

基于多元Logit模型的小微企业信用评级分类研究

赵星然, 邓晓卫, 吕学斌

南京工业大学, 江苏 南京

收稿日期: 2022年9月14日; 录用日期: 2022年9月26日; 发布日期: 2022年10月31日

摘要

首先用Python语言对123家小微企业共约三十七万条交易信息进行分类整理; 然后用主成分分析将多种信息归为四类主要因子, 并引入定性变量; 最后建立多元Logit模型对企业进行由好至差的A、B、C、D四类评级。结果显示: 整体评级准确率达到72%; 交易信息离评级时点越近, 评级的准确率越高; 该模型对D级的评级准确率达到100%, 说明该方法能有效甄别最差级别企业, 为商业银行规避不良贷款发生提供了一种可行性方法。

关键词

信用评级, Python语言, 主成分分析, 多元Logit模型

Research on Credit Rating Classification of Small and Micro Enterprises Based on Multinomial Logit Model

Xingran Zhao, Xiaowei Deng, Xuebin Lv

Nanjing University of Technology, Nanjing Jiangsu

Received: Sep. 14th, 2022; accepted: Sep. 26th, 2022; published: Oct. 31st, 2022

Abstract

First, use Python to classify and sort a total of about 370,000 transaction information from 123 small and micro enterprises; then use principal component analysis to classify a variety of information into four main factors, and introduce qualitative variables; finally, establish Multinomial

Logit model to rank companies from good to bad in four categories: A, B, C and D. The result shows: the overall rating accuracy rate reaches 72%; the closer the transaction information is to the rating time point, the higher the accuracy of the rating; the accuracy rate of the D-level rating of this model reaches 100%, it shows that this method can effectively identify the worst-level enterprises, which provides a feasible method for commercial banks to avoid the occurrence of non-performing loans.

Keywords

Credit Rating, Python, Principal Component Analysis, Multinomial Logit Model

Copyright © 2022 by author(s) and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

1. 引言

小微企业是一国经济发展的基础性力量，它们不仅提供大量就业岗位，而且为经济注入活力，在保持社会稳定方面发挥着巨大作用。特别，美、日等发达国家的经验表明，许多今天的巨型企业都是由当初的小微企业发展起来的。但是小微企业由于成立时间短、几无固定抵押资产、缺少资信等特点，在小微企业的成长过程中一直存在融资难问题。近年来，国家出台了一系列扶持小微企业发展的金融政策，引导商业银行等金融机构为小微企业提供融资支持小微企业发展。金融机构一方面要承担起为小微企业融资，支持小微企业发展的社会责任，但另一方面，又要尽可能的减少小微企业的违约风险，以规避金融系统的借贷风险。因此，如何根据能获得的小微企业信息，给企业进行评级，进而针对不同信用级别的企业制定合理的贷款额度和贷款利息是商业银行十分关注的重要问题。

企业信用等级作为商业银行风险管理中重要的组成部分，为信贷决策提供了重要的依据。关于如何对企业进行信用评级以及哪些因素可能会引起企业违约风险，此前已有大量研究。在传统分类模型方面，Altman 等(2013) [1]认为从统计学的角度考虑，Logit 模型回归似乎更适合于解决违约预测的问题，在因变量是二进制(是/非)的条件下，组别都是离散的、非重叠的和可识别的；孙雨忱(2021) [2]通过构建二元 Logistic 回归的违约率测算模型，测算出不同信用水平下中小企业的预期违约率；高璐冰等(2021) [3]使用熵权 TOPSIS 法计算每个企业的得分，量化信贷风险并确定信用评级，建立二元 Logit 回归模型确定企业违约概率。在机器学习模型方面，乔薇(2011) [4]选取了流动资产周转率、净资产收益率等指标建立了 AHP 模型对中小企业进行多级模糊综合评级；徐晓萍和马文杰(2011) [5]运用判别分析法和决策树模型，对非上市中小企业的违约风险进行了分析，发现二者结合能较好地判断企业违约率；郭妍等(2013) [6]建立二元 Logit 回归模型和 LDA 模型对我国中小企业信贷风险进行实证度量，对企业是否违约进行预测；夏利宇等(2019) [7]将 XGBoost 算法与 Logistic Group Lasso 模型相结合建立了企业是否违约二分类模型，评估了客户违约风险；白羽等(2021) [8]通过建立神经网络模型预测企业未违约率，量化分析了企业信贷风险。

