
FIN 417: QUANTITATIVE RISK MANAGEMENT PROJECT II

Credit Risk and Statistical Learning

Philipp Mayer

Rafael Barroso

SCIPER: 329758

SCIPER: 406607

Ingenierie Mathématique

Mathématiques

École Polytechnique Fédérale de Lausanne

December 16, 2025

Abstract

This project investigates the quantitative modeling of credit risk in a retail loan portfolio using synthetic borrower data. Each applicant is characterized by age, income, and employment status, from which repayment probabilities are generated under two underlying data-generating mechanisms. Logistic regression and support vector machine classifiers are trained to estimate default risk and are evaluated using cross-entropy loss, ROC curves, and area-under-the-curve metrics. Building on these predictive models, several lending strategies are analyzed, including unconditional lending and selective lending based on estimated repayment probabilities. Portfolio-level profit and loss distributions are simulated under a true default model, allowing for the estimation of expected returns, Value-at-Risk, and Expected Shortfall. The results highlight the interaction between predictive accuracy, interest rate structure, and downside risk, illustrating the challenges of achieving profitability under conservative lending policies.

Contents

1	Feature generation	3
1.1	Preliminaries	3
1.2	Empirical Analysis and Discussion	3
2	Default model and data generation	5
2.1	Logistic regression and cross-entropy evaluation	5
2.2	Kernel-based classification via Support Vector Machines	5
2.3	Model evaluation via ROC curves and AUC metric	6
3	Lending strategies	8
3.1	Lend to everyone with 5.5% interest	8
3.2	Lend selectively at 1% interest based on logistic classification	9
3.3	Lend selectively at 1% interest based on SVM classification	10
3.4	Conclusion	10
A	Appendix	11
A.1	Loss functions	11
A.2	Python Code	11

1 Feature generation

1.1 Preliminaries

We begin by providing a brief overview of the methodology used to generate the data that will serve as the basis for our numerical experiments. This synthetic dataset represents the set of loan applicants considered in the analysis and is constructed in a controlled manner so that it reflects key characteristics relevant to credit risk modeling. The generated data will then be used to evaluate and compare different credit lending strategies.

We consider a training sample size of $m = 20,000$ and a test sample size of $n = 10,000$. For each borrower $i = 1, \dots, m + n$, we simulate a random feature vector

$$x^i = (x_1^i, x_2^i, x_3^i) \in \mathbb{R}^3,$$

where the three components are defined as follows:

- $x_1^i \sim \text{Unif}[18, 80]$ represents the borrower's age;
- $x_2^i \sim \text{Unif}[1, 15]$ represents the monthly income, measured in thousands of CHF;
- $x_3^i \sim \text{Ber}(0.1)$ is a binary indicator of employment status, where $x_3^i = 1$ corresponds to a self-employed borrower and $x_3^i = 0$ to a salaried borrower.^a

For each borrower i , we assume that the three components of x^i are mutually independent. Moreover, the feature vectors $\{x^i\}_{i=1}^{m+n}$ are assumed to be independent and identically distributed across borrowers.

^aThat is, $\mathbb{P}(x_3^i = 1) = 0.1$ and $\mathbb{P}(x_3^i = 0) = 0.9$.

After obtaining this synthetic data, we then proceed in computing *empirical means* and *standard deviation* of each feature vector ($\{x^i\}_{i=1}^{m+n}$) in order to summarize their empirical distribution. For each coordinate x_j , $j = 1, 2, 3$, we focus on the first m observations, which constitute the training sample.

The empirical mean of feature j is defined as

$$\hat{\mu}_j = \frac{1}{m} \sum_{i=1}^m x_j^i,$$

and provides an estimate of the average value of the corresponding borrower characteristic in the training population.

Similarly, the empirical variance of feature j is given by

$$\hat{\sigma}_j^2 = \frac{1}{m-1} \sum_{i=1}^m (x_j^i - \hat{\mu}_j)^2,$$

with $\hat{\sigma}_j$ denoting the associated empirical standard deviation.

1.2 Empirical Analysis and Discussion

The resulting quantities of the empirical statistics may be found below on Table 1. On another note, we'd like to give a brief discussion on additional creditor characteristics.

In practice, credit scoring models often rely on a bigger set of borrower characteristics in order to capture

Variable	Empirical mean	Empirical standard deviation
Age	48.4547	17.9461
Income	7.9979	4.0281
Employment status	0.1010	0.3013

Table 1: Empirical means and standard deviations of the feature components computed from the training sample.

different dimensions of ‘credit-worthiness’ beyond basic demographic and income information. In addition to the variables considered in this project, other realistic ones that could be incorporated into a credit scoring framework include the following:

- **Debt-to-income ratio**, which measures the borrower’s existing financial obligations relative to their income and provides a direct indicator of repayment capacity;
- **Past defaults or delinquencies**, reflecting previous difficulties in meeting financial commitments and typically serving as a strong predictor of future default risk;
- **Savings amount**, capturing the level of low volatility readily accessible assets in the possession to the borrower as a buffer against adverse income or expenditure shocks;
- **Length and structure of credit history**, which summarizes past interactions with credit products and repayment behavior, although the relevance of this variable can be limited in the Swiss context (due to the absence of a centralized, publicly accessible credit history system).

2 Default model and data generation

We begin by specifying the data-generating mechanisms that describe the borrower repayment behavior.

Let $\xi^1, \dots, \xi^{m+n} \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(0, 1)$ and define the *sigmoid* function

$$\psi(z) = \frac{1}{1 + e^{-z}}.$$

Using this function, we introduce two repayment probability models $p_s : \mathbb{R}^3 \rightarrow (0, 1)$, for $s \in \{1, 2\}$, given by

•

$$p_1(x) = \psi(13.3 - 0.33x_1 + 3.5x_2 - 3x_3),$$

•

$$p_2(x) = \psi\left(5 - 10 [\mathbb{I}_{(-\infty, 25)} + \mathbb{I}_{(75, \infty)}] x_1 + 1.1x_2 - x_3\right).$$

The two functions differ in structure: while p_1 depends linearly on all features, p_2 introduces a non-linear effect in age through indicator functions that penalize very young and very old borrowers.

Based on these repayment probabilities, we construct two datasets $\{(x^i, y_s^i)\}_{i=1}^{m+n}$, for $s = 1, 2$, by defining the binary repayment indicator

$$y_s^i := \mathbb{I}_{\{\xi^i \leq p_s(x^i)\}}.$$

Here, $y_s^i = 1$ indicates that borrower i repays the loan, and $y_s^i = 0$ corresponds to a default, under model s . These two datasets serve as the basis for the subsequent training and evaluation of statistical learning models.

2.1 Logistic regression and cross-entropy evaluation

We estimate repayment probabilities using logistic regression, which assumes the conditional model

$$\mathbb{P}(y = 1 \mid x) = \psi(\beta_0 + \beta^\top x),$$

where ψ denotes the previously defined sigmoid function. The parameter vector (β_0, β) is estimated on the training data by minimizing the empirical *cross-entropy loss* (see appendix A.1 for a formal definition on this evaluation function).

We report the resulting cross-entropy losses on both the training and test sets. For dataset $s = 1$, the training loss is 0.0328, while the test loss is 0.0318. For dataset $s = 2$, the corresponding losses are 0.4852 on the training set and 0.4799 on the test set.

We suspect that the observable difference in magnitude between the losses for $s = 1$ and $s = 2$ reflects the alignment between the model specification and the previously defined data-generating process. While p_1 is linear in the features and therefore well suited to logistic regression, p_2 adds non-linear effects that are not captured by the model, leading to higher cross-entropy losses.

2.2 Kernel-based classification via Support Vector Machines

We now consider a Support Vector Machine (SVM) classifier applied to standardized feature vectors.

Given training data $\{(x^i, y_s^i)\}_{i=1}^m$, the SVM estimates a decision function of the form

$$f(x) = \sum_{i=1}^m \alpha_i k(x, x^i) + b,$$

where k is a Gaussian kernel

$$k(x, x') = \exp\left(-\frac{1}{10}\|x - x'\|^2\right),$$

and the coefficients α_i are obtained by minimizing a *regularized hinge loss* (see appendix A.1 for more details). To enable probabilistic evaluation, class probabilities are attained/recovered via Platt scaling, giving an estimate $\hat{p}_s^{\text{svm}}(x)$.

We evaluate the SVM models using cross-entropy loss on both the training and test sets. For dataset $s = 1$, the training loss is 0.0402, while the test loss is 0.0386. For dataset $s = 2$, the corresponding losses are 0.2672 on the training set and 0.2501 on the test set.

Compared to logistic regression, the SVM achieves slightly higher cross-entropy loss for dataset $s = 1$, but significantly improves performance for dataset $s = 2$. This reflects the ability of the kernel-based model to capture non-linear structures (regarding feature relationships) in the repayment probability that are not explicitly modeled by a linear classifier.

