

Deep Credit Risk

Machine Learning in Python

Daniel Rösch*
Harald Scheule†

*DANIEL RÖSCH is a professor of business and holds the chair of statistics and risk management at the University of Regensburg, Germany.

†HARALD SCHEULE is a professor of finance at the University of Technology Sydney, Australia.

Deep Credit Risk

Machine Learning in Python

1st edition, 2020

ISBN: 9798617590199

Impressum:
Copyright © 2020 Daniel Rösch, Harald Scheule
All rights reserved.

Daniel Rösch, Harald Scheule
c/o AutorenServices.de
König-Konrad-Str. 22
D-36039 Fulda

Limit of Liability/Disclaimer of Warranty: While the publisher and author have used their best effort in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

Contents

I Principles of Data Learning	11
1 Deep Dive	12
1.1 Boost Your Practical Machine Learning Skills	12
1.2 Credit Risk Information	13
1.3 Hands-on Analysis	14
1.4 Basel, CECL, IFRS 9, DFAST, CCAR and Stress Tests	16
1.5 First Lessons from the COVID-19 Crisis	17
1.5.1 Credit Risk Analytics	17
1.5.2 Machine Learning	19
2 Python Literacy	20
2.1 Synopsis	20
2.2 Installation	20
2.2.1 Anaconda and IDEs	20
2.2.2 Packages	21
2.2.3 Coding Guidelines	22
2.3 First Look	23
2.4 Describing Data	24
2.5 Plotting	25
2.6 Generating New Variables	26
2.7 Transforming Variables	27
2.8 Subsetting Data	28
2.9 Resetting Indexes	30
2.10 Combining Data	30
2.10.1 Concatenating	30
2.10.2 Appending	31
2.10.3 Match Merging	31
2.10.4 Joining	32
2.11 Regression Models	32
2.12 numpy vs pandas	34
2.13 Packages and Basic Settings	36
2.14 Functions	36
2.14.1 versions	36
2.14.2 dataprep	37
2.14.3 woe	38
2.14.4 validation	38
2.14.5 resolutionbias	39
2.15 Sandbox Problems	39
3 Risk-Based Learning	40
3.1 Synopsis	40
3.2 Maximum-Likelihood Estimation	40
3.2.1 Example for Default Modeling	42
3.2.2 Practical Implementation	48
3.3 Bayesian Approaches	49
3.3.1 Distributions	50
3.3.2 Parameter Estimation	51

3.3.3	Example for Default Modeling	52
3.3.4	Analytic Computation	56
3.3.5	Markov-Chain-Monte-Carlo Simulation	58
3.4	Sandbox Problems	60
4	Machine Learning	61
4.1	Synopsis	61
4.2	Terminology	62
4.3	Cost Functions	62
4.4	Information and Entropy	63
4.5	Optimization: Gradient Descent	67
4.6	Learning and Validation	75
4.6.1	Train vs Test Split	75
4.6.2	Bias-Variance Tradeoff	75
4.6.3	Crossvalidation and Tuning	80
4.7	Practical Implementation	80
4.8	P-Value and ML Hacking	83
4.9	Sandbox Problems	83
II	Data Processing and Validation	84
5	Outcome Engineering	85
5.1	Synopsis	85
5.2	Outcomes	85
5.2.1	Default Events	86
5.2.2	Payoff Events	87
5.2.3	Loss Rates Given Default	87
5.2.4	Exposure Conversion Measures	88
5.3	Default Engineering	88
5.3.1	Time-Vintage-Age (TVA) Analysis	88
5.3.2	Multi-Lead Analysis	91
5.3.3	Multi-Period Analysis	92
5.4	LGD Engineering	94
5.4.1	Data Preliminaries	95
5.4.2	Resolution Period	95
5.4.3	Risk-Neutral LGD from Observed Workout Cash Flows	96
5.4.4	LGD Winsorizing	98
5.4.5	Indirect Workout Costs	99
5.4.6	LGD Discount Rates	99
5.4.7	Observed LGD	101
5.4.8	Resolution Bias	103
5.5	Sandbox Problems	108
6	Feature Engineering	109
6.1	Synopsis	109
6.2	Missing Feature Analysis	109
6.2.1	Option 1: Keeping Missing Values	110
6.2.2	Option 2: Deleting Missing Values	111
6.2.3	Option 3: Imputation of Missing Values	111

