

# A Credit Scoring Technology Based on Bayes Discriminant Analysis

Shenghua Zhou

Rizhao Marine Engineering Vocational College, Rizhao Shandong  
Email: 66838031@qq.com

Received: Jul. 19<sup>th</sup>, 2019; accepted: Aug. 2<sup>nd</sup>, 2019; published: Aug. 9<sup>th</sup>, 2019

---

## Abstract

According to FICO Score theory and Bayes Discrimination, the credit scoring model is derived, which ends up as an optimization model with linear objective function and quadratic equality constraints. Finally, compared with Logistic Regression through an example, the result shows that the credit scoring model is effective and can support application scenarios of practical business better.

## Keywords

Credit Scoring, Bayes Discriminant Analysis, Divergence Model, Binning Weight, Score Scaling

---

# 一种基于贝叶斯判别的信用评分方法

周声华

日照航海工程职业学院, 山东 日照  
Email: 66838031@qq.com

收稿日期: 2019年7月19日; 录用日期: 2019年8月2日; 发布日期: 2019年8月9日

---

## 摘 要

本文借鉴了FICO评分的思想, 基于贝叶斯判别定理推导出一套评分模型, 评分模型最终为一个目标函数是线性函数, 约束条件含有二次等式约束的最优化问题。最后, 通过一个实例与Logistic回归做了对比, 实例结果表明模型是有效的, 且模型能够更好的支持实际业务应用场景。

## 关键词

信用评分, 贝叶斯判别, 区分度模型, 分箱权重, 尺度化

Copyright © 2019 by author(s) and Hans Publishers Inc.

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

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



Open Access

## 1. 概述

在信用评分方面, 美国 Fair Isaac Corp 公司的评分(以下简称 FICO 评分[1])是应用最为成熟的一种评分模型, 该评分在 300~850 之间, 信用评分越高, 说明客户的信用风险越小。除此之外, Logistic 回归[2][3]、决策树、神经网络等也是常用的评分方法。

本文借鉴了 FICO 评分的思想, 并基于贝叶斯判别定理推导出了一套信用评分模型(以下简称评分模型), 该模型是一个目标函数为线性函数, 约束条件为二次型的最优化模型。本文评分模型所具有的优点是:

- ① 本文评分模型得到的是非常直观的整数权重, 这对不懂评分技术的业务人员来讲, 能够很方便的对评分结果进行解读和应用。
- ② 当业务人员拒绝客户的信用业务申请时, 可以依据评分结果给予合理的拒绝原因。
- ③ 利用本文评分模型得到的多张评分卡, 可以方便的比较、混合使用。
- ④ 鉴于以往的项目经验: 本文评分模型的稳健性是非常好的, 利用本文方法建立好的评分卡应用3年后仍然有很好的预测性, 而Logistic回归、决策树要逊色很多。

## 2. 基于贝叶斯判别的评分模型

评分模型的建立分为 5 个步骤: 输入变量的筛选、输入变量的分箱、评分模型的求解与评估、评分结果的拟合与尺度化、评分模型的部署。

本文研究的内容主要是: 评分模型的建立、评分结果的拟合与尺度化。

假设有一批信用良好的客户样本(以下简称好客户样本)和信用不良的客户样本(以下简称坏客户样本), 我们要通过这两组样本数据建立评分模型。基于以往业务经验对评分模型做如下要求和假设:

- ① 客户信用评分  $S$  越大, 代表该客户是好客户的概率越大; 反之, 代表是坏客户的概率越大。
- ② 好客户样本评分  $S_g$  服从正态分布  $N(\bar{S}_g, \sigma_g^2)$ ,  $\bar{S}_g$  为好客户样本均值,  $\sigma_g^2$  为好客户样本方差;  $S_g$  的密度函数为  $f_g(S)$ 。
- ③ 坏客户样本评分  $S_b$  服从正态分布  $N(\bar{S}_b, \sigma_b^2)$ ,  $\bar{S}_b$  为坏客户样本均值,  $\sigma_b^2$  为坏客户样本方差;  $S_b$  的密度函数为  $f_b(S)$ 。
- ④  $p(g|S)$ : 信用评分为  $S$  的客户是好客户的概率;  $p(b|S)$ : 信用评分为  $S$  的客户是坏客户的概率;
- ⑤  $odds(S) = p(g|S)/p(b|S)$ : 信用评分  $S$  对应的好、坏客户的概率比。
- ⑥ 评分模型有  $p$  个输入变量, 各分箱组数分别是  $q_1, q_2, \dots, q_p$  个, 各分箱权重分别如下:  
第 1 个输入变量的分箱权重为:  $w_{11}, \dots, w_{1q_1}$

.....