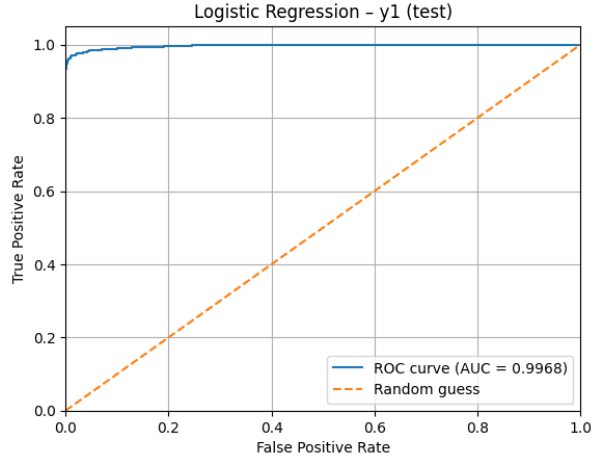
2.3 Model evaluation via ROC curves and AUC metric

We evaluate the predictive performance of the logistic regression and SVM models using Receiver Operating Characteristic (ROC) curves on the *test* set, and summarize each curve through its Area Under the Curve (AUC) (i.e. typical Riemann sum under the curve). A model with AUC close to 1 achieves strong separation between repaid loans ($y = 1$) and defaults ($y = 0$) across a wide range of decision thresholds, where in the other hand an AUC close to 0.5 corresponds to near-random classification.

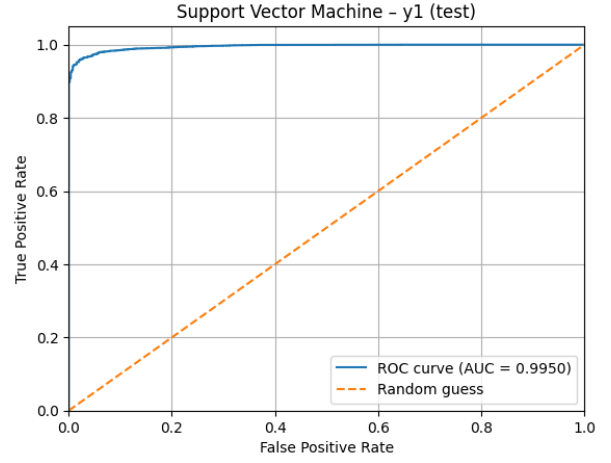
Figure 1 collects the ROC curves for both models and both datasets. For dataset $s = 1$, both curves lie very close to the top-left corner, which imply an excellent discriminative power. The logistic regression model attains an AUC of 0.9968, while the SVM achieves an AUC of 0.9950. The difference is small (about 0.18% relative), which may hint at the fact that both methods perform essentially ‘optimally’ in this setting.

For dataset $s = 2$, the ROC curves are visibly less aligned near the top-left corner, indicating a more challenging classification task. In this case, the logistic regression curve is substantially closer to the diagonal compared to the SVM curve, and the AUC values reflect this gap: logistic regression achieves an AUC of 0.8403, whereas the SVM reaches 0.9596. This corresponds to an improvement of roughly 12.1% in AUC when using the SVM instead of logistic regression on dataset $s = 2$.

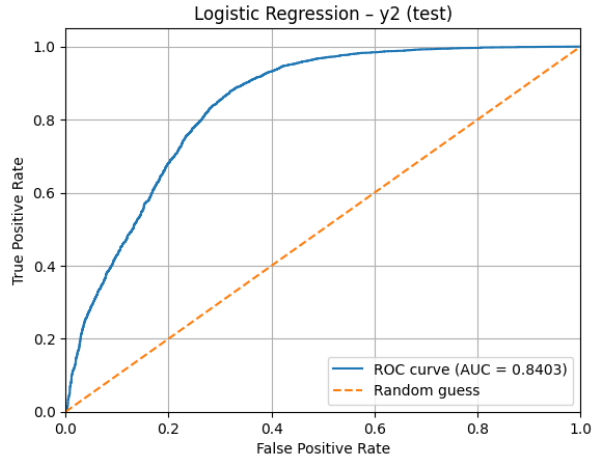
Overall, Figure 1 shows that while logistic regression and SVM perform similarly on dataset $s = 1$, the SVM provides a clear improvement on dataset $s = 2$ in terms of ROC shape and AUC, indicating better separation of repayment and default outcomes in the second setting.



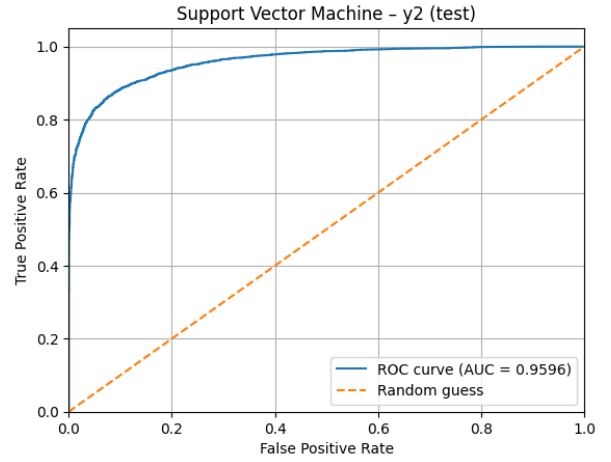
(a) Logistic regression, dataset $s = 1$ (AUC = 0.9968).



(b) SVM, dataset $s = 1$ (AUC = 0.9950).



(c) Logistic regression, dataset $s = 2$ (AUC = 0.8403).



(d) SVM, dataset $s = 2$ (AUC = 0.9596).

Figure 1: ROC curves on the test set for logistic regression and SVM, for both datasets $s = 1$ and $s = 2$. The dashed diagonal in each plot corresponds to random guessing (AUC = 0.5).

3 Lending strategies

For this last part, we only focus on dataset $s = 2$, which we treat as the true data-generating mechanism which underlies the borrower repayment behavior.

We assume that each loan has a fixed principal^a of 1000 CHF and that repayment is all-or-nothing. We consider the following lending policies, applied to the test set borrowers $i = m + 1, \dots, m + n$:

1. lend to all applicants at an interest rate of 5.5%;
2. lend selectively at an interest rate of 1%, but only to applicants satisfying $\hat{p}_2^{\log}(x^i) \geq 95\%$;
3. lend selectively at an interest rate of 1%, but only to applicants satisfying $\hat{p}_2^{\text{SVM}}(x^i) \geq 95\%$.

^aOriginal sum of money borrowed (like a loan) or invested.

Here, $\hat{p}_2^{\log}(x^i)$ and $\hat{p}_2^{\text{SVM}}(x^i)$ denote the estimated repayment probabilities for borrower i obtained from the logistic regression and SVM models trained on dataset $s = 2$, respectively. To assess the risk and profitability of each policy, we simulate 50,000 repayment scenarios. For each borrower i and scenario k , we draw $\xi^{i,k} \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(0, 1)$ and define the repayment indicator

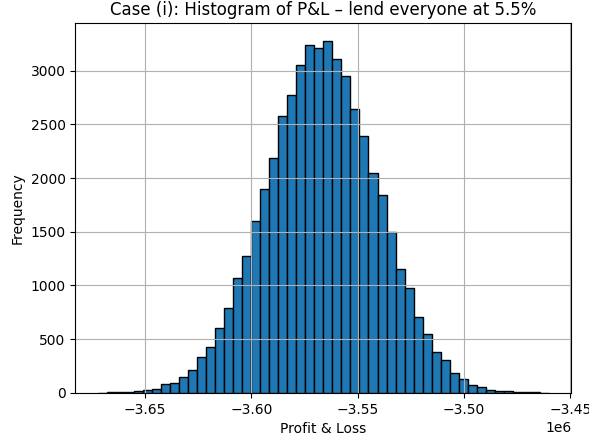
$$D^{i,k} := \mathbb{I}_{\{\xi^{i,k} \leq p_2(x^{m+i})\}},$$

where p_2 is the true repayment probability under dataset $s = 2$. Adding profits and losses across borrowers gives an empirical distribution of portfolio profit and loss (P&L) for each lending strategy.

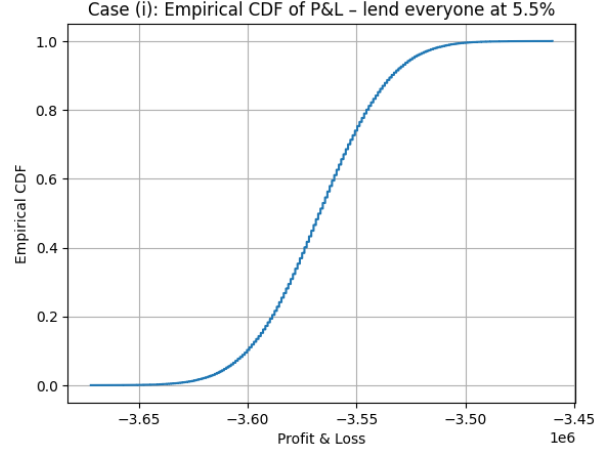
3.1 Lend to everyone with 5.5% interest

Figure 2 shows the empirical distribution of portfolio P&L under the policy of lending to all applicants at an interest rate of 5.5%. The histogram in Figure 2a shows a concentrated distribution centered far in the negative region, which shows that losses occur in the majority of simulated scenarios. This reflects the fact that, while the interest rate is relatively high, defaults among riskier borrowers generate large losses that dominate interest income.