综上，目前的研究主要方向是用二元 Logit 模型和其它包括机器学习方法对企业是否违约的二分类情形进行预测研究。但银行要发放贷款，并且要实施有差异的贷款利率，就需要获得准确的企业信用等级，以按不同等级发放贷款额度和确定贷款利率。目前，对企业进行等级评价这方面的研究较少，特别是对几无固定资产、无担保的小微企业的信用评级研究还未见有。如何针对小微企业所提供的有限信息，通过

数据处理获取相关特征，对企业进行等级评定，进而为商业银行制定贷款策略提供依据，正是本文研究的主题。

本文研究这样一类小微企业：它们无固定资产，也无信誉担保，在向银行提出贷款时，唯一能提供的真实信息是它们一段时间的大量的“进、销项”财务记录。本文首先用 Python 语言对初始数据进行处理、分类，提取有用信息；然后利用因子分析法对大量、多维信息进行降维；最后，建立多元多分类 Logit 模型对企业进行信用评级(按 A、B、C、D 四类评级)，为商业银行更加客观、准确地评估小微企业的信用状况，制定合理贷款额度及贷款利率提供有效参考。

2. 数据处理及变量选取

数据来源于某银行获得的 123 家有信贷记录且已有评级的小微企业的相关经营信息¹，包括：企业进项和销项发票信息、企业类型、企业信用评级等数据，其中进项发票数据共 210,948 条，销项发票数据总共 162,485 条。时间 2017 年 1 月至 2020 年 2 月。

2.1. 数据预处理

首先，运用 Python 语言对 123 家经营信息共约 38 万条原始发票数据进行整合处理²。发票所含直接信息为：企业名称、企业信誉评级、企业是否违约、进项发票金额及税额、销项发票金额及税额、进/销项作废发表、发票号码、开票日期、购/销方单位代号等。引入“有效收入”，“作废发票率”及“负数发票率”定义如下：

$$\text{有效收入} = \sum \text{销项有效发票价税总和} - \sum \text{进项有效发票价税总和},$$

$$\text{作废发票率} = (\text{进项作废发票数量} + \text{销项作废发票数量}) / (\text{进项发票数量} + \text{销项发票数量}),$$

$$\text{负数发票率} = (\text{进项负数发票数量} + \text{销项负数发票数量}) / (\text{进项发票数量} + \text{销项发票数量}).$$

经 Python 语言处理将发票信息归结为以下几个方面：1) 总发票、有效发票、负数发票和作废发票的数量；2) 有效金额，有效税额，有效价税总和以及作废额；3) 进项和销项发票总数；4) 有效收入、作废发票率、负数发票率等。

2.2. 主成分因子提取

由预处理共计获得 19 个变量，它们是 x_1 ：进项发票数量、 x_2 ：进项有效发票数量、 x_3 ：进项作废发票数量、 x_4 ：进项负数发票数量、 x_5 ：进项有效金额、 x_6 ：进项有效税额、 x_7 ：进项无效额、 x_8 ：进项有效价税总和、 x_9 ：销项发票数量、 x_{10} ：销项有效发票数量、 x_{11} ：销项作废发票数量、 x_{12} ：销项负数发票数量、 x_{13} ：销项有效金额、 x_{14} ：销项有效税额、 x_{15} ：销项无效额、 x_{16} ：销项有效价税总和、 x_{17} ：有效收入、 x_{18} ：作废发票率、 x_{19} ：负数发票率。由于变量过多，且一些变量可能存在较强的相关性，不能直接用于评级分析，故进行变量整合，拟采用因子分析法对预处理得到的特征变量进行降维。因子分析使用统计软件 SPSS 24。

首先考察原有变量是否适合进行因子分析，利用 KMO 指标和 Bartlett 检验进行。KMO 指标主要作用是观测相关系数值和偏相关系数值，Bartlett 检验用于检验各变量是否独立。KMO 和 Bartlett 检验结果如表 1 所示。

¹ 数据来源：http://www.mcm.edu.cn/html_cn/node/10405905647c52abfd6377c0311632b5.html。

² 由于代码篇幅过长，此处省略，若有需要可联系作者。

Table 1. KMO and Bartlett test
表 1. KMO 和巴特利特检验

KMO 取样适切性量数		0.681
巴特利特球形度检验	近似卡方	4480.903
	自由度	91
	显著性	0.000

由表 1 可得, $KMO = 0.681$, 根据专家给出的标准, $KMO > 0.5$, 说明所选取的原始变量之间存在强相关性。巴特利特球形度检验原假设为相关系数矩阵是一个单位阵, 变量无相关性, 由结果得显著性为 0.000, 所以该相关阵各变量之间存在相关性, 说明可以做因子分析。

