glmnet::glmnet(
1970   ↪ model on the training folds only
1971   x = X_full[idx_tr, , drop = FALSE],
1972   y = y[idx_tr],
1973   family = "binomial",
1974   alpha = alpha_dml,
1975   lambda = lambda_g,
1976   standardize = TRUE
1977 )
1978
1979 g_hat[idx_va] <- predict(
1980   ↪ risks on the held-out fold
1981   fit_g,
1982   newx = X_full[idx_va, , drop = FALSE],
1983   s = lambda_g,
1984   type = "response"
1985 )[, 1]
1986
1987 eps <- 1e-4
1988   ↪ value to avoid extreme predicted probabilities
1989 m_hat <- pmin(pmax(m_hat, eps), 1 - eps)
1990   ↪ propensities away from zero and one
1991 g_hat <- pmin(pmax(g_hat, eps), 1 - eps)
1992   ↪ risks away from zero and one
1993
1994 ## -----
1995 ## 14.4 Orthogonalized regression: theta via residual-on-residual
1996 ## -----
1997
1998 tilde_y <- y - g_hat
1999   ↪ residualized outcome using cross-fitted g_hat
2000 tilde_d <- d - m_hat
2001   ↪ residualized treatment using cross-fitted m_hat
2002
2003 dml_fit <- lm(tilde_y ~ tilde_d)
2004   ↪ residualized outcome on residualized treatment
2005
2006 theta_dml <- unname(coef(dml_fit)[["tilde_d"]])
2007   ↪ effect estimate in probability units
2008 theta_dml_ci <- confint(dml_fit)[["tilde_d", ]]
2009   ↪ ninety-five percent confidence interval
2010 theta_dml_se <- summary(dml_fit)$coefficients[["tilde_d", "Std. Error"]]
2011   ↪ error
2012 theta_dml_pval <- summary(dml_fit)$coefficients[["tilde_d", "Pr(>|t|)"]]
2013   ↪ for the residual-on-residual coefficient
2014
2015 cat("\n[DML] Estimated ATE of 60-month term on default (train sample):\n")
2016   ↪ header
2017 cat(sprintf(" theta_DML = %.4f (%.2f percentage points)\n",
2018   ↪ in probability units and percentage points
2019     theta_dml, 100 * theta_dml))
2020 cat(sprintf(" 95%% CI      = [%4f, %4f]\n",
2021   ↪ confidence interval endpoints
2022     theta_dml_ci[1], theta_dml_ci[2]))
2023 cat(sprintf(" Std. Error = %.4f,  p-value = %.4g\n",
2024   ↪ error and p-value
2025     theta_dml_se, theta_dml_pval))
2026
2027 ## -----
2028 ## 14.5 Naive benchmark: raw difference in default rates (60m minus 36m)

```