The corresponding empirical CDF in Figure 2b confirms this observation: almost the entire distribution lies below zero, and the probability of achieving a positive portfolio return is effectively none. Although this strategy involves no model risk or selection error, it exposes the lender to heavy downside risk due to the unfiltered default exposure.



(a) Histogram of portfolio P&L under the policy of lending to all applicants at 5.5% interest.



(b) Empirical CDF of portfolio P&L under the policy of lending to all applicants at 5.5% interest.

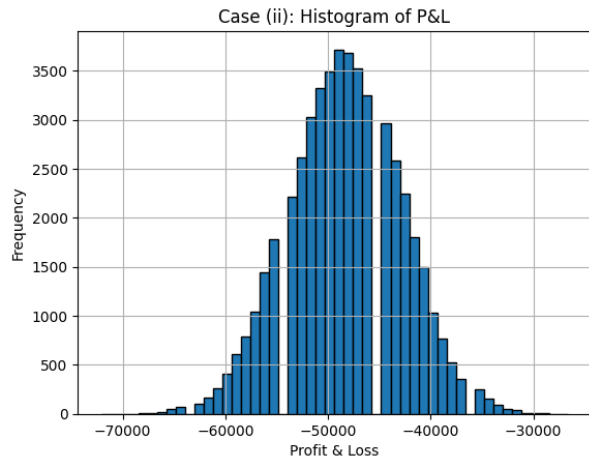
Figure 2: Case 1: Lend to everyone at 5.5% interest.

The $\text{VaR}_{95\%}$ is 3609865 CHF and $\text{ES}_{95\%}$ is 3620303 CHF.

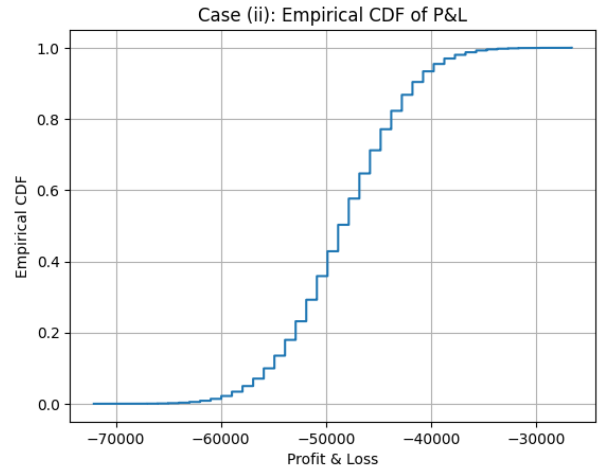
3.2 Lend selectively at 1% interest based on logistic classification

Figure 3 presents the P&L distribution obtained by lending selectively at 1% interest using the logistic regression model. Compared to the unconditional lending strategy, the histogram in Figure 3a shows a noticeable shift of the distribution towards less negative values, reflecting the exclusion of applicants considered high-risk by the model.

But, the distribution is still entirely concentrated in the negative region, and the empirical CDF in Figure 3b indicates that large losses still occur with high probability. This behavior demonstrates the asymmetric risk–return profile of low-interest lending: while the model reduces default frequency, the small gains from successful repayments are insufficient to cancel out losses from the remaining defaults. As a result, the expected portfolio P&L remains negative despite selective lending.



(a) Histogram of portfolio P&L under selective lending at 1% interest based on the logistic regression model.



(b) Empirical CDF of portfolio P&L under selective lending at 1% interest based on the logistic regression model.

Figure 3: Case 2: Lend selectively at 1% interest, but only to applicants with $\hat{p}_2^{\text{log}}(x^i) \geq 95\%$

The $\text{VaR}_{95\%}$ is 56940 CHF and $\text{ES}_{95\%}$ is 58933 CHF.

3.3 Lend selectively at 1% interest based on SVM classification

Figure 4 shows the P&L distribution obtained when selective lending is performed using the SVM classification model. The histogram in Figure 4a is more shifted towards higher values relative to the logistic case, showing improved risk screening. This is consistent with the ‘better’ classification performance of the SVM observed in the ROC and AUC analysis for dataset $s = 2$.

Needless to say, the empirical CDF in Figure 4b reveals that the distribution still places noticeable probability mass on negative outcomes. While the SVM-based strategy reduces downside risk compared to the logistic approach, the low interest rate limits the achievable upside. So, even the most accurate classification model considered does not suffice to render the lending strategy profitable under the chosen economic assumptions.

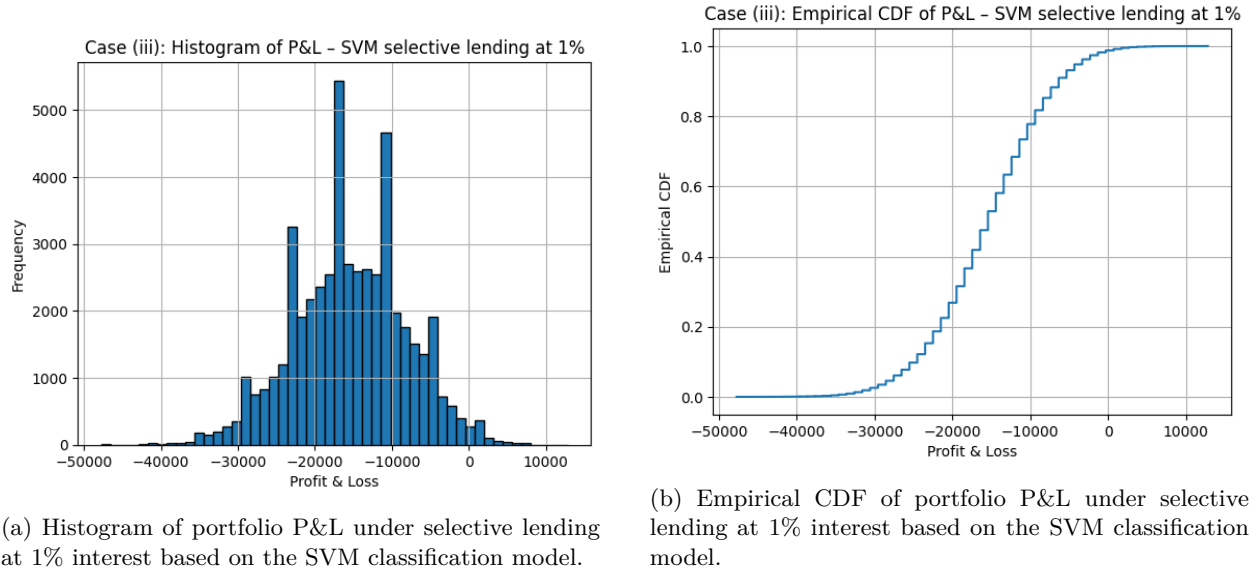


Figure 4: Case 3: Lend selectively at 1% interest, but only to applicants with $\hat{p}_2^{\text{svm}}(x^i) \geq 95\%$

The $\text{VaR}_{95\%}$ is 27550 CHF and $\text{ES}_{95\%}$ is 30319 CHF.

3.4 Conclusion

In conclusion, looking at all lending strategies, we observe significantly high Value at Risks and Expected Shortfall indicating non-profitable lending strategies. We believe this could be explained by multiple factors. First of all, all-or-nothing repayment wouldn’t be a viable contract, nor underwritten by most insurers in the majority of cases. Second of all, we believe that the repayment probabilities computed from the synthetic dataset are not realistic (lending to unemployed people without taking personal fortune into account for instance). We believe that in practice, these factors would lead to stricter risk management rules and thus higher probability of making a profit per contract signed.

Moreover, analyzing the underlying distributions for the synthetic feature generation reveals a lack of consideration for the true population looking for a credit. Especially since most individual characteristics have cross dependencies that influence their realization, it could be useful to consider a more robust approach which is able to capture these relationships.

Finally, we consider the performance of the SVM model to hold certain validity upon the classification task. Although, if a more revised approach is given into the creditor’s generation, a more robust machinery could required in order to capture these new complexities in a more proficient manner.

A Appendix

A.1 Loss functions

We briefly recall the definitions of the loss functions used for training and evaluating the classification models.

Cross-entropy loss. Given predicted probabilities $\hat{p}(x) = \mathbb{P}(y = 1 \mid x)$ and binary labels $y \in \{0, 1\}$, the cross-entropy loss is defined as

$$\ell_{\text{CE}}(y, \hat{p}) = -[y \log(\hat{p}) + (1 - y) \log(1 - \hat{p})].$$

The empirical cross-entropy loss over a dataset $\{(x^i, y^i)\}_{i=1}^N$ is given by

$$\frac{1}{N} \sum_{i=1}^N \ell_{\text{CE}}(y^i, \hat{p}(x^i)).$$

This loss is used throughout the report to evaluate probabilistic predictions from both logistic regression and SVM models.