因子提取方法选用主成分分析法, 首先计算各因子的特征根, 结果显示: 特征根 > 1 的因子总共 4 个, 且该 4 个因子已经可以包含原始变量 86.225% 的信息, 所以提取 4 个公因子。提取的四个主成分因子如下:

第一主成分因子(F_1)在指标进项有效金额、进项有效税额、进项无效额、进项有效价税总和、销项有效金额、销项有效税额、销项有效价税总和有较大载荷, 这些指标主要反映了这些企业的有效资金情况, 所以定义第一主成分因子为“金额因子”; 第二主成分因子(F_2)在指标进项发票数量、进项有效发票数量、进项作废发票数量、销项作废发票数量有较大载荷, 这些指标主要是反映了这些小微企业的各项发票数量, 故定义第二主成分因子为“活力因子”; 第三主成分因子(F_3)在指标销项负数发票数量、负数发票率有较大载荷, 这些指标主要是反映了企业入账记税后购方因故发生退货并退款的情况, 故定义第三主成分因子为“故障因子”; 第四主成分因子(F_4)在指标有效收入有较大载荷, 主要是反映了这些小微企业收入情况, 定义第四主成分因子为“收入因子”³。

通过因子分析对指标降维后, 最终确定了上述 4 个评价指标, 分别是 F_1 : 金额因子, F_2 : 活力因子, F_3 : 故障因子, F_4 : 收入因子。将分析这些变量对企业评级产生的影响。

2.3. 其它变量确定

为保证研究小微企业评级问题的全面性, 除前述经因子分析获得的 4 个主成分因子变量 $F_1 \sim F_4$ 外, 本文还根据 123 家企业特点, 将企业进行了分类, 引入定性变量 F_5 : 企业类别; 定性变量 F_6 : 企业是否违约; 衡量企业发展是否稳定, 引入变量 F_7 : 进稳定度, F_8 : 销稳定度, F_9 : 进销稳定度。这些变量的具体定义如表 2 所示。本文将建立多元 Logit 模型, 通过上述定义的变量 $F_1 \sim F_9$, 研究它们对企业评级的影响。并由此根据获得的权重再对企业进行评级。按照行业惯例, 企业一般按好 \rightarrow 差被评为 A、B、C、D 四级。

3. 模型建立及估计

3.1. 多元 Logit 模型设定

多元 Logit 模型可视为被解释变量任意选定基准组后将其他组别与基准组分别配对构成的多个二元 Logit 模型实施联合估计。模型设定具体如下:

$$G_j = \ln \left(\frac{P(y = j | x)}{P(y = b | x)} \right) = \beta_{j0} + X\beta_j = \beta_{j0} + \sum_{k=1}^K \beta_{jk} x_k \quad (1)$$

³限于篇幅, 各主成分荷载具体结果此处省略, 若有需要可联系作者。

其中, $X = [x_1, x_2, \dots, x_K]$ 是解释变量, $\beta_j = [\beta_{j1}, \dots, \beta_{jK}]^T$ 为解释变量的系数, β_{j0} 为截距项, 类别变量 y 有 $j = 1, 2, 3, \dots, J$ 个类别. b 为选定的基准组, 当 $j = b$ 时, 由于 $\ln 1 = 0$, 故 $\beta_{b0} = \beta_b = 0$, 即 $G_b = 0$. 故在 J 个类别下, 只要估计 $J-1$ 个模型即可.

通过模型估计, 可以得到每种类别选择的预测概率:

$$P(y = j | x) = \frac{e^{\beta_{j0} + \sum_{k=1}^K \beta_{jk} x_k}}{1 + \sum_{j=1}^{J-1} e^{\beta_{j0} + \sum_{k=1}^K \beta_{jk} x_k}} \quad (2)$$

Table 2. Variable definitions

表 2. 变量定义

解释变量	符号	定义
企业类别	F_5	F_{51} = 科技类, F_{52} = 贸易类, F_{53} = 管理类, F_{54} = 综合类
违约变量	F_6	F_{60} = 无违约, F_{61} = 有违约
进稳定度	F_7	$\frac{\text{企业发票数量前3名的销方单位的发票总数}}{\text{进项发票数量}}$
销稳定度	F_8	$\frac{\text{企业发票数量前3名的购方单位的发票总数}}{\text{销项发票数量}}$
进销稳定度	F_9	$\left(\frac{F_7^2 + F_8^2}{2}\right)^{\frac{1}{2}}$
信用等级(因变量)	y	$0 = D, 1 = C, 2 = B, 3 = A$

