

“十二五”国家重点图书出版规划项目



应用统计工程前沿丛书

# 风险模型： 基于R的保险损失预测

孟生旺 著



清华大学出版社



## 内 容 简 介

保险是经营风险的行业,风险的评估和定价是保险公司最为核心的竞争力。本书以保险业为研究对象,讨论了相应的风险模型及其应用,主要包括损失概率、损失次数、损失金额和累积损失的分布模型以及它们的预测模型,同时还探讨了巨灾损失和相依风险的建模问题。在实证研究中,以 R 语言为计算工具,提供了详细的程序代码,方便读者再现完整的计算过程。

本书适合风险管理、保险与精算等相关专业的高年级学生、研究人员或从业人员参考。

本书封面贴有清华大学出版社防伪标签,无标签者不得销售。

版权所有,侵权必究。侵权举报电话:010-62782989 13701121933

### 图书在版编目(CIP)数据

风险模型:基于 R 的保险损失预测/孟生旺著. —北京:清华大学出版社,2017  
(应用统计工程前沿丛书)  
ISBN 978-7-302-48206-2

I. ①风… II. ①孟… III. ①保险业—风险管理—研究 IV. ①F840.323

中国版本图书馆 CIP 数据核字(2017)第 209666 号

责任编辑:魏贺佳

封面设计:傅瑞学

责任校对:刘玉霞

责任印制:杨 艳

出版发行:清华大学出版社

网 址: <http://www.tup.com.cn>, <http://www.wqbook.com>

地 址:北京清华大学学研大厦 A 座 邮 编:100084

社总机:010-62770175 邮 购:010-62786544

投稿与读者服务:010-62776969, [c-service@tup.tsinghua.edu.cn](mailto:c-service@tup.tsinghua.edu.cn)

质量反馈:010-62772015, [zhiliang@tup.tsinghua.edu.cn](mailto:zhiliang@tup.tsinghua.edu.cn)

印 装 者:三河市金元印装有限公司

经 销:全国新华书店

开 本:170mm×230mm 印 张:27.75

字 数:513千字

版 次:2017年9月第1版

印 次:2017年9月第1次印刷

定 价:89.00元

---

产品编号:076146-01

# “应用统计工程前沿丛书”

## 编 委 会

顾问：袁 卫 吴喜之 易丹辉 胡飞芳

主任：赵彦云 金勇进

委员：王晓军 张 波 孟生旺 许王莉 吕晓玲

蒋 妍 李静萍 王 星 肖宇谷

# 为中国的应用统计开拓奋进

（“应用统计工程前沿丛书”代序）

改革开放以来，我国统计事业取得了突飞猛进的发展。市场化、全球化和信息技术与网络经济的蓬勃发展，使统计在经济、社会、管理、医学、生物、农业、工程等领域中的应用迎来了又一春天。2011年2月，国务院学位委员会第28次会议通过了新的《学位授予和人才培养学科目录(2011)》，将统计学上升为一级学科，这是国家对统计学科建设与发展的重大支持，它将全面推动统计学理论方法和应用研究的深入发展。

长期以来，我国统计学科主要在经济学、理学和医学等门类下发展，未来进一步发展，一级统计学科将成为一面旗帜。世界上先进国家的实践充分表明，统计广泛应用在各个学科，在信息网络技术与计算机强大能力的推动下，统计学科发展特别是统计的应用正展示出一种前所未有的时代特征，它将为创造新的人类文明、提升人类发展能力做出新的重要贡献。

新中国把中国从一盘散沙凝聚成高度集中的国家，推行计划经济发展模式。这一时期，统计直接为计划服务，为政府各级管理部门，为企事业单位的计划管理，为市场资源配置，为消费、投资的安排等提供全面系统的服务，因此在经济社会管理中发挥了重要作用。但是，由于权力至上的落后观念和体系机制呆滞，统计的科学性不被重视，统计数据搜集整理的简单化和主观操作造成了很多不良的后果。改革开放之后，市场的作用强化了统计的社会影响和地位，但是，惯性的从上向下的主观思维方式仍然没有彻底的改观，因此，统计的科学应用仍然需要依靠内生发展的强大魅力不断深入和扩大。