第  $p$  个输入变量的分箱权重为:  $w_{p1}, \dots, w_{pq_p}$

记:  $\mathbf{w} = (w_{11}, \dots, w_{1q_1}, w_{21}, \dots, w_{2q_2}, \dots, w_{pq_p})^T$ ;  $T = \sum_{i=1}^p q_i$ : 总共分箱组数。

⑦ 好客户样本  $m$  个, 分别为:  $\mathbf{g}^{(1)}, \dots, \mathbf{g}^{(m)}$ ; 坏客户样本  $n$  个, 分别为:  $\mathbf{b}^{(1)}, \dots, \mathbf{b}^{(n)}$ 。

$$\mathbf{g}^{(r)} = (g_{11}^{(r)}, \dots, g_{1q_1}^{(r)}, \dots, g_{p1}^{(r)}, \dots, g_{pq_p}^{(r)})^T \quad (1 \leq r \leq m)$$

$$\mathbf{b}^{(r)} = (b_{11}^{(r)}, \dots, b_{1q_1}^{(r)}, \dots, b_{p1}^{(r)}, \dots, b_{pq_p}^{(r)})^T \quad (1 \leq r \leq n)$$

注: 对于每个样本第  $i$  ( $1 \leq i \leq p$ ) 个输入变量的分箱取值  $x_{i,1}, x_{i,2}, \dots, x_{i,q_i}$  中, 有且仅有一个分箱值为 1, 其他值为 0, 表示该样本第  $i$  输入变量值落在取值为 1 的分箱区间内。

⑧ 样本中好客户占比  $p_g = m/(m+n)$ , 坏客户占比  $p_b = n/(m+n)$ 。

⑨  $S_{g^{(r)}}$ : 第  $r$  个好客户样本的信用评分 ( $1 \leq r \leq m$ );  $S_{b^{(r)}}$ : 第  $r$  个坏客户样本的信用评分 ( $1 \leq r \leq n$ )。

基于①、②、③的要求和假设, 我们可以画出评分分布示意图如图 1 所示:

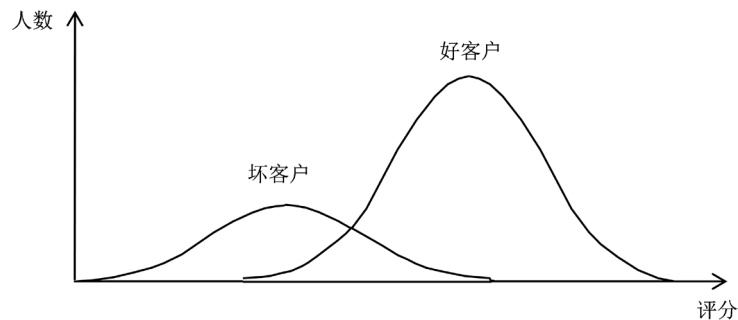


Figure 1. The diagram of score distribution  
图 1. 评分分布示意图

为了能尽量区分好、坏客户两个群体, 显然好客户的信用评分应该尽量大, 坏客户的信用评分应该尽量小, 我们以此建立我们评分模型的目标函数:

$$\max \sum_{r=1}^m S_{g^{(r)}} - \sum_{r=1}^n S_{b^{(r)}} \quad (1)$$

接下来, 我们确定评分模型的约束条件:

首先, 由贝叶斯判别定理[4][5], 可得

$$p(g|S) = p_g \cdot f_g(S), p(b|S) = p_b \cdot f_b(S)$$

$$\ln(\text{odds}(S)) = \ln(p_g/p_b) + \ln(f_g(S)/f_b(S)) \quad (2)$$

为了在拟合阶段中, 评分结果与客户好坏的概率值能建立函数关系, 我们在模型建立时就要考虑评分  $S$  与  $p(g|S), p(b|S)$  的关系。

由于  $\ln(p_g/p_b)$  是个定值, 不妨就假设:  $\ln(\text{odds}(S)) = \ln(p_g/p_b) + S$

$$\text{即 } S = \ln(f_g(S)/f_b(S)) \quad (3)$$

这样就建立了  $S$  与  $p(g|S), p(b|S)$  的函数关系, 且  $S$  越大,  $p(g|S)$  越大,  $p(b|S)$  越小。

又因为  $S_g \sim N(\bar{S}_g, \sigma_g^2)$ ,  $S_b \sim N(\bar{S}_b, \sigma_b^2)$ , 所以

$$\begin{aligned}\ln(f_g(S)/f_b(S)) &= \ln\left(\frac{1}{\sqrt{2\pi}\sigma_g}\exp\left(-\frac{(S-\bar{S}_g)^2}{2\sigma_g^2}\right)\right) - \ln\left(\frac{1}{\sqrt{2\pi}\sigma_b}\exp\left(-\frac{(S-\bar{S}_b)^2}{2\sigma_b^2}\right)\right) \\ &= \left(\frac{1}{2\sigma_b^2} - \frac{1}{2\sigma_g^2}\right) \cdot S^2 + \left(\frac{\bar{S}_g}{\sigma_g^2} - \frac{\bar{S}_b}{\sigma_b^2}\right) \cdot S + \left(\frac{\bar{S}_b^2}{2\sigma_b^2} - \frac{\bar{S}_g^2}{2\sigma_g^2}\right)\end{aligned}$$

由(3)得  $\sigma_g^2 = \sigma_b^2$ ,  $\bar{S}_g - \bar{S}_b = \sigma_g^2$ ,  $\bar{S}_g + \bar{S}_b = 0$ 。

这样我们就可以得到一个初步的评分模型：

$$\begin{aligned}\max \quad & \sum_{r=1}^m S_{g^{(r)}} - \sum_{r=1}^n S_{b^{(r)}} \\ \text{s.t.} \quad & \begin{cases} \sigma_g^2 - \sigma_b^2 = 0 \\ \bar{S}_g - \bar{S}_b = \sigma_g^2 \\ \bar{S}_g + \bar{S}_b = 0 \end{cases}\end{aligned}\quad (4)$$

实际上，好、坏客户两类群体信用评分的方差一般不会完全相等，而且在数据测试中我们发现  $\sigma_g^2$ 、 $\sigma_b^2$  不需要严格相等，效果会更好一些，这样我们可以把约束条件  $\sigma_g^2 - \sigma_b^2 = 0$  去掉，同时  $\bar{S}_g - \bar{S}_b = \sigma_g^2$  改为  $\bar{S}_g - \bar{S}_b = 0.5 \cdot (\sigma_g^2 + \sigma_b^2)$ 。

另外，如果对  $S$  不加约束，由(4)求得的目标函数会异常大，甚至求不出最优解，因此需要对信用评分  $S$  的取值范围加以约束，可以想到的方法有：

- ① 直接将  $S$  约束在某一区间范围内；
- ② 将每个分箱权重约束在某一范围内；
- ③ 设定各分箱权重的平方和小于某个阈值。

这3种方法都是有效的，但我们发现：第3种约束效果要好一些。另外，考虑到不同的评分模型其分箱组数是会不一样的，为了模型的普适性，我们采用“各分箱权重平方和的平均值小于某个阈值”来对  $S$  进行约束。这样，评分模型进一步优化为：

$$\begin{aligned}\max \quad & \sum_{r=1}^m S_{g^{(r)}} - \sum_{r=1}^n S_{b^{(r)}} \\ \text{s.t.} \quad & \begin{cases} \bar{S}_g - \bar{S}_b = 0.5 \cdot (\sigma_g^2 + \sigma_b^2) \\ \bar{S}_g + \bar{S}_b = 0 \\ \frac{1}{T} \cdot \sum_{i=1}^p \sum_{j=1}^{q_i} w_{ij}^2 \leq K \end{cases}\end{aligned}\quad (5)$$