6.2.4	Which Option Should You Consider?	114
6.3	Feature Outlier Analysis	114
6.4	Scaling	116
6.4.1	Feature Ratios	116
6.4.2	Feature Scaling	116
6.5	Non-linear Feature Transformations	118
6.5.1	Option 1: Polynomials	119
6.5.2	Option 2: Splines	120
6.5.3	Option 3: Categorization	121
6.5.4	Option 4: Weight-of-Evidence	122
6.5.5	Impact of Transformations on Fitted Default Rates	123
6.6	Feature Reduction by Aggregation	125
6.7	Feature Reduction by Clustering	128
6.7.1	Distance Measures	128
6.7.2	K-Means Clustering	129
6.7.3	Hierarchical Clustering	132
6.8	Feature Reduction by Principal Component Analysis	139
6.9	Function <code>dataprep</code>	142
6.10	Sandbox Problems	143
7	Feature Selection	144
7.1	Synopsis	144
7.2	Economic Feature Selection	144
7.2.1	Liquidity and Equity	145
7.2.2	Time, Vintage and Age	150
7.2.3	Environment	151
7.2.4	Function <code>dataprep</code>	151
7.3	Univariate Feature Selection	152
7.3.1	Means Test	152
7.3.2	F-Statistic	153
7.3.3	Association	154
7.3.4	WOE Correlations and Information Value	155
7.4	Model-based Feature Selection	157
7.4.1	Manual Selection	157
7.4.2	In (1) and Out (0) Selection	157
7.4.3	Regularization	159
7.5	Synthesis	162
7.6	Sandbox Questions	162
8	Validation	164
8.1	Synopsis	164
8.2	Qualitative Validation	164
8.3	Quantitative Validation	165
8.3.1	Basics	165
8.3.2	Backtesting as Part of Validation	167
8.3.3	Traffic Light Approach	168
8.4	Metrics for Discriminatory Power	168
8.4.1	Confusion Matrix	169
8.4.2	Accuracy Metrics	170
8.4.3	Classification Report	171

8.4.4	ROC Curve	172
8.4.5	Portfolio Dependence	177
8.5	Metrics for Calibration	180
8.5.1	Brier Score	180
8.5.2	OLS R-Squared	181
8.5.3	scikit-learn R-Squared	181
8.5.4	Binomial Test	182
8.5.5	Jeffrey's Prior Test	184
8.5.6	Calibration Curve	186
8.5.7	Hosmer-Lemeshow	187
8.6	Metrics for Stability	188
8.7	Function validation	189
8.8	Other Outcomes	191
8.9	Validation Study	191
8.9.1	Data Preparation and Feature Engineering	191
8.9.2	Fitting of Candidate Models and Validation	192
8.9.3	Comparing ROC Curves Out-of-Time	196
8.9.4	Model Stability	197
8.9.5	Practical Recommendations	199
8.10	Sandbox Problems	200

III Default, Payoff, LGD and EAD Modeling 201

9	Default Modeling	202
9.1	Synopsis	202
9.2	Data Preparation	202
9.3	Random Credit Defaults and Expectations	202
9.4	Default Models for Probability of Defaults	204
9.4.1	Link Function	204
9.4.2	Logit/Probit Models and GLMs	206
9.4.3	GLM/Logit Model	207
9.4.4	GLM/Probit Model	208
9.4.5	Comparison GLM/Logit vs. GLM/Probit Model	209
9.4.6	Multivariate Interactions	211
9.4.7	Comprehensive Model	213
9.4.8	statsmodels vs. scikit-learn	215
9.5	Forecasting PDs	218
9.5.1	Training and Test Sample	218
9.5.2	Fitting for Training Sample	218
9.5.3	Crisis Prediction for Test Sample	219
9.6	Crisis PDs	221
9.6.1	Asymptotic Single Risk Factor	221
9.6.2	Point-in-Time PDs	224
9.6.3	Through-the-Cycle PDs	226
9.6.4	Basel Capital	228
9.7	Stress-Testing	229
9.7.1	EBA and FRB Stress-Testing	229
9.7.2	Scenario-Based Stress-Testing	230
9.7.3	Parameter-Based Stress-Testing	232