近年来，全球化进一步加速了经济结构的转型与效率的提高。事实上，一国的稳步可持续发展离不开扎实的基础。在当今的信息化网络化时代，信息基础设施及其运用效率成为基础的基础，伴随而来的是统计在搜集数据、整理数据、分析数据上发挥的重要基础性作用。电子金融、电子政务、电子商务、网上购物、微博等一系列以网络信息技术为支撑的经济社会活动创造了大数据的新时代，计算机科学、数据库技术、大数据统计分析成为新时代发展的耀眼之星，统计学理论方法在海量数据挖掘分析、高维分析和复杂系统模型分析，以及时空的统计图示图解分析等方面正显示出强劲发展的能量，应该讲现时期是统计应用最好的发展机遇，它将大大提高人类发展的创造力、生产力，造福社会、造福人类。

## 二

在发展非凡的年代，谁能插上翅膀自由翱翔，谁能潜下海底自由鱼跃，统计学科需当仁不让，测度方位、穿透迷雾、指引方向、科学决策，助国家繁荣昌盛，立世界之林，这是当今中国人民大学统计学科建设的基本认知和理念。中国人民大学统计学科成立于1950年，已有60多年的发展历程，为共和国建设培养了大批优秀人才。他们广泛分布在政府部门以及银行、保险、证券、数据调查与咨询等商业企业，发挥了骨干作用。几代人大统计学人的辉煌历程和奉献，铸就了中国人民大学应用统计的特色，其作为国家应用统计重点学科、教育部重点研究基地和国家统计局重点研究基地，在融入世界一流队列、开拓中国应用、培养高精尖应用统计人才、全方位支持国家建设和发展中，做出了重要的贡献。

今天，中国人民大学统计学科布局不仅深入经济社会发展领域和保险精算与金融风险管理工作领域，而且已经扩展到人文社会科学的许多领域，如法律、新闻、政治学、伦理学、教育学、心理学、文献计量等，展示出应用统计在量化人文社会科学研究中的重要作用。同时，我们也在生物、医学与公共健康领域开展了深入的统计交叉应用研究。建设扎实的概率论与数理统计基础，发展强大的应用统计是中国人民大学统计学院继往开来的基本目标。

## 三

为了系统总结和凝练中国人民大学在统计学各个领域的科研成果，引领和推动我国统计学学科建设，提高统计学在人文社会科学与自然科学各领域科学研究，以及和管理、决策支持等方面应用的科学化和普及水平，促进统计学及其交叉学科人才培养，我们组织编写了这套“应用统计工程前沿丛书”。丛书选题覆盖应用统计学的主要分支领域，如人文、社会、政治、经济、金融、管理、法律、教育、生物、卫生、网络、数据挖掘等，力求在科学性、应用性、创新性、前沿性和可读性上形成特色。

丛书针对各领域的实际问题，着重统计学方法、模型的创新、设计和应用。在应用领域的具体统计问题研究上，积极发展统计应用流程科学，强调应用背景描述清晰，基础问题明确，发挥对微观数据、大量数据归纳探索与挖掘的统计方法作用，发展标准化的统计思维方法，创建应用领域的重要统计模型，深入解决问题，推动应用领域适应信息社会的高速发展。我们首次提出应用统计工程一词。工程是将自然科学原理应用到工农业生产部门中去而形成的各学科的总称。“工程”是科学的某种应用，通过这一应用，使自然界的物质和能源的特性能够通过各种结构、机器、产品、系统和过程，以最短的时间和少而精的人力做出高效、可靠且对人类有用的东西。我们强调应用统计的工程性，也就是强调统计的实际应用价值、科学流程与先进的统计应用技术。

丛书要反映统计学科多个前沿领域的科研进展，反映信息化和网络化背景下在诸

多统计学应⽤领域产生的新的统计学问题及其方法和模型的发展,以及在⼈⽂社会科学各个领域的开创性应⽤研究。丛书选题覆盖了应⽤统计学的各主要分⽀学科和主要新兴应⽤领域,系统总结和凝练应⽤统计的专门技术方法,引领和推动我国大数据中的统计科学方法及其应⽤,提高网络信息统计处理与网络经济活动与经营活动的统计科学分析能⼒,提高统计学在企业经营管理、市场营销、科学决策,以及全面提升综合竞争⼒⽅⾯的作用,提高统计学在宏观经济产业政策、货币政策、收⼊分配政策等重大政策制定与效果分析,以及全面提升我国国际竞争⼒和国家软实⼒⽅⾯的作用。