3.2. 模型估计

针对本文研究的问题共有四个类别, 即模型(1)中 $j = 0, 1, 2, 3$, 考虑到类别变量, 解释变量共十一个, 故 $K = 11$. 因变量选取 $y = 3$ (A 级企业)为基准组, 自变量中的定性变量 F_5 选取 F_{54} = 综合类为基准组, F_6 选取 F_{61} = 有违约为基准组. 企业样本数为 123 个, 估计软件为 RStudio4.0.5. 估计结果由表 3 给出.

4. 企业评级预测分析

由多元 Logit 模型估计得到的结果(表 3), 可获得模型(1)中各解释变量的权重 β_j , 再把企业的各解释变量值代入模型(1), 即可求得企业的机会比率对数值 $G_j, j = 0, 1, 2, 3$, 根据公式(2)可以得到企业的信誉评级为 A、B、C、D(即 y 取值为 3, 2, 1, 0)的概率如下.

$$P(y = 3 | x) = \frac{e^{G_3}}{e^{G_0} + e^{G_1} + e^{G_2} + e^{G_3}}$$

$$P(y = 2 | x) = \frac{e^{G_2}}{e^{G_0} + e^{G_1} + e^{G_2} + e^{G_3}}$$

$$P(y = 1 | x) = \frac{e^{G_1}}{e^{G_0} + e^{G_1} + e^{G_2} + e^{G_3}}$$

$$P(y = 0 | x) = \frac{e^{G_0}}{e^{G_0} + e^{G_1} + e^{G_2} + e^{G_3}}$$

规定: 上述预测概率最高的类别就是该企业最终被评定的类别, 由此计算得到 123 家企业的评级,

再将此结果与该企业原有评级结果进行比较,得到由多元 Logit 模型对 123 家企业评级预测的正确率,结果由表 4 给出。

从表 4 给出的评级正确率结果可以看出,基于多元 Logit 模型在企业评级预测中,两级(A、D 级)分辨的正确率较高。企业信用评级为 D 的预测准确率达到 100.0%,表明仅利用从企业获得的进、销项信息,通过建立多元 Logit 模型可以有效甄别较差企业,从而可以大概率地为银行规避向这类企业发放贷款的风险。在企业信用评级为 A 上分类准确率达到 70.4%,说明对于信用情况好的企业,该模型也能够较为准确地识别。模型总体分类准确率为 66.7%,模型对总体变异的解释能力一般。

Table 3. Multivariate Logit model estimation result table of enterprise rating

表 3. 企业评级的多元 Logit 模型估计结果表

变量	信誉等级					
	B (y = 2)		C (y = 1)		D (y = 0)	
	β	P值	β	P值	β	P值
F_1	-0.477	0.389	-2.024	0.489	-1262.384	0.916
F_2	-1.221*	0.091	-0.651*	0.097	31.557	0.932
F_3	-1.335	0.215	0.269	0.664	11.658	0.972
F_4	-0.680	0.411	0.648	0.544	212.271*	0.097
F_{51} = 科技类	1.397*	0.095	1.675*	0.072	21.256	0.928
F_{52} = 贸易类	1.684	0.130	1.727	0.135	-0.411	0.926
F_{53} = 管理类	1.700**	0.041	1.605*	0.063	-31.020	0.964
F_{60} = 无违约	-14.493**	0.047	-15.810***	0.000	-497.053*	0.086
F_7	-8.186	0.323	-11.514*	0.084	-725.432	0.928
F_8	0.745	0.922	-4.078	0.549	-1169.233	0.926
F_9	8.259*	0.091	18.377	0.196	1957.680	0.926

注:表中数据右上角的***, **, *表示估计的该参数分别在 1%, 5%, 10%的水平下显著。

Table 4. Predicted winning rate table based on full data enterprise rating classification

表 4. 基于全数据企业评级分类预测胜率表

实际 评级	预测评级				正确百分比
	0	1	2	3	
0 (24)	24	0	0	0	100.0%
1 (34)	0	15	14	5	44.1%
2 (38)	0	10	24	4	63.2%
3 (27)	0	2	6	19	70.4%
平均值					66.7%