```

2014 ## -----
2015
2016 mean_d1 <- mean(y[d == 1])                                # compute the default
2017   ↪ rate among treated loans
2017 mean_d0 <- mean(y[d == 0])                                # compute the default
2018   ↪ rate among control loans
2018 theta_naive <- mean_d1 - mean_d0                          # compute the raw
2019   ↪ difference in means
2020
2020 se_naive <- sqrt(                                         # compute a
2021   ↪ two-sample standard error under independence
2021     var(y[d == 1]) / sum(d == 1) +
2022     var(y[d == 0]) / sum(d == 0)
2023 )
2024
2025 ci_naive <- theta_naive + c(-1, 1) * 1.96 * se_naive      # compute a
2026   ↪ ninety-five percent confidence interval
2027
2027 cat("\n[Naive] Difference in default rates (60m minus 36m):\n") # print a report
2028   ↪ header
2028 cat(sprintf("  theta_naive = %.4f (%.2f percentage points)\n",
2029   ↪ estimate in probability units and percentage points
2029     theta_naive, 100 * theta_naive))
2030 cat(sprintf("  95% CI      = [% .4f, % .4f]\n",
2031   ↪ confidence interval endpoints
2031     ci_naive[1], ci_naive[2]))
2032
2033 dml_comparison <- dplyr::tibble(                               # build a compact
2034   ↪ table comparing naive and DML estimates
2034     estimator = c("Naive diff-in-means", "DML (cross-fit elastic net"),
2035     theta      = c(theta_naive, theta_dml),
2036     ci_low     = c(ci_naive[1], theta_dml_ci[1]),
2037     ci_high    = c(ci_naive[2], theta_dml_ci[2])
2038 )
2039
2040 cat("\n==== DML vs naive effect of 60-month term on default (train sample) ===\n") # print a header
2041   ↪ for the comparison table
2041 print(dml_comparison)                                         # print the
2042   ↪ comparison table
2043
2043 ## =====
2044 ## 15. DML extension II: Effect of high interest rate on default
2045 ## =====
2046
2047 ## -----
2048 ## 15.1 Define high-rate treatment and build DML dataset (train only)
2049 ## -----
2050
2051 high_rate_thr <- quantile(train_ml_imp$int_rate, 0.75, na.rm = TRUE)      # compute the
2052   ↪ seventy-fifth percentile of interest rate in the training sample
2053 cat("\n[High-rate DML] 75th percentile threshold for int_rate:",           # report the
2053     round(high_rate_thr, 2), "%\n")
2053   ↪ threshold that defines the high-rate treatment
2054
2055 dml_hr_df <- train_ml_imp %>%                                     # start from the
2056   ↪ training sample to keep evaluation data separated
2056     dplyr::select(-is_train) %>%                                    # drop the split flag
2057   ↪ since it is not a covariate
2057     dplyr::mutate(
2058       default = as.numeric(default),                                # store the outcome
2058       ↪ as numeric zero or one

```

```

2059   high_rate = as.numeric(int_rate >= high_rate_thr)           # store the treatment
2060   ↪   as numeric one when interest rate is above the threshold
2061 )
2062 n_hr <- nrow(dml_hr_df)                                         # store the DML
2063 ↪   sample size for the high-rate analysis
2064 y_hr <- dml_hr_df$default                                       # extract the outcome
2065 ↪   vector
2066 d_hr <- dml_hr_df$high_rate                                      # extract the
2067 ↪   treatment vector
2068
2069 ## -----
2070 ## 15.2 Cross-fitting setup (precompute X once; fixed lambdas for speed)
2071 ## -----
2072 K_hr <- 5                                                       # set the number of
2073 ↪   folds used for cross-fitting
2074 set.seed(base_seed + 11)                                         # set a seed so fold
2075 ↪   assignment is reproducible
2076
2077 fold_id_hr <- sample(rep(seq_len(K_hr), length.out = n_hr))      # assign roughly
2078 ↪   balanced fold labels
2079
2080 X_full_hr <- model.matrix(                                         # build the one-hot
2081   ↪   encoded covariate matrix and exclude outcome and treatment
2082   ~ . - default - high_rate,
2083   data = dml_hr_df
2084 )
2085
2086 alpha_dml_hr <- if (exists("alpha_best")) alpha_best else 0.75    # reuse the selected
2087 ↪   elastic-net mixing value when available
2088
2089 set.seed(base_seed + 12)                                           # set a seed so
2090 ↪   global lambda selection is reproducible
2091
2092 cv_g_hr_full <- glmnet::cv.glmnet(                                    # select one lambda
2093   ↪   for the outcome model on the full training sample
2094   x = X_full_hr,
2095   y = y_hr,
2096   family = "binomial",
2097   alpha = alpha_dml_hr,
2098   nfolds = 5,
2099   standardize = TRUE,
2100   type.measure = "deviance"
2101 )
2102
2103 cv_m_hr_full <- glmnet::cv.glmnet(                                    # select one lambda
2104   ↪   for the treatment model on the full training sample
2105   x = X_full_hr,
2106   y = d_hr,
2107   family = "binomial",
2108   alpha = alpha_dml_hr,
2109   nfolds = 5,
2110   standardize = TRUE,
2111   type.measure = "deviance"
2112 )
2113
2114 lambda_g_hr <- cv_g_hr_full$lambda.min                                # store the selected
2115 ↪   lambda for the outcome model
2116 lambda_m_hr <- cv_m_hr_full$lambda.min                                 # store the selected
2117 ↪   lambda for the treatment model

```