本套丛书主要面向统计学及其交叉学科领域的科研人员、研究生和高年级本科生,以及在实际工作中需要应⽤统计学理论与方法的各领域专业人士。丛书在理论方法与应⽤领域深⼊结合研究上,强调增加关键点的细节内容,突出以统计知识为核心的应⽤领域的统计知识体系建设。丛书在内容上力求拥有清晰的逻辑结构;对方法、概念和统计问题的描述增加相关概念知识和应⽤背景及交叉学科知识运⽤的铺垫;同时给出相关参考文献或推荐阅读⽬录,以帮助有兴趣的读者进⼊深⼊学习。奉献给相关专业的读者能读懂并能够学以致⽤的应⽤统计,这是本丛书追求的重要⽬标之一。

赵彦云 吕晓玲

2014年12月

# | 前言

保险是经营风险的行业,风险的评估和定价是保险公司最核心的竞争力。风险的内涵十分丰富,可以从不同的角度进行划分和归类。以保险风险为例,可以分为财产风险、人身风险、责任风险、信用风险等。本书所谓的风险,主要是指保险风险,或者更具体地说,是指保险损失的风险。保险损失具体表现为损失概率、损失次数和损失金额的大小,相应地,风险模型也就包括损失概率模型、损失次数模型、损失金额模型和累积损失模型。本书讨论的风险模型虽然以财产与责任保险业务为主要背景,但也可以扩展到信用风险评估和金融风险管理等领域,具有更加广泛的应用价值。

作者在中国人民大学统计学院为风险管理与精算专业的研究生讲授“风险模型”课程已有十余年,在此期间先后完成了包括国家社会科学基金重大项目、国家自然科学基金面上项目、教育部人文社会科学重点研究基地重大项目在内的十余项风险管理与精算方向的研究课题,取得了一定的研究成果。本书就是结合作者十余年的“风险模型”教学经验和部分课题的研究成果撰写而成。

全书共由十三章内容构成,主要介绍了风险模型的理论性质、数据拟合方法以及基于R的实际应用,适合风险管理、保险和精算等相关专业的研究生以及精算师、风险管理师等专业人士参考。

在写作过程中,注重内容的完整性、系统性和前沿性,强调理论模型在解决实际风险管理问题中的应用。为了方便读者重现有关实证分析的具体过程,提供了完整的R程序代码和数据集,可以通过书中提供的链接地址下载。

本书的部分内容是作者主持完成的下述科研项目的阶段性成果:国家社会科学基金重大项目“巨灾保险的精算统计模型及其应用研究”(16ZDA052),教育部人文社会科学重点研究基地重大项目“基于大数据的精算统计模型与风险管理问题研究”(16JJD910001)。

对于本书可能存在的任何缺陷,作者负有不可推卸之责任,欢迎各位读者批评指正,以期再版时得以修正。今后如有补充或更新材料,将及时在作者的新浪博客上(<http://blog.sina.com.cn/mengshw>)发布。

孟生旺

中国人民大学统计学院教授,博士生导师

中国人民大学应用统计科学研究中心研究员

甘肃省“飞天学者”特聘计划兰州财经大学讲座教授

第 1 章 风险度量	1
1.1 描述随机变量的函数	2
1.1.1 分布函数	2
1.1.2 概率密度函数	4
1.1.3 生存函数	5
1.1.4 概率母函数	6
1.1.5 矩母函数	7
1.1.6 危险率函数	8
1.2 常用的风险度量方法	9
1.2.1 VaR	10
1.2.2 TVaR	14
1.2.3 基于扭曲变换的风险度量	19
第 2 章 损失金额分布模型	31
2.1 常用的损失金额分布	32
2.1.1 正态分布	32
2.1.2 指数分布	33
2.1.3 伽马分布	35
2.1.4 逆高斯分布	37
2.1.5 对数正态分布	39
2.1.6 帕累托分布	41
2.1.7 韦布尔分布	42
2.2 新分布的生成	44
2.2.1 函数变换	44
2.2.2 混合分布	50
2.3 免赔额的影响	53
2.4 赔偿限额的影响	60
2.5 通货膨胀的影响	65