9.8	Low Default Portfolios	236
9.8.1	Independent Default Events	236
9.8.2	Margin of Conservatism	239
9.8.3	Dependent Defaults	240
9.8.4	Scaling to the Mean	242
9.9	Sandbox Problems	243
10	Payoff Modeling	244
10.1	Synopsis	244
10.2	Data Preparation	244
10.3	Payoff Models	244
10.3.1	IFRS 9	246
10.3.2	Selection Control	247
10.3.3	Other Models	251
10.4	Sandbox Problems	251
11	LGD Modeling	253
11.1	Synopsis	253
11.2	Data Preparation	253
11.3	Linear Regression	254
11.4	Transformed Linear Regression	258
11.5	Fractional Response Regression	260
11.6	Beta Regression	263
11.7	Forecasting LGDs	266
11.7.1	Training and Test Sample	266
11.7.2	Fitting for Training Sample	267
11.7.3	Crisis Prediction for Test Sample	268
11.8	Sandbox Problems	269
12	Exposure Modeling	270
12.1	Synopsis	270
12.2	Data Preparation	271
12.3	Non-distressed Exposures	271
12.3.1	Computation of Credit Conversion Measures	271
12.3.2	Fitting of Credit Conversion Measures	272
12.3.3	Fitting Exposures	275
12.4	Exposures at Default	276
12.4.1	Computation of Credit Conversion Measures	276
12.4.2	Fitting of Credit Conversion Measures	278
12.4.3	Fitting Exposures	280
12.5	Sandbox Problems	281
IV	Machine Learning for PD and LGD Forecasting	282
13	Standalone Techniques	283
13.1	Synopsis	283
13.2	Data Preparation	283
13.3	Logistic Regression	284
13.3.1	Classical Logistic Regression	284

13.3.2	Logistic Regression with Regularization	286
13.3.3	Logistic Regression with Regularization and Hyperparameter Tuning	288
13.4	K-Nearest Neighbors	296
13.4.1	Idea	296
13.4.2	Practical Implementation	302
13.4.3	Hyperparameter Tuning	304
13.5	Naive Bayes	306
13.5.1	Idea	306
13.5.2	Practical Implementation	307
13.6	Decision Trees	308
13.6.1	Idea	308
13.6.2	Practical Implementation	309
13.6.3	Hyperparameter Tuning	311
13.6.4	Visualization of Trees	313
13.7	Support Vector Machines	314
13.7.1	Idea	314
13.7.2	Practical Implementation	314
13.8	Synthesis	315
13.9	Sandbox Problems	316
14	Neural Networks and Deep Learning	317
14.1	Synopsis	317
14.2	Neural Networks and Deep Learning	317
14.2.1	Idea	317
14.2.2	Simple Network without Hidden Layer	318
14.2.3	Neural Network with Hidden Layers and Non-Linearity	323
14.2.4	Practical Implementation	326
14.3	Synthesis	333
14.4	Sandbox Problems	334
15	Ensemble Techniques	335
15.1	Synopsis	335
15.2	Data Preparation	335
15.3	Bagging	336
15.3.1	Idea	336
15.3.2	Practical Implementation	337
15.4	Boosting	338
15.4.1	Idea	338
15.4.2	Practical Implementation	339
15.5	Random Forests	340
15.5.1	Idea	340
15.5.2	Practical Implementation	340
15.5.3	Hyperparameter Tuning	341
15.6	Boosted Trees	343
15.6.1	Summary of Approaches	343
15.6.2	Adaptive Boosted Trees	344
15.6.3	Stochastic Gradient Boosting	345
15.6.4	Light GBM	346
15.7	Voting Classifier	347
15.8	Synthesis	348