```

2106 m_hat_hr <- numeric(n_hr)                                     # allocate
2107   ↪  out-of-fold predicted propensities
2108 g_hat_hr <- numeric(n_hr)                                     # allocate
2109   ↪  out-of-fold predicted default risks
2110
2111 ## -----
2112 ## 15.3 Cross-fitted nuisance predictions: m_hat_hr and g_hat_hr
2113 ## -----
2114 for (k in seq_len(K_hr)) {                                       # loop over folds for
2115   ↪  cross-fitting
2116   cat("[High-rate DML] Fold", k, "of", K_hr, "...\\n")          # print progress for
2117   ↪  long runs
2118
2119 idx_tr <- which(fold_id_hr != k)                                # training indices
2120   ↪  for this fold
2121 idx_va <- which(fold_id_hr == k)                                 # validation indices
2122   ↪  for this fold
2123
2124 fit_m_hr <- glmnet::glmnet(                                         # fit the treatment
2125   ↪  model on training folds only
2126   x = X_full_hr[idx_tr, , drop = FALSE],
2127   y = d_hr[idx_tr],
2128   family = "binomial",
2129   alpha = alpha_dml_hr,
2130   lambda = lambda_m_hr,
2131   standardize = TRUE
2132 )
2133
2134 m_hat_hr[idx_va] <- predict(                                         # predict
2135   ↪  propensities on the held-out fold
2136   fit_m_hr,
2137   newx = X_full_hr[idx_va, , drop = FALSE],
2138   s = lambda_m_hr,
2139   type = "response"
2140 )[, 1]
2141
2142 fit_g_hr <- glmnet::glmnet(                                         # fit the outcome
2143   ↪  model on training folds only
2144   x = X_full_hr[idx_tr, , drop = FALSE],
2145   y = y_hr[idx_tr],
2146   family = "binomial",
2147   alpha = alpha_dml_hr,
2148   lambda = lambda_g_hr,
2149   standardize = TRUE
2150 )
2151
2152 g_hat_hr[idx_va] <- predict(                                         # predict baseline
2153   ↪  risks on the held-out fold
2154   fit_g_hr,
2155   newx = X_full_hr[idx_va, , drop = FALSE],
2156   s = lambda_g_hr,
2157   type = "response"
2158 )[, 1]
2159 }
2160
2161 eps_hr <- 1e-4                                                 # choose a small
2162   ↪  value to avoid extreme predicted probabilities
2163 m_hat_hr <- pmin(pmax(m_hat_hr, eps_hr), 1 - eps_hr)          # clamp propensities
2164   ↪  away from zero and one

```