<b>第 3 章 损失次数分布模型</b> .....	70
3.1 $(a, b, 0)$ 分布类 .....	71
3.1.1 泊松分布 .....	71
3.1.2 二项分布 .....	74
3.1.3 负二项分布 .....	76
3.1.4 几何分布 .....	79
3.2 $(a, b, 1)$ 分布类 .....	80
3.2.1 零截断分布 .....	81
3.2.2 零调整分布 .....	83
3.3 零膨胀分布 .....	84
3.4 复合分布 .....	85
3.4.1 复合分布的概率计算 .....	86
3.4.2 复合分布的比较 .....	89
3.5 混合分布 .....	95
3.6 免赔额对损失次数模型的影响 .....	99
3.6.1 免赔额对 $(a, b, 0)$ 分布类的影响 .....	100
3.6.2 免赔额对 $(a, b, 1)$ 分布类的影响 .....	100
3.6.3 免赔额对复合分布的影响 .....	101
<b>第 4 章 累积损失分布模型</b> .....	103
4.1 集体风险模型 .....	104
4.1.1 精确计算 .....	105
4.1.2 参数近似 .....	111
4.1.3 Panjer 递推法 .....	117
4.1.4 傅里叶近似 .....	126
4.1.5 随机模拟 .....	131
4.2 个体风险模型 .....	138
4.2.1 卷积法 .....	138
4.2.2 参数近似法 .....	141
4.2.3 复合泊松近似法 .....	143
<b>第 5 章 损失分布模型的参数估计</b> .....	147
5.1 参数估计 .....	148

5.1.1	极大似然法	148
5.1.2	矩估计法	156
5.1.3	分位数配比法	157
5.1.4	最小距离法	158
5.2	模型的评价和比较	162
<b>第 6 章</b>	<b>巨灾损失模型</b>	<b>166</b>
6.1	广义极值分布	167
6.1.1	极值分布函数	169
6.1.2	极大吸引域	171
6.1.3	区块最大化方法	172
6.2	广义帕累托分布	173
6.2.1	分布函数	173
6.2.2	超额损失的分布	174
6.2.3	更大阈值下超额损失的分布	177
6.2.4	尾部生存函数	178
6.2.5	风险度量	178
6.2.6	参数的极大似然估计	179
6.2.7	尾部指数的 Hill 估计	180
6.2.8	尾部生存函数的 Hill 估计	182
6.3	偏正态分布和偏 $t$ 分布	189
<b>第 7 章</b>	<b>损失预测的广义线性模型</b>	<b>195</b>
7.1	广义线性模型的结构	196
7.1.1	指数分布族	197
7.1.2	连接函数	203
7.2	模型的参数估计方法	204
7.2.1	极大似然估计	204
7.2.2	牛顿迭代法	206
7.2.3	迭代加权最小二乘法	207
7.2.4	牛顿迭代法与迭代加权最小二乘法的比较	212
7.2.5	离散参数的估计	212
7.2.6	参数估计值的标准误	213
7.3	模型的比较与诊断	213

7.3.1	偏差 .....	214
7.3.2	模型比较 .....	218
7.3.3	伪判定系数 .....	221
7.3.4	残差 .....	223
7.3.5	Cook 距离 .....	225
7.3.6	连接函数的诊断 .....	225
<b>第 8 章</b>	<b>损失金额预测模型 .....</b>	<b>227</b>
8.1	线性回归模型 .....	228
8.1.1	模型设定 .....	228
8.1.2	参数估计 .....	230
8.1.3	连接函数 .....	230
8.1.4	模拟数据分析 .....	231
8.2	损失金额预测的伽马回归 .....	234
8.2.1	模型设定 .....	234
8.2.2	迭代加权最小二乘估计 .....	235
8.2.3	模拟数据分析 .....	236
8.3	损失金额预测的逆高斯回归 .....	238
8.3.1	模型设定 .....	240
8.3.2	迭代加权最小二乘估计 .....	241
8.3.3	模拟数据分析 .....	241
8.3.4	GAMLSS 的应用 .....	244
8.4	有限赔款预测模型 .....	248
8.5	混合损失金额预测模型 .....	252
8.6	应用案例 .....	256
8.6.1	数据介绍 .....	256
8.6.2	描述性分析 .....	259
8.6.3	案均赔款的预测模型 .....	261
8.6.4	案均赔款对数的预测模型 .....	266
<b>第 9 章</b>	<b>损失概率预测模型 .....</b>	<b>271</b>
9.1	基于个体观察数据的损失概率预测 .....	273
9.1.1	伯努利分布 .....	273
9.1.2	伯努利分布假设下的逻辑斯谛回归 .....	273