其中， $T$  为分箱组数， $K$  为阈值(在本文实例计算中， $K = 2$  效果比较理想)。

### 3. 评分模型的参数推导

下面我们进行具体的参数推导，由前面的假设我们可以得出：

第  $r$  个好客户样本的信用评分：

$$S_{g^{(r)}} = (\mathbf{g}^{(r)})^T \cdot \mathbf{w} = \sum_{i=1}^p \sum_{j=1}^{q_i} g_{ij}^{(r)} \cdot w_{ij} \quad (1 \leq r \leq m)$$

第  $r$  个坏客户样本的信用评分：

$$S_{b^{(r)}} = (\mathbf{b}^{(r)})^T \cdot \mathbf{w} = \sum_{i=1}^p \sum_{j=1}^{q_i} b_{ij}^{(r)} \cdot w_{ij} \quad (1 \leq r \leq n)$$

好客户样本信用评分之和:

$$\sum_{r=1}^m S_{g^{(r)}} = \sum_{r=1}^m (\mathbf{g}^{(r)})^T \cdot \mathbf{w} = \left( \sum_{r=1}^m (\mathbf{g}^{(r)})^T \right) \cdot \mathbf{w}$$

坏客户样本信用评分之和:

$$\sum_{r=1}^n S_{b^{(r)}} = \sum_{r=1}^n (\mathbf{b}^{(r)})^T \cdot \mathbf{w} = \left( \sum_{r=1}^n (\mathbf{b}^{(r)})^T \right) \cdot \mathbf{w}$$

好客户样本信用评分平均值:

$$\bar{S}_g = \frac{1}{m} \sum_{r=1}^m S_{g^{(r)}} = \left( \frac{1}{m} \sum_{r=1}^m (\mathbf{g}^{(r)})^T \right) \cdot \mathbf{w}$$

坏客户样本信用评分平均值:

$$\bar{S}_b = \frac{1}{n} \sum_{r=1}^n S_{b^{(r)}} = \left( \frac{1}{n} \sum_{r=1}^n (\mathbf{b}^{(r)})^T \right) \cdot \mathbf{w}$$

好客户样本信用评分方差:

$$\begin{aligned} \sigma_g^2 &= E(S_g - \bar{S}_g)^2 = \frac{1}{m} \sum_{r=1}^m (S_{g^{(r)}} - \bar{S}_g)^2 = \frac{1}{m} \sum_{r=1}^m \left( \left( \mathbf{g}^{(r)} - \frac{1}{m} \sum_{r=1}^m \mathbf{g}^{(r)} \right)^T \cdot \mathbf{w} \right)^2 \\ &= \frac{1}{m} \sum_{r=1}^m \left( \mathbf{w}^T \cdot \left( \mathbf{g}^{(r)} - \frac{1}{m} \sum_{r=1}^m \mathbf{g}^{(r)} \right) \cdot \left( \mathbf{g}^{(r)} - \frac{1}{m} \sum_{r=1}^m \mathbf{g}^{(r)} \right)^T \cdot \mathbf{w} \right) \\ &= \frac{1}{m} \mathbf{w}^T \cdot \sum_{r=1}^m \left( \left( \mathbf{g}^{(r)} - \frac{1}{m} \sum_{r=1}^m \mathbf{g}^{(r)} \right) \cdot \left( \mathbf{g}^{(r)} - \frac{1}{m} \sum_{r=1}^m \mathbf{g}^{(r)} \right)^T \right) \cdot \mathbf{w} \end{aligned}$$

同理, 坏客户样本信用评分方差:

$$\sigma_b^2 = E(S_b - \bar{S}_b)^2 = \frac{1}{n} \sum_{r=1}^n (S_{b^{(r)}} - \bar{S}_b)^2 = \frac{1}{n} \mathbf{w}^T \cdot \sum_{r=1}^n \left( \left( \mathbf{b}^{(r)} - \frac{1}{n} \sum_{r=1}^n \mathbf{b}^{(r)} \right) \cdot \left( \mathbf{b}^{(r)} - \frac{1}{n} \sum_{r=1}^n \mathbf{b}^{(r)} \right)^T \right) \cdot \mathbf{w}$$