15.9	Sandbox Problems	349
16	Machine Learning for LGD	350
16.1	Synopsis	350
16.2	Data Preparation	350
16.3	Regression	351
16.3.1	Linear Regression	351
16.3.2	Regression with Regularization	355
16.4	K-Nearest Neighbors	361
16.5	Decision Trees	364
16.6	Random Forests	368
16.7	Boosted Trees	371
16.7.1	Adaptive Boosted Trees	371
16.7.2	Light GBM	372
16.8	Support Vector Machine	373
16.9	Neural Networks	374
16.10	Voting Regressor	378
16.11	Synthesis	380
16.12	Sandbox Problems	381
V	Synthesis: Lifetime Modeling, IFRS 9/CECL, Loan Pricing and Credit Portfolio Risk	382
17	Multi-period Modeling	383
17.1	Synopsis	383
17.2	Outcomes by Age	383
17.3	Roll Rate Analysis	385
17.3.1	Rating Classification Criteria	385
17.3.2	Rating Class Formation	386
17.3.3	Single-Period Rating Migrations	388
17.3.4	Multi-Period Rating Migrations	390
17.3.5	Multi-Period Cumulative Default Rates	390
17.3.6	Multi-Period Marginal Default Rates	391
17.4	Age Feature Models	393
17.4.1	PDs by Age	393
17.4.2	Multi-Period Forecasting of PDs	395
17.5	Survival Time Models	398
17.5.1	Probability Density Function, Survival Probability and Hazard Rate	398
17.5.2	Cross-Sectional Dataset	399
17.5.3	Cox Proportional Hazard Model	400
17.6	Other Risk Measures	407
17.7	Sandbox Problems	407
18	Expected Credit Losses	408
18.1	Synopsis	408
18.2	Expected Loss Concepts	408
18.2.1	Basel vs. CECL/IFRS 9	408
18.2.2	One-Period Expected Losses	409

18.2.3	Lifetime Expected Losses	409
18.3	Credit Risk Modeling of Age and Time	410
18.3.1	Default Rates by Age and Time	411
18.3.2	PDs by Age and Time	411
18.4	Multi-period Forecasting of Time-Varying Features	414
18.4.1	Time-Varying Features	414
18.4.2	Time Series Tests	414
18.4.3	Vector Autoregressions	415
18.4.4	Model Fitting	417
18.4.5	Multi-period Forecasting	417
18.5	Multi-period Forecasting of PDs	420
18.6	Computing Expected Lifetime Loss	423
18.6.1	Expected Losses	425
18.6.2	Expected Lifetime Losses	425
18.7	IFRS 9 Significant Increase in Credit Risk	425
18.8	Loan Pricing and Other Economic Models	428
18.9	Sandbox Problems	432
19	Unexpected Credit Losses	434
19.1	Synopsis	434
19.2	Unexpected Loss or Credit-Value-at-Risk	434
19.3	Basel Calibrations	435
19.4	Asset Correlation	439
19.4.1	Probit-Linear Model without Features	439
19.4.2	Probit-Linear Regression with Features	443
19.5	Credit Portfolio Loss Distributions	445
19.5.1	Infinitely Granular Portfolio	445
19.5.2	Limited Granularity: Numerical Integration	448
19.5.3	Limited Granularity: Monte-Carlo Simulation	451
19.5.4	Comparison of Approaches	454
19.6	Applications	455
19.6.1	Expected Loss	456
19.6.2	Value-at-Risk	456
19.6.3	Expected Shortfall	457
19.7	Sandbox Problems	459
20	Outlook	460
20.1	Where Do We Stand Today?	460
20.2	Roles of Machines	460
20.3	Where Next?	461
20.4	Keep in Touch	461
	About the Authors	463
	Bibliography	464

Part I

Principles of Data Learning

1 Deep Dive

1.1 Boost Your Practical Machine Learning Skills

“Deep Credit Risk — Machine Learning in Python” aims at starters and pros alike to enable you to:

- Understand the role of liquidity, equity and many other key banking features;
- Engineer and select features;
- Predict defaults, payoffs, loss rates and exposures;
- Predict downturn and crisis outcomes using pre-crisis features;
- Understand the implications of COVID-19;
- Apply innovative sampling techniques for model training and validation;
- Deep learn from Logit Classifiers to Random Forests and Neural Networks
- Do unsupervised Clustering, Principal Components and Bayesian Techniques;
- Build multi-period models for CECL, IFRS 9 and CCAR;
- Build credit portfolio correlation models for value-at-risk and expected shortfall; and
- Run over 1,500 lines of pandas, statsmodels and scikit-learn Python code
- Access real credit data and much more ...