```

2155 g_hat_hr <- pmin(pmax(g_hat_hr, eps_hr), 1 - eps_hr) # clamp risks away
2156   ↵ from zero and one
2157 
2158 ## -----
2159 ## 15.4 Orthogonalised regression: DML estimate for high-rate effect
2160 ## -----
2161 tilde_y_hr <- y_hr - g_hat_hr # compute
2162   ↵ residualized outcome using cross-fitted g_hat_hr
2163 tilde_d_hr <- d_hr - m_hat_hr # compute
2164   ↵ residualized treatment using cross-fitted m_hat_hr
2165 
2166 dml_hr_fit <- lm(tilde_y_hr ~ tilde_d_hr) # regress
2167   ↵ residualized outcome on residualized treatment
2168 
2169 theta_dml_hr <- unname(coef(dml_hr_fit)[ "tilde_d_hr" ]) # store the DML
2170   ↵ effect estimate in probability units
2171 theta_dml_hr_ci <- confint(dml_hr_fit)[ "tilde_d_hr", ] # store the
2172   ↵ ninety-five percent confidence interval
2173 theta_dml_hr_se <- summary(dml_hr_fit)$coefficients[ "tilde_d_hr", "Std. Error"] # store the
2174   ↵ standard error
2175 theta_dml_hr_pval <- summary(dml_hr_fit)$coefficients[ "tilde_d_hr", "Pr(>|t|)"] # store the
2176   ↵ p-value
2177 
2178 cat("\n[High-rate DML] Estimated ATE of high interest rate on default:\n") # print a report
2179   ↵ header
2180 cat(sprintf(" theta_DML_highrate = %.4f (%.2f percentage points)\n", # print the estimate
2181   ↵ in probability units and percentage points
2182     theta_dml_hr, 100 * theta_dml_hr))
2183 cat(sprintf(" 95%% CI           = [% .4f, % .4f]\n", # print the
2184   ↵ confidence interval endpoints
2185     theta_dml_hr_ci[1], theta_dml_hr_ci[2]))
2186 cat(sprintf(" Std. Error        = %.4f, p-value = %.4g\n", # print the standard
2187   ↵ error and p-value
2188     theta_dml_hr_se, theta_dml_hr_pval))
2189 
2190 ## -----
2191 ## 15.5 Naive benchmark: raw difference in default rates (high minus low)
2192 ## -----
2193 
2194 mean_hr1 <- mean(y_hr[d_hr == 1]) # compute the default
2195   ↵ rate among high-rate loans
2196 mean_hr0 <- mean(y_hr[d_hr == 0]) # compute the default
2197   ↵ rate among low-rate loans
2198 
2199 theta_naive_hr <- mean_hr1 - mean_hr0 # compute the raw
2200   ↵ difference in means
2201 
2202 se_naive_hr <- sqrt( # compute a
2203   ↵ two-sample standard error under independence
2204     var(y_hr[d_hr == 1]) / sum(d_hr == 1) +
2205       var(y_hr[d_hr == 0]) / sum(d_hr == 0)
2206 )
2207 
2208 ci_naive_hr <- theta_naive_hr + c(-1, 1) * 1.96 * se_naive_hr # compute a
2209   ↵ ninety-five percent confidence interval
2210 
2211 cat("\n[High-rate naive] Difference in default rates (high minus low):\n") # print a report
2212   ↵ header
2213 cat(sprintf(" theta_naive_highrate = %.4f (%.2f percentage points)\n", # print the naive
2214   ↵ estimate in probability units and percentage points

```