$$\mathbf{g} = \sum_{r=1}^m \mathbf{g}^{(r)}, \quad \bar{\mathbf{g}} = \frac{1}{m} \sum_{r=1}^m \mathbf{g}^{(r)}, \quad \mathbf{H}_g = \frac{1}{m} \sum_{r=1}^m (\mathbf{g}^{(r)} - \bar{\mathbf{g}})(\mathbf{g}^{(r)} - \bar{\mathbf{g}})^T,$$

$$\mathbf{b} = \sum_{r=1}^n \mathbf{b}^{(r)}, \quad \bar{\mathbf{b}} = \frac{1}{n} \sum_{r=1}^n \mathbf{b}^{(r)}, \quad \mathbf{H}_b = \frac{1}{n} \sum_{r=1}^n (\mathbf{b}^{(r)} - \bar{\mathbf{b}})(\mathbf{b}^{(r)} - \bar{\mathbf{b}})^T$$

$$\text{则 } \sum_{r=1}^m S_{g^{(r)}} = \mathbf{g}^T \cdot \mathbf{w}, \quad \bar{S}_g = \bar{\mathbf{g}}^T \cdot \mathbf{w}, \quad \sigma_g^2 = \mathbf{w}^T \cdot \mathbf{H}_g \cdot \mathbf{w}, \quad \sum_{r=1}^n S_{b^{(r)}} = \mathbf{b}^T \cdot \mathbf{w}, \quad \bar{S}_b = \bar{\mathbf{b}}^T \cdot \mathbf{w}, \quad \sigma_b^2 = \mathbf{w}^T \cdot \mathbf{H}_b \cdot \mathbf{w}$$

最后, 我们的评分模型就可以表示为:

$$\begin{aligned} &\max (\mathbf{g}^T - \mathbf{b}^T) \cdot \mathbf{w} \\ &\text{s.t.} \begin{cases} (\bar{\mathbf{g}}^T - \bar{\mathbf{b}}^T) \cdot \mathbf{w} = 0.5 \cdot \mathbf{w}^T \cdot (\mathbf{H}_g + \mathbf{H}_b) \cdot \mathbf{w} \\ (\bar{\mathbf{g}}^T + \bar{\mathbf{b}}^T) \cdot \mathbf{w} = 0 \\ \frac{1}{T} \cdot \mathbf{w}^T \cdot \mathbf{w} \leq K \end{cases} \end{aligned} \quad (6)$$

$T$  为分箱组数,  $K$  为阈值,  $\mathbf{w}$  是我们要求解的分箱权重向量, 模型是一个二次型最优化问题[6]。

#### 4. 评分结果的拟合与尺度化

在模型建立中, 我们假设:

$$\ln(odds(S)) = \ln(p_g/p_b) + S$$

因此, 我们采取  $S, \ln(odds(S))$  进行线性拟合。

在信用评分的实际应用中, 我们往往对某一具体的好坏概率比  $odds(S_0)$  特别重视, 期望该  $odds(S_0)$  对应某个评分  $S'_0$ , 不仅如此, 还要求信用评分  $S'$  每增加一个固定值  $\Delta S'_0$ , 好坏概率比  $odds(S')$  就增加一个  $odds(S_0)$ 。例如: 我们期望好坏概率比为 100 时对应的信用评分为 500 分, 且信用评分每增加 20 分, 好坏概率比就增加 100, 当信用评分为 700 分时, 可以推算出好坏概率比为 1100。

在此做如下假设:

$S'$ :  $S$  对应的尺度化后评分。

$odds(S_0)$ : 预设的好坏概率比。

$S_0$ :  $odds(S_0)$  对应的尺度化前评分。

$S'_0$ :  $odds(S_0)$  对应的尺度化后评分。

$\Delta S'_0$ : 表示尺度化后的评分值每增加  $\Delta S'_0$ , 好坏概率比就增加 1 个  $odds(S_0)$ 。

$w' = (w'_{11}, \dots, w'_{1q_1}, \dots, w'_{p1}, \dots, w'_{pq_p})^T$ : 尺度化后的各变量分箱权重。