9.1.3	迭代加权最小二乘估计	275
9.1.4	模拟数据分析	276
9.1.5	不同风险暴露时期的处理	278
9.2	基于汇总数据的损失概率预测	282
9.2.1	二项分布	282
9.2.2	二项分布假设下的逻辑斯谛回归	283
9.2.3	迭代加权最小二乘估计	285
9.2.4	模拟数据分析	286
9.3	损失概率预测模型的解释	290
9.4	损失概率预测模型的评价	292
9.4.1	偏差	292
9.4.2	分类表	292
9.4.3	Hosmer-Lemeshow 统计量	295
9.5	其他连接函数	296
9.6	过离散问题	299
9.7	应用案例	300
<b>第 10 章</b>	<b>损失次数预测模型</b>	<b>306</b>
10.1	泊松回归模型	307
10.1.1	泊松分布	307
10.1.2	模型设定	308
10.1.3	迭代加权最小二乘估计	309
10.1.4	抵消项	310
10.1.5	模型参数的解释	311
10.1.6	模拟分析	311
10.2	过离散损失次数预测模型	314
10.2.1	负二项 I 型分布	315
10.2.2	负二项 II 型分布	317
10.2.3	迭代加权最小二乘估计	318
10.2.4	模型参数的解释	319
10.2.5	模拟分析	320
10.3	零截断与零膨胀损失次数预测模型	322
10.3.1	零截断回归模型	322
10.3.2	零膨胀回归模型	324

10.3.3	零调整回归模型 .....	329
10.4	混合损失次数预测模型 .....	333
10.5	应用案例 .....	335
10.5.1	描述性分析 .....	335
10.5.2	索赔频率预测模型 .....	337
<b>第 11 章</b>	<b>累积损失的预测模型 .....</b>	<b>342</b>
11.1	Tweedie 回归 .....	343
11.2	零调整逆高斯回归 .....	351
11.3	应用案例 .....	356
11.3.1	描述性分析 .....	356
11.3.2	纯保费的预测模型 .....	358
<b>第 12 章</b>	<b>相依风险模型 .....</b>	<b>366</b>
12.1	Copula .....	367
12.2	生存 Copula .....	372
12.3	相依性的度量 .....	374
12.3.1	线性相关系数 .....	374
12.3.2	秩相关系数 .....	375
12.3.3	尾部相依指数 .....	376
12.4	常见的 Copula 函数 .....	377
12.4.1	正态 Copula .....	377
12.4.2	t-Copula .....	377
12.4.3	Clayton Copula .....	377
12.4.4	Frank Copula .....	378
12.4.5	Gumbel Copula .....	378
12.4.6	FGM Copula .....	379
12.4.7	厚尾 Copula .....	379
12.5	阿基米德 Copula .....	379
12.6	Copula 的随机模拟 .....	381
12.7	Copula 的参数估计 .....	386
12.8	Copula 的应用 .....	388
<b>第 13 章</b>	<b>贝叶斯风险模型 .....</b>	<b>396</b>
13.1	先验分布的选择 .....	397

13.2	MCMC 方法简介 .....	399
13.2.1	Gibbs 抽样 .....	400
13.2.2	Metropolis-Hastings 算法 .....	400
13.2.3	Hamiltonian Monte Carlo 算法 .....	401
13.2.4	收敛性的诊断 .....	401
13.3	模型评价 .....	403
13.4	贝叶斯模型的应用 .....	403
索引	.....	420
参考文献	.....	423

# 第1章 风险度量

保险是经营风险的行业,风险的评估和度量是保险公司最核心的竞争力。风险的内涵十分丰富,可以从不同的角度进行划分和归类,以保险风险为例,可以分为财产风险、人身风险、责任风险、信用风险等。本书所谓的风险,主要是指保险风险,或者更具体地说,是指保险损失的风险。

风险通常被定义为事件发生结果的不确定性。对于保险而言,风险是指保险损失的不确定性,具体表现为保险事故发生与否的不确定性,事故发生时间的不确定性,事故发生地点的不确定性,事故发生次数的不确定性,以及损失金额的不确定性。