```

2197     theta_naive_hr, 100 * theta_naive_hr))
2198 cat(sprintf(" 95%% CI           = [%.4f, %.4f]\n",
2199   ↪ confidence interval endpoints
2200   ci_naive_hr[1], ci_naive_hr[2]))                                     # print the naive
2201
2202 ## -----
2203 ## 15.6 Compact comparison table for reporting
2204 ## -----
2205
2206 dml_highrate_comparison <- dplyr::tibble(                                # build a
2207   ↪ report-ready comparison table
2208   estimator = c("Naive diff-in-means (high vs low rate)",
2209     "DML (cross-fit elastic net)"),
2210   theta      = c(theta_naive_hr, theta_dml_hr),
2211   ci_low     = c(ci_naive_hr[1], theta_dml_hr_ci[1]),
2212   ci_high    = c(ci_naive_hr[2], theta_dml_hr_ci[2])
2213 )
2214
2215 cat("\n==== Effect of high interest rate on default: naive vs DML ===\n")  # print a header for
2216   ↪ the comparison table
2217 print(dml_highrate_comparison)                                              # print the
2218   ↪ comparison table
2219
2220 ## =====
2221 ## 16. DML Combined comparison table
2222 ## =====
2223
2224 dml_all_prob <- tibble::tibble(                                             # build one table
2225   ↪ that stacks both treatments and both estimators
2226   treatment = c(                                                               # treatment label
2227     ↪ used in the paper table
2228     "60-month term (vs 36-month)",
2229     "60-month term (vs 36-month)",
2230     "High interest (>= 75th pct)",
2231     "High interest (>= 75th pct")
2232   ),
2233   estimator = c(                                                               # estimator label
2234     ↪ used in the paper table
2235     "Naive diff-in-means",
2236     "DML (elastic net, K-fold)",
2237     "Naive diff-in-means",
2238     "DML (elastic net, K-fold)"
2239   ),
2240   theta    = c(theta_naive,          theta_dml,          theta_naive_hr,          theta_dml_hr),      # effect
2241   ↪ estimate in probability units
2242   ci_low   = c(ci_naive[1],        theta_dml_ci[1], ci_naive_hr[1],        theta_dml_hr_ci[1]),  # confidence interval lower bound in probability units
2243   ↪ confidence interval lower bound in probability units
2244   ci_high  = c(ci_naive[2],        theta_dml_ci[2], ci_naive_hr[2],        theta_dml_hr_ci[2]),  # confidence interval upper bound in probability units
2245   ↪ confidence interval upper bound in probability units
2246 ) %>%
2247 dplyr::mutate(
2248   theta_pct = 100 * theta,                                                 # convert the effect
2249   ↪ estimate to percentage points
2250   ci_low_pct = 100 * ci_low,                                              # convert the lower
2251   ↪ confidence bound to percentage points
2252   ci_high_pct = 100 * ci_high,                                              # convert the upper
2253   ↪ confidence bound to percentage points
2254   effect_ci = sprintf("%.2f [% .2f, %.2f]", theta_pct, ci_low_pct, ci_high_pct) # create a
2255   ↪ compact string for reporting effect and interval
2256 )

```

```

2244 cat("\n==== DML extension: naive vs DML for both treatments ===\n")          # print a marker to
2245   ↪ separate output blocks in the console
2246 print()                                         # print a compact
2247   ↪ view for a quick check before exporting
2248 dml_all_prob %>%
2249   dplyr::select(treatment, estimator, effect_ci)
2250
2251 dml_all_tex <- dml_all_prob %>%
2252   ↪ exact table to export to LaTeX
2253 dplyr::transmute(
2254   Treatment = treatment,                         # treatment label
2255   ↪ column used in the manuscript
2256   Estimator = estimator,                        # estimator label
2257   ↪ column used in the manuscript
2258   `Effect (pp) [95% CI]` = effect_ci          # compact effect and
2259   ↪ confidence interval in percentage points
2260 ) %>%
2261 as.data.frame()                                # convert to a base
2262   ↪ data frame for stargazer
2263
2264 stargazer::stargazer(                           # export the LaTeX
2265   ↪ table to file
2266   dml_all_tex,                                 # table object to
2267   ↪ export
2268   type = "latex",                             # output format
2269   summary = FALSE,                            # suppress summary
2270   ↪ statistics
2271   rownames = FALSE,                           # omit row names
2272   no.space = TRUE,                            # use tighter LaTeX
2273   ↪ formatting
2274   digits = 2,                                # set digits option
2275   ↪ even though effect_ci is a string
2276   table.placement = "!htbp",                  # set the LaTeX
2277   ↪ float placement preference
2278   title = "Naive vs DML estimates for two treatments", # table caption
2279   ↪ shown in the manuscript
2280   label = "tab:dml_naive_vs_dml",            # LaTeX label for
2281   ↪ referencing
2282   out = "table_dml_naive_vs_dml.tex"          # output file name
2283 )
2284
2285 ## =====
2286 ## End of file!
2287 ## =====

```