进一步假设  $S, \ln(odds(S))$  拟合得到的线性方程为:

$$\ln(odds(S)) = b_0 + b_1 \cdot S \quad (b_0, b_1 \text{ 是系数}) \quad (7)$$

则我们可以得出如下关系:

$$\ln(odds(S')) = \ln(odds(S)) = b_0 + b_1 \cdot S \quad (8)$$

$$\frac{\ln(odds(S')) - \ln(odds(S'_0))}{S' - S'_0} = \frac{\ln(2 \cdot odds(S'_0)) - \ln(odds(S'_0))}{\Delta S'_0}$$

$$\text{即 } \ln(odds(S')) = \frac{\ln 2}{\Delta S'_0} \cdot S' + \ln(odds(S'_0)) - \frac{\ln 2}{\Delta S'_0} \cdot S'_0 \quad (9)$$

由(8)进一步可得:

$$\ln(odds(S'_0)) = \ln(odds(S_0)) = b_0 + b_1 \cdot S_0 \quad (10)$$

将(8)、(10)代入(9), 得:

$$b_0 + b_1 \cdot S = \frac{\ln 2}{\Delta S'_0} \cdot S' + b_0 + b_1 \cdot S_0 - \frac{\ln 2}{\Delta S'_0} \cdot S'_0$$

$$S' = \Delta S'_0 \cdot \frac{b_1}{\ln 2} \cdot S + S'_0 - \Delta S'_0 \cdot \frac{b_1}{\ln 2} \cdot S_0$$

令  $c_1 = \Delta S'_0 \cdot \frac{b_1}{\ln 2}$ ,  $c_0 = S'_0 - c_1 \cdot S_0$ , 则

$$S' = c_0 + c_1 \cdot S \quad (11)$$

$$\ln(odds(S')) = b_0 - \frac{b_1 \cdot c_0}{c_1} + \frac{b_1}{c_1} \cdot S' \quad (12)$$

式(11)就是尺度化评分  $S'$  与原始评分  $S$  的尺度化关系, 式(12)就是尺度化后评分  $S'$  与好坏概率比的关系。

需要补充说明是：利用式(7)进行拟合时，实际上并不知道每个原始评分  $S$  对应的好坏客户概率比，但是我们可以对原始评分进行排序分组，然后取每个分组原始评分的中间值作为  $S$ ，每个组的好坏客户数比作为  $odds(S)$ ，这样就可以进行拟合了。另外，考虑到按原始评分排序分组以后， $S$  值最大的几个组里可能没有坏客户， $S$  值最小的几个组里可能没有好客户，所以要剔除这些“特殊”组，然后再进行拟合。

最后，我们将尺度化评分  $S'$  拆分到每个变量分箱中。拆分时要遵循如下 2 条原则：

- ① 每个变量的各分箱权重非负。
- ② 各样本的尺度化分箱权重之和仍为  $S'$ 。

记第  $i$  个变量的尺度化前最小分箱权重  $\min(w_i) = \min_{1 \leq j \leq q_i}(w_{ij})$  ( $1 \leq i \leq p$ )，则

$$w'_{ij} = c_1 \cdot (w_{ij} + |\min(w_i)|) + \frac{c_0 - c_1 \cdot \sum_{i=1}^p |\min(w_i)|}{p} \quad (1 \leq i \leq p, 1 \leq j \leq q_i) \quad (13)$$

其中  $w_{ij} + |\min(w_i)|$  是为了将变量的最小分箱权重由负值变为零值，乘以  $c_1$  表示的是每个分箱权重的尺度化也服从  $S$  到  $S'$  的线性关系，加上  $\left(c_0 - c_1 \cdot \sum_{i=1}^p |\min(w_i)|\right) / p$  是为了保证尺度化后的评分值仍然等于尺度化后的分箱权重之和。

例如：假设  $c_0 = 115.8, c_1 = 23, \Delta S'_0 = 20$ ，有 3 个分箱变量，尺度化过程可用如表 1 所示：

**Table 1.** Weight scaling for binning variables  
**表1.** 分箱权重尺度化步骤说明表

变量	$w_{ij}$	$ \min(w_i) $	$c_i$	$c_1 \cdot (w_{ij} +  \min(w_i) )$	$\left(c_0 - c_1 \cdot \sum_{i=1}^p  \min(w_i) \right) / p$	$w'_{ij}$
变量 1	2.1			76	11	87
	0.1	1.2	23	30	11	41
	-1.2			0	11	11
变量 2	-0.5			0	11	11
	0.5	0.5	23	23	11	34
	1			67	11	78
变量 3	0.1	1.9	23	46	11	57
	-1.9			0	11	11