This book has five parts. In the first part, “Principles of Data Learning”, we study the single mortgage dataset underlying this book, refresh our Python skills and compare Risk-based and Machine Learning. In the second part, “Data Processing and Validation”, we develop the basic components and cover Output Engineering, Feature Engineering, Feature Selection and Validation. In the third part, “Default, Payoff, LGD and EAD Modeling”, we develop a complete suite of statistical techniques, including models for probabilities of default and payoff (PD and PP), loss rates given default (LGD) and loan exposures at default (EAD). In the fourth part, “Machine Learning for PD and LGD Forecasting”, we develop a suite of machine learning (ML) concepts, which have the advantage of greater machine processing and hence, greater scalability and flexibility. We cover K Nearest Neighbor, Naive Bayes, Bagging and Boosting, Trees including Random Forests, Neural Networks and Deep Learning including Support Vector Machines. In the fifth part, “Synthesis: Lifetime Modeling, IFRS 9/CECL, Loan Pricing and Credit Portfolio Risk”, we bring everything together and calculate expected losses for loan loss provisioning and loan pricing as well as unexpected credit losses for economic and Basel capital.

We apply rigorous standards as we train methodologies using pre-crisis information and test/validate models with crisis and post-crisis information.

This may be the first comprehensive book in the area of credit risk analytics using Python. Python is an up-and-coming programming language that has become popular in many areas of data analytics in a very short time. Python is popular as it is free to use and offers some advantages for Machine Learning. We start with a first introduction to Python in the next chapter.

This book is aimed at credit analysts in financial institutions, fin-techs and prudential regulators, as well as credit risk researchers. Data, codes and more are available on the accompanying website: www.deepcreditrisk.com.

1.2 Credit Risk Information

Commercial banks keep data in a number of repositories that correspond to steps in the lending process: origination data, performance data, modification data payoff/retention data, default data/workout data and maturity data. These sources may include:

- Origination/underwriting data: relates to the origination or underwriting of loans;
- Performance data: relates to multiple reviews, which are generally conducted monthly, quarterly or annually and cover the period from the origination time to the latest review;
- Modification data: relates to the modification time;
- Payoff/retention data: relates to the payoff time;
- Maturity data: relates to the maturity time and covers administrative matters such as release of collateral and various accounting activities;
- Default/workout data: relates to the period from default to resolution and collects all cash flows that are paid or received by the bank during this time. Times are continuous and resolution periods can cover multiple periods (e.g., up to ten years).

External data may include:

- Macroeconomic information: time varying information that is identical at a given period for a number of borrowers. Macroeconomic information may be stratified by country, state, statistical areas and other units;
- Population statistics: these are generally time-varying;
- There are many other data sources which may include business filings, data from social networks, data from external experts such as ratings agencies, property appraisers, activity profiles of payment systems or transport systems.

The general aim is to construct a panel data. For the panel data, we observe information features and risk outcomes like default, payoff, loss rates and exposures for loans ($i = 1, \dots, l$) and time periods ($t = 1, \dots, T$). Generally speaking, there are a number of steps needed before we can process data.

First, raw data needs to be measured. Credit risk models are developed for scaling: estimated models are used to measure the risk of new borrowers and loans as well as borrowers and loans for which the risk factors have changed since the last review. Any information that is used in a model should be observable and hence measurable for the borrowers and loans. For the data sources, this means that the sensor that has been used in the past to collect the raw information should be available to collect the same information going forward. The sensor information is then transferred to numbers or strings, collected in an array and processed. Should the sensing or the transformation to numbers or strings fail, values will be missing and have to be treated in a second stage.

Second, identifiers that are common across number arrays and other credit relevant information need to be created. Different identifiers may exist. Two very important ones are the ultimate credit

subject (examples are borrowers or loans) and time. Examples for time in lending are application time, origination time, observation time, payoff time, default time and maturity time.

In identifying the ultimate credit subject, one may consider the loan or the borrowers and the question is whether the risk should be measured at the loan-level or the borrower-level. The answers to these questions are often bank-specific and an economic analysis is required for most consumer and corporate loans. Things get more complicated for a few high-net-worth individuals and large firms as these maintain many accounts and credit-risk relevant relationships with other subjects. Even the definition of borrower is not trivial as the borrower may be part of a larger holding structure, family or benefit from guarantees and an analysis at the parent level may be sensible.