Hinge loss. For labels $y \in \{-1, +1\}$ and a real-valued decision function $f(x)$, the hinge loss is defined as

$$\ell_{\text{hinge}}(y, f(x)) = \max\{0, 1 - yf(x)\}.$$

In the SVM framework, model parameters are obtained by minimizing a regularized empirical hinge loss of the form

$$\frac{1}{m} \sum_{i=1}^m \ell_{\text{hinge}}(y^i, f(x^i)) + \lambda \|f\|^2,$$

where $\lambda > 0$ is a regularization parameter. Although the SVM is trained using hinge loss, probabilistic outputs are obtained via Platt scaling, allowing model performance to be evaluated using cross-entropy loss.

A.2 Python Code

This appendix contains the Python code used to generate the synthetic data, estimate the statistical learning models, and evaluate the different lending strategies discussed in the main text. The implementation follows a modular design and closely mirrors the methodology described in the report. Executing the provided scripts reproduces all tables and figures presented in the main sections.

```
1  # %% [markdown]
2  # # Project 2: Credit Risk and Statistical Learning
3
4  # %% [markdown]
5  # **Names of all group members:**
6  # - Rafael Barroso (rafael.barroso@epfl.ch)
7  # - Philipp Mayer (philipp.mayer@epfl.ch)
8
9  # %%
10 ## IMPORT NECESSARY LIBRARIES
11 import numpy as np
12 import pandas as pd
13 import matplotlib.pyplot as plt
14
15 from sklearn.linear_model import LogisticRegression
16 from sklearn.metrics import log_loss
17 from sklearn.preprocessing import StandardScaler
18 from sklearn.svm import SVC
```

```

19 from sklearn.metrics import roc_curve, auc
20
21
22 # %% [markdown]
23 # ## 1: Feature generation
24
25 # %%
26 def gen_features(m: int, n: int, random_seed: int = 0) -> pd.DataFrame:
27     """
28     Generate (m + n) feature vectors for loan applicants.
29
30     Each applicant i has:
31         x_i1: age ~ Uniform[18, 80]
32         x_i2: monthly income (CHF) ~ Uniform[1, 15]
33         x_i3: employment status ~ Bernoulli(p = 0.1); so we label as
34             1 = self-employed
35             0 = salaried
36
37     All three coordinates are assumed independent.
38
39     Parameters
40     -----
41     m : int
42         Number of training samples.
43     n : int
44         Number of test samples.
45     random_seed : int or None, optional
46
47     Returns
48     -----
49     features_df : pandas.DataFrame
50         DataFrame with (m + n) rows and 3 columns:
51         - "age"
52         - "income"
53         - "employment"
54     """
55
56     total_samples = m + n
57     rng = np.random.default_rng(seed=random_seed) # Use numpy's random generator
58
59     ## Simulate features
60     ages = rng.integers(18, 80, size=total_samples)
61     incomes = rng.uniform(low=1.0, high=15.0, size=total_samples)
62     employment_status = rng.binomial(1, 0.1, size=total_samples)
63
64     # Create DataFrame
65     features_df = pd.DataFrame({
66         "age": ages,
67         "income": incomes,
68         "employment": employment_status
69     })
70
71     return features_df
72
73 def compute_empirical_stats(features: pd.DataFrame, m: int) -> tuple[pd.Series, pd.Series]:
74     """
75     Compute empirical mean and variance of each feature using the training set.
76
77     Parameters
78     -----
79     features : pandas.DataFrame
80         DataFrame containing features for all applicants.
81     m : int
82         Number of training samples.
83
84     Returns
85     -----
86     means : pandas.Series
87         Empirical means of each feature.
88     variances : pandas.Series

```

```

89         Empirical variances of each feature.
90         """
91
92         training_features = features.iloc[:m]
93
94         means = training_features.mean()
95         stds = training_features.std()
96
97         return means, stds
98
99     # %%
100     m, n = 20000, 10000
101     features = gen_features(m, n, random_seed=0)
102     means, stds = compute_empirical_stats(features, m)
103     print(features.head())
104     print("Empirical Means:\n", means)
105     print("Empirical Standard Deviations:\n", stds)
106
107     # %%
108     (features['employment'] == 0).sum()
109
110     # %% [markdown]
111     # dayum that's a lot of unemployed people lol
112
113     # %% [markdown]
114     # Some other realistic variables that might be used for credit scoring are:
115     # - Debt to income ratio
116     # - Past defaults/delinquencies
117     # - Savings ammount
118     # - Existing credit history (Although for Switzerland this is not usefull)
119
120     # %% [markdown]
121     # ## 2: **Default model and data generation**
122
123     # %%
124     def sigmoid(z: np.ndarray) -> np.ndarray:
125         """
126         Standard logistic (sigmoid) function:
127              $\sigma(z) = 1 / (1 + \exp(-z))$ 
128         """
129         return 1.0 / (1.0 + np.exp(-z))
130
131     def p1_probability(features: pd.DataFrame) -> np.ndarray:
132         """
133         Compute repayment probabilities  $p_1(x)$  for each row in `features`.
134
135         Model 1:
136              $p_1(x) = \text{sigmoid}(13.3 - 0.33 * \text{age} + 3.5 * \text{income} - 3 * \text{employment})$ 
137
138         Here:
139             - age = features["age"]
140             - income = features["income"]
141             - employment = features["employment"] (0 or 1)
142
143         Returns
144         -----
145         probs : np.ndarray
146             Array of repayment probabilities  $p_1(x_i)$  with shape (n_samples,).
147         """
148         age = features["age"].to_numpy()
149         income = features["income"].to_numpy()
150         emp = features["employment"].to_numpy()
151
152         # Linear predictor for model 1
153         z = 13.3 - 0.33 * age + 3.5 * income - 3.0 * emp
154
155         # Apply sigmoid to obtain probabilities in (0,1)
156         probs = sigmoid(z)
157         return probs
158

```

```

159 def p2_probability(features: pd.DataFrame) -> np.ndarray:
160     """
161     Compute repayment probabilities  $p_2(x)$  for each row in `features`.
162
163     Model 2 (piecewise in age via indicator functions):
164
165     Let:
166         young = 1 if age <= 25, else 0
167         old   = 1 if age >= 75, else 0
168
169     Then:
170          $p_2(x) = \text{sigmoid}(-5 - 10 * \text{young} - 10 * \text{old} + 1.1 * \text{income} - 1.0 * \text{employment})$ 
171
172     Returns
173     -----
174     probs : np.ndarray
175         Array of repayment probabilities  $p_2(x_i)$  with shape (n_samples,).
176     """
177     age = features["age"].to_numpy()
178     income = features["income"].to_numpy()
179     emp = features["employment"].to_numpy()
180
181     # Indicator variables for age
182     # young = 1 if age in (18, 25], old = 1 if age in [75, 80]
183     young = (age <= 25.0).astype(float)
184     old = (age >= 75.0).astype(float)
185
186     # Linear predictor for model 2
187     z = -5.0 - 10.0 * young - 10.0 * old + 1.1 * income - 1.0 * emp
188
189     probs = sigmoid(z)
190     return probs
191
192 def generate_repayment_data(features: pd.DataFrame, random_seed: int = 0) -> pd.DataFrame:
193     """
194     Given a matrix of features X (as a DataFrame), generate repayment
195     outcomes  $y_1$  and  $y_2$  for both models using a common Uniform(0,1)
196     random variable  $\epsilon_i$  for each applicant  $i$ .
197
198     Parameters
199     -----
200     features : pd.DataFrame
201         Must contain at least the columns:
202         "age", "income", "employment".
203         It should have (m + n) rows.
204     random_seed : int or None, optional
205
206     Returns
207     -----
208     data : pd.DataFrame
209         Copy of `features` with 4 new columns:
210         - "p1_true":  $p_1(x_i)$ 
211         - "p2_true":  $p_2(x_i)$ 
212         - "y1": repayment indicator under model 1
213         - "y2": repayment indicator under model 2
214     """
215     n_samples = len(features)
216
217     # Set up RNG
218     rng = np.random.default_rng(seed=random_seed)
219
220     # 1) Compute true repayment probabilities under each model
221     p1 = p1_probability(features) # shape (n_samples,)
222     p2 = p2_probability(features) # shape (n_samples,)
223
224     # 2) Draw a single Uniform(0,1) value  $\epsilon_i$  for each applicant  $i$ 
225     eps = rng.uniform(low=0.0, high=1.0, size=n_samples)
226
227     # 3) Generate repayment indicators using the same  $\epsilon_i$  for both models
228     y1 = (eps <= p1).astype(int) # 1 = repay, 0 = default