## 5. 数据测试与对比

数据来源：SPSS自带的bankloan.sav数据，包含：517位拖欠贷款客户(坏客户)，183位不拖欠贷款客户(好客户)。

输出变量：default (1：坏客户；0：好客户)。

输入变量及分箱结果：见表2。

尺度化要求：500 分对应的好坏概率比是 100:1，且尺度化后的评分每增加 20 分，好坏概率比增加 100。

阈值设置： $K = 2$ 。

利用样本数据计算结果如表 3 所示：

**Table 2.** The binning of input variables  
**表2.** 输入变量及分箱表

变量	<i>address</i> 地址	<i>age</i> 年龄	<i>creddebt</i> 信用卡欠款	<i>debtinc</i> 贷款收入比	<i>ed</i> 学历	<i>employ</i> 工龄	<i>income</i> 收入	<i>othdebt</i> 其他债务
分箱1					ed = 1		(-, 27)	
分箱2			(-, 1.0615)	(-, 5]	ed = 2	(-, 3)	[27, 37)	
分箱3	(-, 7)	(-, 34)	[1.0615, 1.9607)	[5, 8.5)	ed = 3	[3, 7)	[37, 56)	(-, 1.926)
分箱4	[7, +)	[34, +)	[1.9607, +)	[8.5, 13.8)	ed = 4	[7, 12)	[56, +)	[1.926, +)
分箱5				(13.8, +)	ed = 5	[12, +)		

**Table 3.** Calculating results of the scoring model  
**表3.** 评分模型计算结果

变量	变量分箱	<i>g</i>	<i>b</i>	$\bar{g}$	$\bar{b}$	原始权重	尺度化权重
<i>address</i>	(-, 7)	223	117	0.43	0.64	-1.23	36
	[7, +)	294	66	0.57	0.36	1.16	58
<i>age</i>	(-, 34)	218	110	0.42	0.60	-0.26	36
	[34, +)	299	73	0.58	0.40	0.19	40
<i>creddebt</i>	(-, 1.0615)	322	75	0.62	0.41	2.15	76
	[1.0615, 1.9607)	92	40	0.18	0.22	-0.09	55
	[1.9607, +)	103	68	0.20	0.37	-2.14	36
<i>debtinc</i>	(-, 5]	152	17	0.29	0.09	2.01	86
	[5, 8.5)	145	28	0.28	0.15	1.17	78
	[8.5, 13.8)	130	45	0.25	0.25	0.09	68
	(13.8, +)	90	93	0.17	0.51	-3.35	36
<i>ed</i>	ed = 1	293	79	0.57	0.43	1.34	57
	ed = 2	139	59	0.27	0.32	-0.15	43
	ed = 3	57	30	0.11	0.16	-0.94	36
	ed = 4	24	14	0.05	0.08	-0.32	41
	ed = 5	4	1	0.01	0.01	-0.01	44
<i>employ</i>	(-, 3)	79	76	0.15	0.42	-3.37	36
	[3, 7)	119	52	0.23	0.28	-0.49	63
	[7, 12)	140	30	0.27	0.16	1.51	82
	[12, +)	179	25	0.35	0.14	2.26	89
<i>income</i>	(-, 27)	153	77	0.30	0.42	-0.56	36
	[27, 37)	108	39	0.21	0.21	-0.20	39
	[37, 56)	116	35	0.22	0.19	-0.22	39
	[56, +)	140	32	0.27	0.17	0.91	49
<i>othdebt</i>	(-, 1.926)	265	77	0.51	0.42	-0.31	36
	[1.926, +)	252	106	0.49	0.58	0.24	41



尺度化评分  $S'$  与  $\ln(odds(S'))$  如图 2 所示:

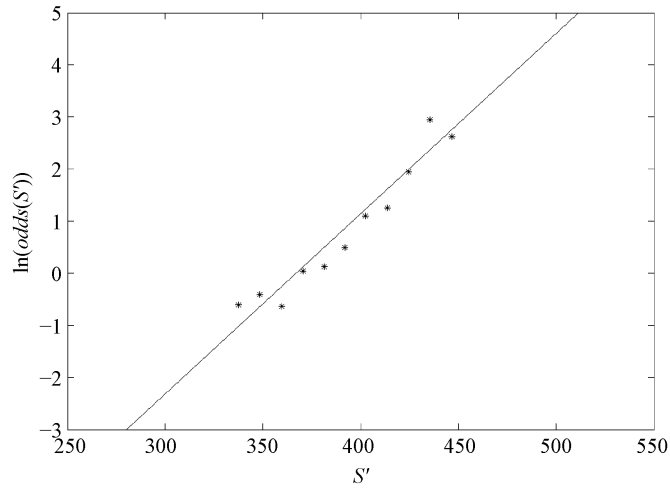


Figure 2. Linear fitting chart by  $S'$  and  $\ln(odds(S'))$

图2.  $S'$  与  $\ln(odds(S'))$  拟合直线图

本文评分模型与Logistic回归ROC曲线比较, 如图3所示:

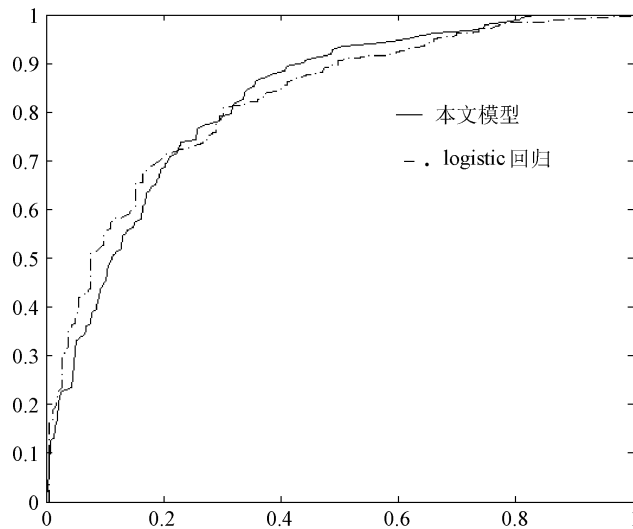


Figure 3. ROC curve comparison chart

图3. ROC曲线对比图

可以看出: 文中的评分模型跟 Logistic 回归模型相比也是一种非常有效的评分方法。另外, 基于以往的项目经验: 本文评分模型的稳健性是非常好的, 利用本文方法建立的评分卡应用 3 年后仍然有很好的预测性, 而 Logistic 回归、决策树要逊色很多。

### 参考文献

- [1] 陈建. 信用评分模型技术与应用[M]. 北京: 中国财政经济出版社, 2005: 1-286.
- [2] 杨静. 信用评分卡的建立与应用[D]: [硕士学位论文]. 天津: 天津商业大学, 2018.

- 
- [3] 石勇, 孟凡. 信用评分基本理论及其应用[J]. 大数据, 2017(1): 24-31.
- [4] 马达. 基于贝叶斯的判别理论及其算法实现[D]: [硕士学位论文]. 北京: 中国地质大学, 2011.
- [5] 宋云鹏, 武钰. 数据挖掘技术在信用评分中的应用研究[J]. 征信, 2013(10): 24-28.
- [6] 袁亚湘. 非线性优化计算方法[M]. 北京: 科学出版社, 2018.

**知网检索的两种方式:**

1. 打开知网首页: <http://cnki.net/>, 点击页面中“外文资源总库 CNKI SCHOLAR”, 跳转至: <http://scholar.cnki.net/new>, 搜索框内直接输入文章标题, 即可查询;  
或点击“高级检索”, 下拉列表框选择: [ISSN], 输入期刊 ISSN: 2163-1476, 即可查询。
2. 通过知网首页 <http://cnki.net/>顶部“旧版入口”进入知网旧版: <http://www.cnki.net/old/>, 左侧选择“国际文献总库”进入, 搜索框直接输入文章标题, 即可查询。

投稿请点击: <http://www.hanspub.org/Submission.aspx>

期刊邮箱: [orf@hanspub.org](mailto:orf@hanspub.org)