Third, the resulting array of credit information is converted into a dataframe in a format that allocates rows as combinations of ultimate credit subject and time. Time is generally observation time but in special applications, such as building score cards for loan applications, it may also be loan origination time. Information is analyzed by computing relations. These relations can be reciprocal or one-way. In most of this book we analyze the impact of independent variables (also known as covariates, risk factors, explanatory variables, features and right-hand side variables) on dependent variables (also known as responses, outputs, outcomes and left-hand side variables).

1.3 Hands-on Analysis

We are working with one comprehensive dataset on US mortgages. The data is collected and provided by International Financial Research (see www.internationalfinancialresearch.org). The dataset is a randomized selection of mortgage loan-level data collected from the portfolios underlying US RMBS securitization portfolios.

We are further working with module `dcr.py`, which contains a number of credit risk functions that are repetitively used throughout this book. These functions make a great addition to any Python library.

To get started:

- Download data `dcr.csv` from www.deepcreditrisk.com;
- Download module `dcr.py` from www.deepcreditrisk.com;
- Place data and file in directory `'C:/TMP'`. You can save both files in a different directory but need to then change all paths in `dcr.py` from `'C:/TMP'` to your preferred path;
- Save all codes (e.g., Jupyter Notebooks) in the same directory.

The dataset is in panel form, reporting origination and performance observations for 5,000 residential US mortgage borrowers over 60 periods. The central variables are `id` and `time`. The data size was chosen to allow for efficient computing of a number of machine learning techniques. A larger dataset is available at www.deepcreditrisk.com.

The quarterly periods have been de-identified and observation periods run from 1 to 60. The observation period starts at the beginning of the millennium and includes the Global Financial Crisis (GFC) in period 27 approximately. Origination times prior to the start of the observation period have negative numbers. For example, `orig_time=-2` implies an origination time of two period prior to the observation window. Within this dataset we observe business cycles with economic upturns and downturns.

Outcome observations are observed over a time period after features are recorded. For example, default, payoff and status events are observed one period after features in the same row. LGD and related recoveries are observed between the default and resolution time.

As in the real world, loans may originate before the start of the observation period. This is possible if loans are transferred between banks and investors as in securitization. The loan observations may thus be censored as the loans mature or borrowers refinance.

The information is organized in the following order:

- Borrower IDs;
- Time stamps;
- Information features at observation time;
- Information features at loan origination;
- Outcome observations.

More specifically, key variables are:

- `id`: borrower id;
- `time`: time stamp of observation;
- `orig_time`: time stamp for origination;
- `first_time`: time stamp for first observation;
- `mat_time`: time stamp for maturity;
- `res_time`: time stamp for resolution;
- `balance_time`: outstanding balance at observation time;
- `LTV_time`: loan to value ratio at observation time, in %;
- `interest_rate_time`: interest rate at observation time, in %;
- `rate_time`: risk-free rate, in %;
- `hpi_time`: house price index at observation time, base year=100;
- `gdp_time`: GDP growth at observation time, in %;
- `uer_time`: unemployment rate at observation time, in %;
- `REtype_CO_orig_time`: real estate type — condominium: 1, otherwise: 0;
- `REtype_PU_orig_time`: real estate type — planned urban developments: 1, otherwise: 0;
- `REtype_SF_orig_time`: real estate type — single family home: 1, otherwise: 0;
- `investor_orig_time`: investor borrower: 1, otherwise: 0;
- `balance_orig_time`: outstanding balance at origination time;
- `FICO_orig_time`: FICO score at origination time, in %;
- `LTV_orig_time`: loan to value ratio at origination time, in %;
- `interest_rate_orig_time`: interest rate at origination time, in %;
- `state_orig_time`: US state in which the property is located;
- `hpi_orig_time`: house price index at observation time, base year=100;
- `default_time`: default observation at observation time;
- `payoff_time`: payoff observation at observation time;
- `status_time`: default (1), payoff (2) and non-default/non-payoff (0) observation at observation time;
- `lgd_time`: LGD, at default time, assuming no discounting of cash flows;
- `recovery_res`: sum of all cash flows received during resolution period.

Python language is case-sensitive. Please ensure that you use the identical upper- and lowercase character as indicated. We make use of all variables, the most important of which used in this book are `default_time` (outcome variable) as well as `LTV_time`, `FICO_orig_time` and `gdp_time` (feature variables).