注:表中括号中的数为实际评为该级别的企业数,后同。

由于小微企业受行业、环境影响较大,企业变化发展很快,所以越靠近评级时点的数据对企业评级预测应该更准确。故本文进一步只选取 2019 年 1 月~2020 年 2 月的数据进行建模预测,得各级别企业预

测正确率如表 5 所示。

由表 5 可得，当获取数据的时间段更靠近预测时间点时，模型分类的准确率有所提高，特别，D 级预测正确率仍保持 100%，其次亦是 A 级；中间两级 B、C 级的预测准确率有所提高。这个结果与相关部门对中小企业信用等级评价结果有效期统一为三年，相关协会在有效期内对企业每年进行一次复查等政策是完全吻合的。所以用于企业信誉评级的相关数据具有时效性，数据离评价时间越近，对企业评价的准确度越高。

Table 5. Predicted winning rate table based on 2019~2020 data enterprise rating classification

表 5. 基于 2019 年~2020 年数据企业评级分类预测胜率表

实测	预测				正确百分比
	0	1	2	3	
0 (24)	24	0	0	0	100.0%
1 (34)	0	20	9	5	58.8%
2 (38)	0	10	26	2	68.4%
3 (27)	0	2	6	19	70.4%
平均值					72.4%

5. 结论及政策建议

本文以 2017 年 1 月~2020 年 2 月某银行获得的 123 家有信贷记录且已有评级的小微企业营业数据(进、销项数据)为样本，以 Python 语言为工具，结合多因子分析方法，建立多元 Logit 模型研究了这一类小微企业的信用评级问题，结果显示：

1) 基于多元 Logit 模型对企业进行评级，该方法对两端企业(A、D 级)的评级准确率更高。特别，对 D 级的评级准确率达到 100%且结果稳健，这为商业银行在有限信息下甄别较差企业、有效规避不良贷款发生提供了一种有效方法。

2) 当用三年数据进行企业信用评级时，总体评级准确率为 66.7%；若改用企业最近一年的数据进行评级，准确率提高了 5.7 个百分点，达到 72.4%。这表明用于企业信誉评级的相关数据具有时效性，交易信息离评级时点越近，评级的准确率越高。这一结论与相关部门对企业评级的信息规范要求相吻合。

本文为在有限信息下针对无抵押、无担保的小微企业评级提供了一种方法。该方法虽然在 D 级评级上具有较好优势且结果稳健，但在中间级别 B、C 级的评级准确率上还有待进一步提高。由于本研究只获得 123 个小微企业数据，数据有限，因此在尝试用其它方法建模时效果欠佳。故，未来在可获得更多数据时，将尝试拓展其它建模方法，特别是近年兴起的机器学习方法用于企业评级预测研究，以获取更好预测效果，为金融机构合理放贷、规避风险提供有效策略。

基金项目

国家自然科学基金项目(11801267)。

2022 年江苏省研究生科研与实践创新计划项目(SJCX22_0410)。

参考文献

- [1] Roggi, O. and Altman, E. (2012) Managing and Measuring Risk: Emerging Global Standards and Regulations after the Financial Crisis. World Scientific, Singapore, 343-455.

-
- [2] 孙雨忱. 信息不对称下银行对中小微企业的最优信贷策略研究——基于 Logistic 回归的违约率测算模型[J]. 金融发展研究, 2021, 4(6): 78-84.
- [3] 高璐冰, 赵国庆, 侯家璇. 基于商业银行视角下中小微企业信贷决策研究[J]. 中国商论, 2021(9): 55-57.
- [4] 乔薇. 中小企业信用评级指标体系与模型的构建[J]. 开封大学学报, 2011, 25(4): 89-93.
- [5] 徐晓萍, 马文杰. 非上市中小企业贷款违约率的定量分析——基于判别分析法和决策树模型的分析[J]. 金融研究, 2011(3): 111-120.
- [6] 郭妍, 张立光, 刘佳. 中小企业信贷风险度量模型研究——基于山东省的实证分析[J]. 东岳论丛, 2013, 34(7): 58-61.
- [7] 夏利宇, 张勇, 鲁强, 汤广瑞. 结合 XGBoost 算法和 Logistic 回归的信用评级方法[J]. 征信, 2019, 37(11): 56-59.
- [8] 白羽, 三郎斯基, 王晓妍, 孙少飞, 王恒友. 基于层次分析法量化中小微企业信贷风险[J]. 北京建筑大学学报, 2021, 37(2): 93-97.