```

```

229     y2 = (eps <= p2).astype(int)
230
231     # 4) Build output DataFrame (do not modify input in-place)
232     data = features.copy()
233     data["p1_true"] = p1
234     data["p2_true"] = p2
235     data["y1"] = y1
236     data["y2"] = y2
237
238     return data
239
240 # %%
241 data = generate_repayment_data(features, random_seed=0)
242 print(data.head(-10))
243
244 # %% [markdown]
245 # ### (A): **Fit Logistic Regression model**
246
247 # %%
248 ## HELPER FUNCTION TO SPLIT TRAIN/TEST SETS
249 def tts_split(
250     data: pd.DataFrame,
251     target_col: str,
252     m: int,
253     feature_cols: list[str] = None,
254 ):
255     """
256     Split the full DataFrame into training and test sets.
257
258     Parameters
259     -----
260     data : pd.DataFrame
261         Full dataset with (m + n) rows.
262     target_col : str
263         Name of the target column to predict ("y1" or "y2").
264     m : int
265         Number of training samples (first m rows are training).
266     feature_cols : list of str or None, optional
267         Names of feature columns. If None, we default to:
268         ["age", "income", "employment"].
269
270     Returns
271     -----
272     X_train : np.ndarray of shape (m, d)
273     y_train : np.ndarray of shape (m,)
274     X_test : np.ndarray of shape (n, d)
275     y_test : np.ndarray of shape (n,)
276     """
277     # If no specific feature list is given, use the standard three
278     if feature_cols is None:
279         feature_cols = ["age", "income", "employment"]
280
281     # Extract feature matrix and target vector as NumPy arrays
282     X = data[feature_cols].to_numpy()
283     y = data[target_col].to_numpy()
284
285     # First m rows: training set
286     X_train = X[:m, :]
287     y_train = y[:m]
288
289     # Remaining rows: test set
290     X_test = X[m:, :]
291     y_test = y[m:]
292
293     return X_train, y_train, X_test, y_test
294
295 # %% [markdown]
296 # can also use `train_test_split` prebuilt sklearn function but anyways...
297
298 # %%

```

```

299 def fit_logistic_regression(X_train: np.ndarray, y_train: np.ndarray) -> LogisticRegression:
300     """
301     Fit a simple logistic regression model using scikit-learn.
302
303     We do NOT do any fancy preprocessing here:
304     - no feature scaling,
305     - default L2 regularization,
306     - just more iterations to ensure convergence.
307
308     Parameters
309     -----
310     X_train : np.ndarray
311         Training feature matrix of shape (m, d).
312     y_train : np.ndarray
313         Training labels (0/1) of shape (m,).
314
315     Returns
316     -----
317     model : LogisticRegression
318         Fitted logistic regression model.
319     """
320     model = LogisticRegression(
321         solver="lbfgs", # standard solver
322         max_iter=1000, # more iterations just to be safe
323     )
324     model.fit(X_train, y_train)
325     return model
326
327 def cross_entropy_loss(y_true: np.ndarray, p_hat: np.ndarray) -> float:
328     """
329     Compute the (average) cross-entropy loss for binary labels.
330
331     Formula:
332         
$$L = - (1/N) * \sum_{i=1}^N [ y_i * \log(p_i) + (1 - y_i) * \log(1 - p_i) ]$$

333
334     where:
335     -  $y_i$  are the true labels in  $\{0,1\}$ 
336     -  $p_i$  are the predicted probabilities  $P(y=1 | x_i)$ 
337
338     Parameters
339     -----
340     y_true : np.ndarray
341         Vector of true labels in  $\{0,1\}$ .
342     p_hat : np.ndarray
343         Vector of predicted probabilities in  $[0,1]$ .
344
345     Returns
346     -----
347     loss : float
348         Cross-entropy loss (smaller is better, if you know how to use it hehe).
349     """
350     # Small epsilon to avoid log(0)
351     eps = 1e-15
352
353     # Clip predicted probabilities to stay away from 0 and 1
354     p_hat_clipped = np.clip(p_hat, eps, 1.0 - eps)
355
356     # Compute each term  $y_i * \log(p_i) + (1 - y_i) * \log(1 - p_i)$ 
357     term1 = y_true * np.log(p_hat_clipped)
358     term2 = (1 - y_true) * np.log(1.0 - p_hat_clipped)
359
360     # Average cross-entropy
361     loss = -np.mean(term1 + term2)
362     return loss
363
364 def evaluate_logistic_model(model: LogisticRegression, X_train: np.ndarray, y_train:
↪ np.ndarray, X_test: np.ndarray, y_test: np.ndarray):
365     """
366     Compute predicted probabilities and cross-entropy losses
367     for both TRAINING and TEST sets.

```



```

368
369     Parameters
370     -----
371     model : LogisticRegression
372         Fitted logistic regression.
373     X_train, y_train : np.ndarray
374         Training data.
375     X_test, y_test : np.ndarray
376         Test data.
377
378     Returns
379     -----
380     results : dict
381         Dictionary with:
382         - "p_train": predicted probabilities on train set
383         - "p_test" : predicted probabilities on test set
384         - "ce_train": cross-entropy on train set
385         - "ce_test" : cross-entropy on test set
386
387     """
388     # predict_proba returns a matrix with 2 columns:
389     # [:,0] = P(y=0 | x), [:,1] = P(y=1 | x)
390     p_train = model.predict_proba(X_train)[:, 1]
391     p_test = model.predict_proba(X_test)[:, 1]
392
393     # Compute cross-entropy losses
394     ce_train = cross_entropy_loss(y_train, p_train)
395     ce_test = cross_entropy_loss(y_test, p_test)
396
397     results = {
398         "p_train": p_train,
399         "p_test": p_test,
400         "ce_train": ce_train,
401         "ce_test": ce_test,
402     }
403     return results
404
405 # %%
406 X1_train_log1, y1_train_log1, X1_test_log1, y1_test_log1 = tts_split(data, target_col="y1",
407     ↪ m=m)
408 log_model1 = fit_logistic_regression(X1_train_log1, y1_train_log1)
409 log_results1 = evaluate_logistic_model(log_model1, X1_train_log1, y1_train_log1, X1_test_log1,
410     ↪ y1_test_log1)
411 print("Logistic Regression Model 1 Cross-Entropy Loss:")
412 print(" Train:", log_results1["ce_train"])
413 print(" Test :", log_results1["ce_test"])
414
415 # %%
416 X1_train_log2, y1_train_log2, X1_test_log2, y1_test_log2 = tts_split(data, target_col="y2",
417     ↪ m=m)
418 log_model2 = fit_logistic_regression(X1_train_log2, y1_train_log2)
419 log_results2 = evaluate_logistic_model(log_model2, X1_train_log2, y1_train_log2, X1_test_log2,
420     ↪ y1_test_log2)
421 print("Logistic Regression Model 2 Cross-Entropy Loss:")
422 print(" Train:", log_results2["ce_train"])
423 print(" Test :", log_results2["ce_test"])
424
425 # %% [markdown]
426 # ### (B): **Fit Support Vector Machine model**
427
428 # %%
429 def standardize_train_test(X_train: np.ndarray, X_test: np.ndarray):
430     """
431     Standardize features:
432
433     For each feature j:
434         mean_j = average of column j on the TRAINING set
435         std_j = standard deviation of column j on the TRAINING set
436
437     Then:
438         X_std[:, j] = (X[:, j] - mean_j) / std_j
439

```

```

434
435 IMPORTANT:
436 -----
437 - We fit the scaler on the TRAINING data only.
438 - We then apply the same transformation to the TEST data.
439
440 Parameters
441 -----
442 X_train : np.ndarray
443 Training feature matrix of shape (m, d).
444 X_test : np.ndarray
445 Test feature matrix of shape (n, d).
446
447 Returns
448 -----
449 X_train_std : np.ndarray
450 Standardized training features.
451 X_test_std : np.ndarray
452 Standardized test features.
453 scaler : StandardScaler
454 Fitted scaler (in case you want to reuse it).
455 """
456 scaler = StandardScaler()
457 X_train_std = scaler.fit_transform(X_train)
458 X_test_std = scaler.transform(X_test)
459
460 return X_train_std, X_test_std, scaler
461
462 def fit_svm_classifier(X_train_std: np.ndarray, y_train: np.ndarray, C: float = 1.0, kernel:
↳ str = "rbf", gamma: str = "scale") -> SVC:
463 """
464 Fit a Support Vector Machine (SVM) classifier with probability outputs.
465
466 We keep the setup simple:
467 - RBF kernel (non-linear)
468 - Default C and gamma (you can change them later if needed)
469 - probability=True so we can call predict_proba and compute cross-entropy
470
471 Parameters
472 -----
473 X_train_std : np.ndarray
474 Standardized training features of shape (m, d).
475 y_train : np.ndarray
476 Binary training labels (0/1) of shape (m,).
477 C : float, optional
478 Regularization parameter. Default = 1.0.
479 kernel : str, optional
480 Kernel type. Default = "rbf".
481 gamma : str, optional
482 Kernel coefficient. "scale" is scikit-learn's default.
483
484 Returns
485 -----
486 model : SVC
487 Fitted SVM classifier.
488 """
489 model = SVC(
490     C=C,
491     kernel=kernel,
492     gamma=gamma,
493     probability=True, # needed to get probabilities
494 )
495 model.fit(X_train_std, y_train)
496 return model
497
498 def evaluate_svm_model(model: SVC, X_train_std: np.ndarray, y_train: np.ndarray, X_test_std:
↳ np.ndarray, y_test: np.ndarray):
499 """
500 Compute predicted probabilities and cross-entropy losses
501 for both TRAINING and TEST sets using the already defined