A central variable is the loan to value ratio (LTV). The data includes LTV at loan origination time (`LTV_orig_time`) and at observation time (`LTV_time`). The two ratios are related. `LTV_time` includes `LTV_orig_time`, the change in house prices (measured by the house price index, HPI) and amortization measured by the loan balance at loan origination (`balance_orig_time`) and at observation time (`balance_time`). The formula for the calculation is:

$$LTV_time = \frac{\text{balance_time}}{\text{house price at time}}$$

The house price at time can be approximated as follows:

$$\text{house price at time} = \text{house price at origination} * \frac{\text{hpi_time}}{\text{hpi_orig_time}}$$

With:

$$\text{house price at origination} = \frac{\text{balance_orig_time}}{LTV_orig_time}$$

We later convert the LTV ratio to home equity for economic discussions. Most variables are generally observed for all observations. Exceptions are `lqd_time`, `recovery_res` and `res_time` which are only observed for `default_time=1` and if the resolution process is complete.

1.4 Basel, CECL, IFRS 9, DFAST, CCAR and Stress Tests

Deep credit risk analytics is the basis of a number of regulations for financial institutions that all have an impact on the amount of capital an institution holds. Important regulations are:

- Basel: minimum amount of required Tier I and Tier II capital. Basel may include the various reforms (Basel I to Basel III), and a number of nationally issued guidance notes;
- Current Expected Credit Loss (CECL), IFRS 9: loan loss provisioning and eligible amount of available Tier I capital;
- National stress tests (e.g., Dodd-Frank Act Stress Test (DFAST) or Federal Reserve Bank (FRB) stress tests in the US): requirement of additional capital buffers;
- Comprehensive Capital Analysis and Review (CCAR): requirement of additional capital buffers.

All frameworks are related and large banks require the estimation of PDs, LGDs and EADs. However, the estimation details may differ. Basel requires through-the-cycle PDs, Downturn EADs and Downturn LGDs. CCAR, CECL and IFRS 9 require Lifetime PDs and EADs/LGDs that are based on current economic circumstances and are forward-looking, i.e., take future expectations into account. DFAST and FRB stress tests require stressed PDs, LGDs and EADs.

It is possible to meet all requirements within the frameworks presented in this book. Further, financial institutions may compute the minimum amount of economic capital. Economic capital is important

for setting loan prices and provide incentives that are aligned with stakeholder values. Banks that work off credit prices that are too low attract risks that are too high through adverse selection. Bank that use prices that are too high lose market share and are unable to offset fixed costs.

1.5 First Lessons from the COVID-19 Crisis

1.5.1 Credit Risk Analytics

At the time of writing this book the COVID-19 Crisis has started and it was evident that there will be a long-lasting impact on many areas of human life including bank lending and credit risk. There are a number of aspects that are important for credit risk analytics and we believe that we cover some in this book.

The crisis has led to a sudden shock to employment and incomes that is unprecedented in most banks' data including the Global Financial Crisis (GFC). Many employees have been made redundant, stood down (i.e., kept employed at zero income), or lost a part of their entitlements. The follow-on effects were the inability of private and corporate borrowers to make rent payments and private and corporate investors to service loans. Governments have tried to offset liquidity constraints by distributing money (sometimes called helicopter money), lowering taxes and offering liquidity facilities. We discuss liquidity to some degree in this book by measuring feature `cep_time` (see the chapter on feature selection). There are many aspects and proxies for liquidity and banks often consider working capital ratios, debt to income ratios and more sophisticated measures from transaction data.

Equity next to liquidity is a central aspect. The COVID-19 crisis had at the beginning limited impact on asset and hence wealth values. Many believe that there will be also a shock to asset markets. We discuss equity in this book using feature `equity_time`.

Time effects are very relevant, this includes the origination time, age and the state of the economy. There is most "meat on the bone" in credit risk models to carefully collect and interpret information. For example, the pandemic is global as many systems that we have created. The impact on the travel industry, global supply chains for many labor, finance, commodities, manufactured goods and service will become key features in explaining credit risk outcomes yet most models built are poor on this end. Many identified features are missing in today's models and can be added to the models.

That said the correct use of methodologies is key from a model design perspective. For example, models built during economic upturns struggle predicting credit risk outcomes during downturns.