随机变量是描述不确定性的常用工具,所以保险损失也可以用随机变量进行描述。为此,本章首先介绍描述随机变量的有关函数,包括分布函数、概率密度函数、生存函数、概率母函数、矩母函数和危险率函数,然后介绍一些常用的风险度量方法,包括 VaR、TVaR 和基于扭曲变换的风险度量方法。描述随机变量的这些函数都可以完整刻画损失的分布情况,而风险度量则是对这些函数的一种高度概括,它通过一个实值来反映风险的大小,可以更加容易地应用于实际的风险管理。

## 1.1 描述随机变量的函数

对于保险而言,损失随机变量一般是非负的,可以分为连续型变量(如损失金额)和离散型变量(如损失次数)两大类。当然,也存在一些混合型损失随机变量,如保单的累积损失,一方面在零点有一个较高的概率堆积,另一方面在大于零的部分又是连续的。无论是哪种类型的损失随机变量,都可以用一个函数进行描述。本节主要介绍刻画损失随机变量的常用函数,如分布函数、概率密度函数、生存函数、概率母函数、矩母函数、危险率函数,这些函数是建立风险模型的基本工具。

### 1.1.1 分布函数

令  $X$  表示损失随机变量,则其分布函数定义为

$$F(x) = \Pr(X \leq x)$$

上式表明,损失随机变量  $X$  的分布函数就是  $X$  小于或等于  $x$  的概率。

**【例 1-1】** 随机变量  $X$  的取值范围为  $(10, 30, 40, 70, 90)$ , 取每个值的概率均为  $1/5$ , 求  $X$  的分布函数。

**【解】** 根据分布函数的定义,随机变量  $X$  的分布函数如下:

$$F(10) = \Pr(X \leq 10) = 1/5 = 0.2$$

$$F(30) = \Pr(X \leq 30) = 2/5 = 0.4$$

$$F(40) = \Pr(X \leq 40) = 3/5 = 0.6$$



$$F(70) = \Pr(X \leq 70) = 4/5 = 0.8$$

$$F(90) = \Pr(X \leq 90) = 5/5 = 1$$

绘制分布函数的 R 程序代码如下,其中函数 `ecdf()` 表示经验累积分布函数 (empirical cumulative distribution function)。输出结果如图 1-1 所示,其中横轴表示随机变量  $x$  的取值,纵轴表示分布函数  $F(x)$  的取值。

```
# 随机变量的取值
x = c(10, 30, 40, 70, 90)

# 绘制分布函数
plot(ecdf(x), main = '', ylab = 'F(x)')
```

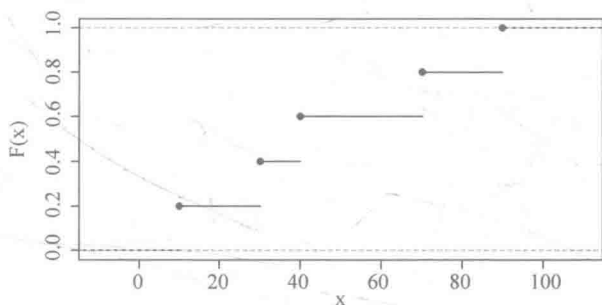


图 1-1 随机变量的分布函数\*

**【例 1-2】** 假设损失服从伽马分布,形状参数为  $\text{shape}=2$ ,比率参数为  $\text{rate}=1/3$ ,绘制其分布函数,并计算损失小于 10 的概率和损失大于 4 的概率。

**【解】** 在 R 程序中,可以用函数 `pgamma()` 计算伽马分布的分布函数。在该函数中,第一个参数默认为形状参数  $\text{shape}$ ,第二个参数默认为比率参数  $\text{rate}$ ,相应地,伽马分布的均值等于形状参数除以比率参数。

绘制分布函数的 R 程序代码如下,输出结果如图 1-2 所示,其中横轴表示随机变量  $x$  的取值,纵轴表示分布函数  $F(x)$  的取值。

```
curve(pgamma(x, 2, 1/3), lwd = 2, xlim = c(0, 20), ylab = 'F(x)')
```

伽马分布的第二个参数也可以表示为尺度参数  $\text{scale}$ ,此时在程序代码中必须明确声明  $\text{scale}=1/\text{rate}$ 。譬如,在本例的程序代码中,如果使用  $\text{scale}$  参数,则伽马分布函数应该表示为 `pgamma(x, shape=2, scale=3)`。计算损失概率的程序代

\* 本图为上述代码运行后直接输出的结果。本书中为了使程序代码简洁明了、重点突出,均省略了对图片、字体等进行修饰的代码,后文类似情况不再赘述。