```

```

502     `cross_entropy_loss` function.
503
504     Parameters
505     -----
506     model : SVC
507         Fitted SVM classifier.
508     X_train_std, y_train : np.ndarray
509         Standardized training data and labels.
510     X_test_std, y_test : np.ndarray
511         Standardized test data and labels.
512
513     Returns
514     -----
515     results : dict
516         Dictionary with:
517         - "p_train": predicted probabilities on train set
518         - "p_test" : predicted probabilities on test set
519         - "ce_train": cross-entropy on train set
520         - "ce_test" : cross-entropy on test set
521     """
522     # As with logistic regression, predict_proba returns:
523     # [:,0] = P(y=0 | x), [:,1] = P(y=1 | x)
524     p_train = model.predict_proba(X_train_std)[:, 1]
525     p_test = model.predict_proba(X_test_std)[:, 1]
526
527     ce_train = cross_entropy_loss(y_train, p_train)
528     ce_test = cross_entropy_loss(y_test, p_test)
529
530     results = {
531         "p_train": p_train,
532         "p_test": p_test,
533         "ce_train": ce_train,
534         "ce_test": ce_test,
535     }
536     return results
537
538     # %%
539     X2_train_svm1, y2_train_svm1, X2_test_svm1, y2_test_svm1 = tts_split(data, target_col="y1",
540     ↪ m=m)
541     X_train, X_test, scaler = standardize_train_test(X2_train_svm1, X2_test_svm1)
542     svm_model1 = fit_svm_classifier(X2_train_svm1, y2_train_svm1)
543     results_svm1 = evaluate_svm_model(svm_model1, X2_train_svm1, y2_train_svm1, X2_test_svm1,
544     ↪ y2_test_svm1)
545     print("SVM Model 1 Cross-Entropy Loss:")
546     print("  Train:", results_svm1["ce_train"])
547     print("  Test :", results_svm1["ce_test"])
548
549     # %%
550     X2_train_svm2, y2_train_svm2, X2_test_svm2, y2_test_svm2 = tts_split(data, target_col="y2",
551     ↪ m=m)
552     X_train, X_test, scaler = standardize_train_test(X2_train_svm2, X2_test_svm2)
553     svm_model2 = fit_svm_classifier(X2_train_svm2, y2_train_svm2)
554     results_svm2 = evaluate_svm_model(svm_model2, X2_train_svm2, y2_train_svm2, X2_test_svm2,
555     ↪ y2_test_svm2)
556     print("SVM Model 2 Cross-Entropy Loss:")
557     print("  Train:", results_svm2["ce_train"])
558     print("  Test :", results_svm2["ce_test"])
559
560     # %%
561     ## HELPER FUNCTIONS FOR ROC AUC ANALYSIS
562     def compute_roc_auc(y_true: np.ndarray, p_hat: np.ndarray):
563         """
564         Compute ROC curve and AUC given true labels and predicted probabilities.
565
566         Parameters
567         -----
568         y_true : np.ndarray
569             Binary true labels (0/1), shape (n_samples,)
570         p_hat : np.ndarray
571             Predicted probabilities P(y=1|x), shape (n_samples,) (i.e. prob of y=1 given x)

```

```

568
569     Returns
570     -----
571     fpr : np.ndarray
572         False positive rate values for different thresholds.
573     tpr : np.ndarray
574         True positive rate values for different thresholds.
575     roc_auc : float
576         Area under the ROC curve.
577     """
578     # roc_curve returns (FPR, TPR, thresholds)
579     fpr, tpr, _ = roc_curve(y_true, p_hat) # we want only true positive rate and false
580     ↪ positive rate, no need thresholds
581     roc_auc = auc(fpr, tpr)
582     return fpr, tpr, roc_auc
583
584 def plot_roc_curve(fpr, tpr, roc_auc: float, title: str):
585     """
586     Plot ROC curve with diagonal reference line and AUC in the legend.
587
588     Parameters
589     -----
590     fpr : np.ndarray
591         False positive rate values.
592     tpr : np.ndarray
593         True positive rate values.
594     roc_auc : float
595         Area under the ROC curve.
596     title (trivial) : str
597     """
598     plt.figure()
599     # ROC curve
600     plt.plot(fpr, tpr, label=f"ROC curve (AUC = {roc_auc:.4f})")
601     # Diagonal baseline (naive or random classifier)
602     plt.plot([0, 1], [0, 1], linestyle="--", label="Random guess")
603
604     plt.xlim([0.0, 1.0])
605     plt.ylim([0.0, 1.05])
606     plt.xlabel("False Positive Rate")
607     plt.ylabel("True Positive Rate")
608     plt.title(title)
609     plt.legend(loc="lower right")
610     plt.grid(True)
611     plt.show()
612
613 # %%
614 fpr_log_y1, tpr_log_y1, auc_log_y1 = compute_roc_auc(y1_test_log1, log_results1["p_test"])
615 plot_roc_curve(fpr_log_y1, tpr_log_y1, auc_log_y1, "Logistic Regression - y1 (test)")
616
617 # %% [markdown]
618 # This is a NICE curve. Usually we want it to be as L shaped as possible which would translate
619 ↪ to having a perfect true positive rate. In other words, we barely missed any positive
620 ↪ classifications.
621
622 # %%
623 fpr_log_y2, tpr_log_y2, auc_log_y2 = compute_roc_auc(y1_test_log2, log_results2["p_test"])
624 plot_roc_curve(fpr_log_y2, tpr_log_y2, auc_log_y2, "Logistic Regression - y2 (test)")
625
626 # %% [markdown]
627 # Less 'nicer'.
628
629 # %%
630 fpr_svm_y1, tpr_svm_y1, auc_svm_y1 = compute_roc_auc(y2_test_svm1, results_svm1["p_test"])
631 plot_roc_curve(fpr_svm_y1, tpr_svm_y1, auc_svm_y1, "Support Vector Machine - y1 (test)")
632
633 # %% [markdown]
634 # Compared to the logarithmic regression model, we may start our assumptions by stating that
635 ↪ since we get a slightly 'bigger' area under the curve (which we have yet to discuss). In
636 ↪ other words, we have a better *rate* at classifying correct positives using log model.
637 #

```

```

633 # Needless to say, this is also a pretty good classification.
634
635 # %%
636 fpr_svm_y2, tpr_svm_y2, auc_svm_y2 = compute_roc_auc(y2_test_svm2, results_svm2["p_test"])
637 plot_roc_curve(fpr_svm_y2, tpr_svm_y2, auc_svm_y2, "Support Vector Machine - y2 (test)")
638
639 # %% [markdown]
640 # Here we actually see the efficiency of the svm model. Although for the first probabilistic
641   ↳ test the logarithmic regression yields superior metrics, here the svm model clearly comes
642   ↳ out on top. As a preliminary result, we might state that if we take both probabilistic
643   ↳ approaches into consideration, the svm model performs better on average (taking into
644   ↳ account both probs.)
645
646 # %% [markdown]
647 # ## 3: **Lending Strategies**
648
649 # %%
650 ## HELPER FUNCTIONS FOR P&L ANALYSIS
651 def compute_pl_expectation(pl_samples: np.ndarray) -> float:
652     """
653     Compute the EXPECTED profit & loss (P&L).
654
655     Here, P&L is a random variable:
656     - positive values -> profit
657     - negative values -> loss
658
659     The expectation is simply:
660      $E[P\&L] \text{ is approximately } = (1/N) * \sum_i pl_i$ 
661
662     Parameters
663     -----
664     pl_samples : np.ndarray
665         1D array of simulated P&L values.
666
667     Returns
668     -----
669     expected_pl : float
670         Sample mean of the P&L.
671     """
672     expected_pl = np.mean(pl_samples)
673     return expected_pl
674
675 def plot_pl_histogram(pl_samples: np.ndarray,
676                       title: str = "Histogram of Profit & Loss",
677                       bins: int = 50):
678     """
679     Plot a histogram of simulated P&L values.
680
681     Parameters
682     -----
683     pl_samples : np.ndarray
684         1D array of simulated P&L values.
685     title : str, optional
686         Title for the plot.
687     bins : int, optional
688         Number of bins in the histogram.
689
690     Notes
691     -----
692     - The vertical axis shows how many simulations fall into each
693       P&L interval.
694     """
695     plt.figure()
696     plt.hist(pl_samples, bins=bins, edgecolor="black")
697     plt.xlabel("Profit & Loss")
698     plt.ylabel("Frequency")
699     plt.title(title)
700     plt.grid(True)
701     plt.show()
702