In this book, we try to provide and showcase solutions for many questions. The following table shows potential solutions to some challenges that may come from the COVID-19 crisis impact:

Challenge	Solution	Chapter
Calculating Crisis PDs without downturn data	Model-based measurement of crisis PDs, Parameter-based stress-testing (Margin of conservatism, Bayesian approach)	Chapter 9 Default Modeling — Crisis PDs/Stress-testing, Chapter 3 Risk-based Modeling — Bayesian Approaches

Challenge	Solution	Chapter
Calculating Crisis PDs with downturn data	Scenario-based stress-testing, Parameter-based stress-testing (Regime-switching models)	Chapter 9 Default Modeling — Stress-testing
Liquidity as a driver of default	Estimation of models with liquidity as feature; Inclusion of additional liquidity feature (e.g., income over non-discretionary expenses)	Chapter 5 Outcome Engineering, Chapter 9 Default Modeling
Equity as a driver of default	Estimation of models with equity as feature; Inclusion of additional equity feature (e.g., home improvement/deterioration)	Chapter 5 Outcome Engineering, Chapter 9 Default Modeling
Impact of time effects	TVA Analysis: control for vintage and age effects through dummy variables or other features that describe the origination process and for time effects through macroeconomic features	Chapter 5 Outcome Engineering, Chapter 9 Default Modeling
Low default portfolios	Most prudent estimators/Margin of conservatism	Chapter 9 Default Modeling — Low default portfolios
Validation of pre-crisis models	Backtesting: split training and validation sample along time dimension	Chapter 8 Validation
Ability of machine learning models to predict defaults for severe downturns	Backtesting of machine learning approaches	Chapter 13 Standalone Techniques, Chapter 14 Neural Networks and Deep Learning, Chapter 15 Ensemble Techniques, Chapter 16 Machine Learning for LGD
Adequacy of model estimates for application	Comparison of Basel capital with expected loss	Chapter 9 Default Modeling — Crisis PDs

Our approach includes four main elements. First, model calibrations should include the latest credit data. Second, features that are identified as drivers of credit risk outcomes should be included in the models. For example liquidity, equity and network effects should be added to the feature space. Third, validation should focus on backtesting and include a train-test split along the time line. This is known as backtesting. Fourth, the adequacy of model estimates for applications like capital adequacy, loan loss

provisioning and loan pricing needs to be vetted.

1.5.2 Machine Learning

Fintech is short for “financial technology” and is a worldwide innovation boom that has created multi-billion dollar industries in many countries. Fintech is seen as the future of financial services and challenged existing financial service providers like banks, insurers and pension funds by lower variable costs. The services are generally offered digitally and online via mobile apps. Back-end operations, including credit risk analytics are often to a higher degree automatized. This industry has been at the forefront of applying advanced technique in credit risk analytics

The outbreak of the pandemic COVID-19 led to an unprecedented decline in many financial markets globally. The pandemic was to a large degree unexpected but two issues have become already apparent. First, fintech firms have in many cases lost more market value than existing banks as they have issued loans to higher risk borrowers and often at lower lending standards. Second, advanced econometric techniques had a limited ability to predict the crisis.

In defense of the fintech industry, we would like to stress that the growth — like with many other innovations — was driven by investor over-excitement, which perhaps led to lending standards that were below model standards. Further, credit risk generally realizes in time-lags and we have not seen the full impact of the pandemic on credit risk. Traditionally, economic downturns result in unemployment, a lower wealth and liquidity of borrowers declining housing and other collateral values. The outcomes are greater future likelihoods and magnitudes of credit losses. Whilst advanced models may have failed in an immediate crisis prediction, they may adjust quicker to the new risk levels than traditional models.

It is possible that an upcoming credit crisis will be interpreted as a failure of fintechs and advanced credit risk approaches. We take a neutral position in this book and judge all techniques based on pandemic levels. We rely on the Global Financial Crisis, as a severe economic shock, and analyze the performance of these advanced approaches as if they had been built prior to the crisis and use-tested in the aftermath.

The fintech and banking industry are under great stress and the industry is likely to transform as a consequence. However, we believe that new business models will continue to challenge existing lenders by using new technology. This book will help to provide a common understanding and enable credit risk analysts to benefit from opportunities that will come with it.