```

```

699 def estimate_var_es(pl_samples: np.ndarray,
700                     alpha: float = 0.95) -> tuple[float, float]:
701     """
702     Estimate Value-at-Risk (VaR) and Expected Shortfall (ES)
703     at confidence level `alpha` based on simulated P&L.
704
705     Convention used here:
706     -----
707     - P&L > 0 : profit
708     - P&L < 0 : loss
709
710     We define LOSS = -P&L.
711     -> large positive LOSS corresponds to large negative P&L.
712
713     For losses L:
714
715         VaR_alpha = quantile_alpha(L)
716         ES_alpha = average of losses that are WORSE than VaR_alpha, i.e.
717
718             ES_alpha = E[L | L >= VaR_alpha]
719
720     Parameters
721     -----
722     pl_samples : np.ndarray
723         1D array of simulated P&L values.
724     alpha : float, optional
725         Confidence level (e.g. 0.95 for 95% VaR and ES).
726
727     Returns
728     -----
729     var_alpha : float
730         Estimated VaR at level `alpha` (in loss units).
731     es_alpha : float
732         Estimated ES at level `alpha` (in loss units).
733     """
734     # Convert P&L to losses: L = -P&L
735     losses = -pl_samples
736
737     # 1) VaR is the alpha-quantile of the loss distribution
738     var_alpha = np.quantile(losses, alpha)
739
740     # 2) ES is the average loss in the tail beyond VaR
741     tail_losses = losses[losses >= var_alpha]
742     es_alpha = tail_losses.mean() if tail_losses.size > 0 else var_alpha
743
744     return var_alpha, es_alpha
745
746 def plot_pl_distribution(pl_samples: np.ndarray,
747                         title: str = "Empirical CDF of Profit & Loss"):
748     """
749     Plot the empirical distribution function (CDF) of P&L.
750
751     Steps:
752     -----
753     1) Sort the P&L samples.
754     2) For each sorted value  $x_{(i)}$ , define  $F(x_{(i)}) = i / N$ ,
755        where  $N$  is the number of samples.
756     3) Plot  $x_{(i)}$  vs  $F(x_{(i)})$ .
757
758     Parameters
759     -----
760     pl_samples : np.ndarray
761         1D array of simulated P&L values.
762     title : str, optional
763         Title for the plot.
764
765     Notes
766     -----
767     - The CDF shows, for any value  $x$ , the proportion of scenarios
768       with P&L  $\leq x$ .

```

```

769     """
770     # Sort P&L values
771     pl_sorted = np.sort(pl_samples)
772
773     # Empirical CDF values from 1/N to N/N
774     n = pl_sorted.size
775     cdf = np.arange(1, n + 1) / n
776
777     plt.figure()
778     plt.step(pl_sorted, cdf, where="post")
779     plt.xlabel("Profit & Loss")
780     plt.ylabel("Empirical CDF")
781     plt.title(title)
782     plt.grid(True)
783     plt.show()
784
785     def simulate_repayment_scenarios_p2(
786         data: pd.DataFrame,
787         m: int,
788         num_scenarios: int = 50000,
789         random_seed: int | None = 123,
790     ) -> np.ndarray:
791         """
792         Simulate repayment scenarios for the TEST set using the true  $p_2(x)$ .
793
794         We only use the TEST set:
795              $i = m, \dots, m+n-1$ 
796
797         For each test borrower  $i$  and each scenario  $k$ , we draw:
798              $D_{\{i,k\}} = 1\{U_{\{i,k\}} \leq p_2(x_i)\}$ ,
799         where  $U_{\{i,k\}} \sim \text{Unif}(0,1)$  are independent.
800
801         Parameters
802         -----
803         data : pd.DataFrame
804             Full dataset with at least the column "p2_true".
805         m : int
806             Number of TRAINING observations (first  $m$  rows). The rest are TEST.
807         num_scenarios : int, optional
808             Number of scenarios  $K$  to simulate.
809         random_seed : int or None, optional
810             Seed for reproducibility.
811
812         Returns
813         -----
814         D : np.ndarray
815             2D array of shape  $(n_{\text{test}}, \text{num\_scenarios})$  with entries in  $\{0,1\}$ ,
816             where  $n_{\text{test}} = \text{len}(\text{data}) - m$ .
817              $D[i,k] = 1$  means borrower  $i$  (in test set) repays in scenario  $k$ .
818         """
819         rng = np.random.default_rng(seed=random_seed)
820
821         # Extract true  $p_2(x)$  for test borrowers
822         p2_all = data["p2_true"].to_numpy() # shape (m + n,)
823         p2_test = p2_all[m:] # shape (n,)
824
825         n_test = p2_test.size
826
827         # Draw  $U_{\{i,k\}} \sim \text{Unif}(0,1)$  for all  $i,k$ 
828         # Shape: (n_test, num_scenarios)
829         U = rng.uniform(low=0.0, high=1.0, size=(n_test, num_scenarios))
830
831         # Repayment indicator: 1 if  $U \leq p_2(x)$ , else 0
832         D = (U <= p2_test[:, np.newaxis]).astype(int)
833
834         return D
835
836     # %%
837     D_test = simulate_repayment_scenarios_p2(data, m=m, num_scenarios=50000, random_seed=123)
838

```

```

839 # %%
840 # Case (i)
841 def compute_pl_lend_everyone(
842     D: np.ndarray,
843     loan_amount: float = 1000.0,
844     interest_rate: float = 0.055,
845 ) -> np.ndarray:
846     """
847     Compute scenario P&L for policy (i): lend to EVERY test borrower
848     at a fixed interest rate.
849
850     We work scenario by scenario. For each borrower i and scenario k:
851
852         If  $D_{\{i,k\}} = 1$  (repayment):
853             - We get back principal + interest:  $\text{loan\_amount} * (1 + r)$ 
854             - We had to put up the principal:  $\text{loan\_amount}$ 
855             -> Profit on this loan =  $\text{loan\_amount} * r$ 
856
857         If  $D_{\{i,k\}} = 0$  (default):
858             - We lose the principal:  $-\text{loan\_amount}$ 
859             -> Profit on this loan =  $-\text{loan\_amount}$ 
860
861     So, for each loan:
862          $\text{P\&L}_{\{i,k\}} = \text{loan\_amount} * [ r * D_{\{i,k\}} - 1 * (1 - D_{\{i,k\}}) ]$ 
863
864     Total P&L in scenario k = sum over all test borrowers i.
865
866     Parameters
867     -----
868     D : np.ndarray
869         Repayment indicator matrix of shape (n_test, num_scenarios),
870         with entries 0 or 1.
871     loan_amount : float, optional
872         Principal per loan (CHF). Default: 1000.
873     interest_rate : float, optional
874         Interest rate
875
876     Returns
877     -----
878     pl_scenarios : np.ndarray
879         1D array of length num_scenarios, where each entry is the
880         TOTAL P&L across all test loans in that scenario.
881     """
882     r = interest_rate
883
884     # Compute P&L per loan and scenario:
885     # shape (n_test, num_scenarios)
886     pl_per_loan = loan_amount * (r * D - (1 - D))
887
888     # Sum across loans (axis=0) to get total P&L per scenario
889     pl_scenarios = pl_per_loan.sum(axis=0)
890
891     return pl_scenarios
892
893 # %%
894 pl_case_i = compute_pl_lend_everyone(
895     D_test,
896     loan_amount=1000.0,
897     interest_rate=0.055,
898 )
899
900 exp_pl_i = compute_pl_expectation(pl_case_i)
901 var_95_i, es_95_i = estimate_var_es(pl_case_i, alpha=0.95)
902
903 print("E[P&L] =", exp_pl_i)
904 print("VaR_95 (loss units) =", var_95_i)
905 print("ES_95 (loss units) =", es_95_i)
906
907 plot_pl_histogram(pl_case_i, "Case (i): Histogram of P&L - lend everyone at 5.5%", bins=50)
908 plot_pl_distribution(pl_case_i, "Case (i): Empirical CDF of P&L - lend everyone at 5.5%")

```



```

909
910 # %%
911 # case (ii)
912 def compute_lending_indicator_logistic_y2(
913     data: pd.DataFrame,
914     m: int,
915     log_model_y2,
916     feature_cols: list[str] | None = None,
917     threshold: float = 0.95,
918 ) -> np.ndarray:
919     """
920     Compute a 0/1 lending indicator for EACH TEST applicant based on
921     the logistic regression model fitted for y2.
922
923     We do:
924
925         1) Extract TEST features  $X_{test}$  (rows  $m, \dots, m+n-1$ ).
926         2) Compute predicted probabilities:
927              $p\_hat = P(y_2 = 1 \mid x, \text{logistic model})$ .
928         3) Define lending rule:
929              $lend\_i = 1$  if  $p\_hat\_i \geq \text{threshold}$ , else 0.
930
931     Parameters
932     -----
933     data : pd.DataFrame
934         Full dataset with at least the feature columns.
935     m : int
936         Number of training observations (first m rows). The rest are TEST.
937     log_model_y2 :
938         Fitted LogisticRegression model for target y2.
939         This is the model you obtained earlier with fit_logistic_regression.
940     feature_cols : list of str or None, optional
941         Names of feature columns. If None, defaults to:
942         ["age", "income", "employment"]
943     threshold : float, optional
944         Lending threshold on predicted probability (default: 0.95).
945
946     Returns
947     -----
948     lend_indicator : np.ndarray
949         1D array of length  $n_{test} = \text{len}(\text{data}) - m$ , with entries in {0,1}.
950          $lend\_indicator[i] = 1$  means "lend to test borrower i".
951     """
952     if feature_cols is None:
953         feature_cols = ["age", "income", "employment"]
954
955     # Extract feature matrix X for all borrowers
956     X_all = data[feature_cols].to_numpy()
957
958     # Split into train and test manually (we only need test here)
959     X_test = X_all[m:, :] # rows  $m, \dots, m+n-1$ 
960
961     # Predicted probabilities on TEST set:
962     # predict_proba returns  $[:,0] = P(y=0)$ ,  $[:,1] = P(y=1)$ 
963     p_hat_test = log_model_y2.predict_proba(X_test)[:, 1]
964
965     # Lending decision: lend if predicted prob  $\geq$  threshold
966     lend_indicator = (p_hat_test >= threshold).astype(int)
967
968     return lend_indicator
969
970 def compute_pl_lend_selective(
971     D: np.ndarray,
972     lend_indicator: np.ndarray,
973     loan_amount: float = 1000.0,
974     interest_rate: float = 0.01,
975 ) -> np.ndarray:
976     """
977     Compute scenario P&L for a selective lending policy.
978

```

```

979     Inputs:
980     -----
981     - D: 2D array of shape (n_test, num_scenarios)
982       D[i,k] = 1 if test borrower i repays in scenario k
983       = 0 if test borrower i defaults in scenario k
984
985     - lend_indicator: 1D array of shape (n_test,)
986       lend_indicator[i] = 1 if we decide to lend to borrower i
987       = 0 otherwise
988
989     P&L per loan and scenario:
990
991     If lend_indicator[i] = 1:
992
993       If D[i,k] = 1 (repayment):
994         P&L_{i,k} = + loan_amount * interest_rate
995
996       If D[i,k] = 0 (default):
997         P&L_{i,k} = - loan_amount
998
999     If lend_indicator[i] = 0:
1000       P&L_{i,k} = 0
1001
1002     We compute this vectorized for all i,k.
1003
1004     Parameters
1005     -----
1006     D : np.ndarray
1007       Repayment indicator matrix (n_test, num_scenarios), entries 0/1.
1008     lend_indicator : np.ndarray
1009       Lending indicator vector (n_test,), entries 0/1.
1010     loan_amount : float, optional
1011       Principal per loan, default 1000 CHF.
1012     interest_rate : float, optional
1013       Interest rate
1014
1015     Returns
1016     -----
1017     pl_scenarios : np.ndarray
1018       1D array of length num_scenarios, total P&L per scenario.
1019     """
1020     r = interest_rate
1021
1022     # Compute P&L per loan & scenario as if we lent to everyone:
1023     # shape (n_test, num_scenarios)
1024     pl_per_loan_full = loan_amount * (r * D - (1 - D))
1025
1026     # Apply the lending indicator: if lend_indicator[i] = 0,
1027     # then the entire row i becomes 0.
1028     # We do this by multiplying each row i by lend_indicator[i].
1029     # Expand lend_indicator to shape (n_test, 1) so broadcasting works.
1030     lend_matrix = lend_indicator[:, np.newaxis] # shape (n_test, 1)
1031     pl_per_loan = pl_per_loan_full * lend_matrix
1032
1033     # Total P&L in each scenario: sum across all test borrowers (axis=0)
1034     pl_scenarios = pl_per_loan.sum(axis=0)
1035
1036     return pl_scenarios
1037
1038     # %%
1039     lend_indicator_log2 = compute_lending_indicator_logistic_y2(
1040         data, m=m, log_model_y2=log_model2, threshold=0.95
1041     )
1042     pl_case_ii = compute_pl_lend_selective(
1043         D_test,
1044         lend_indicator=lend_indicator_log2,
1045         loan_amount=1000.0,
1046         interest_rate=0.01,
1047     )
1048

```

```

1049 exp_pl_ii = compute_pl_expectation(pl_case_ii)
1050 var_95_ii, es_95_ii = estimate_var_es(pl_case_ii, alpha=0.95)
1051
1052 print("Case (ii): E[P&L] =", exp_pl_ii)
1053 print("Case (ii): VaR_95 =", var_95_ii)
1054 print("Case (ii): ES_95 =", es_95_ii)
1055
1056 plot_pl_histogram(pl_case_ii, "Case (ii): Histogram of P&L", bins=50)
1057 plot_pl_distribution(pl_case_ii, "Case (ii): Empirical CDF of P&L")
1058
1059 # %%
1060 # case (iii)
1061 def compute_lending_indicator_svm_y2(
1062     data: pd.DataFrame,
1063     m: int,
1064     svm_model_y2: SVC,
1065     scaler_y2: StandardScaler,
1066     feature_cols: list[str] | None = None,
1067     threshold: float = 0.95,
1068 ) -> np.ndarray:
1069     """
1070     Compute a 0/1 lending indicator for EACH TEST applicant based on
1071     the SVM model fitted for y2.
1072
1073     Steps:
1074     -----
1075     1) Extract feature matrix for all borrowers.
1076     2) Take the TEST rows (indices m,...,m+n-1).
1077     3) Standardize them using the SAME scaler used to train the SVM.
1078     4) Compute predicted probabilities:
1079         p_hat_svm2(x_i) = P(y2 = 1 | x_i, SVM model)
1080     5) Lending rule on TEST set:
1081         lend_i = 1 if p_hat_svm2(x_i) >= threshold, else 0.
1082
1083     Parameters
1084     -----
1085     data : pd.DataFrame
1086         Full dataset with at least the feature columns.
1087     m : int
1088         Number of TRAINING observations (first m rows). The rest are TEST.
1089     svm_model_y2 : SVC
1090         Fitted SVM classifier for target y2 with probability=True.
1091     scaler_y2 : StandardScaler
1092         Fitted scaler used to standardize the features for svm_model_y2.
1093     feature_cols : list of str or None, optional
1094         Names of feature columns.
1095         If None, defaults to ["age", "income", "employment"].
1096     threshold : float, optional
1097         Lending threshold on predicted probability (default 0.95).
1098
1099     Returns
1100     -----
1101     lend_indicator : np.ndarray
1102         1D array of length n_test = len(data) - m, with entries in {0,1}.
1103         lend_indicator[i] = 1 means "lend to test borrower i".
1104     """
1105     if feature_cols is None:
1106         feature_cols = ["age", "income", "employment"]
1107
1108     # All features
1109     X_all = data[feature_cols].to_numpy()
1110
1111     # TEST features: rows m,...,m+n-1
1112     X_test = X_all[m:, :]
1113
1114     # Standardize using the SAME scaler as in training
1115     X_test_std = scaler_y2.transform(X_test)
1116
1117     # Probabilities on TEST set
1118     # predict_proba returns[:,0] = P(y=0),[:,1] = P(y=1)

```

```

1119     p_hat_test = svm_model_y2.predict_proba(X_test_std)[: , 1]
1120
1121     # Lending decision: lend if predicted prob >= threshold
1122     lend_indicator = (p_hat_test >= threshold).astype(int)
1123
1124     return lend_indicator
1125
1126 # %%
1127 lend_indicator_svm2 = compute_lending_indicator_svm_y2(
1128     data=data,
1129     m=m,
1130     svm_model_y2=svm_model2,
1131     scaler_y2=scaler, # the StandardScaler used for SVM training
1132     feature_cols=["age", "income", "employment"],
1133     threshold=0.95, # lend only if predicted prob >= 95%
1134 )
1135
1136 pl_case_iii = compute_pl_lend_selective(
1137     D=D_test,
1138     lend_indicator=lend_indicator_svm2,
1139     loan_amount=1000.0,
1140     interest_rate=0.01, # 1% interest
1141 )
1142
1143 exp_pl_iii = compute_pl_expectation(pl_case_iii)
1144 print(f"Case (iii) - E[P&L]   {exp_pl_iii:.2f} CHF")
1145
1146 plot_pl_histogram(
1147     pl_case_iii,
1148     title="Case (iii): Histogram of P&L - SVM selective lending at 1%",
1149     bins=50,
1150 )
1151
1152 plot_pl_distribution(
1153     pl_case_iii,
1154     title="Case (iii): Empirical CDF of P&L - SVM selective lending at 1%",
1155 )
1156
1157 var_95_iii, es_95_iii = estimate_var_es(pl_case_iii, alpha=0.95)
1158 print(f"Case (iii) - 95% VaR (loss units): {var_95_iii:.2f} CHF")
1159 print(f"Case (iii) - 95% ES  (loss units): {es_95_iii:.2f} CHF")
1160
1161

```