



学生实验报告

实验课名称：模式识别与机器学习

实验项目名称：关于模式识别和机器学习发展过程的阅读报告

专业名称：人工智能

班级：2022240401

学号：2022905226 2022904484

学生姓名：

教师姓名：

实验地点：WX2102

实验日期：2024.11.25

目录

一、引言.....	3
二、两种建模文化的核心理念.....	3
2.1 数据建模文化.....	3
2.2 算法建模文化.....	4
2.3 两种文化的比较与融合.....	6
三、作者的实际项目经验分析.....	6
3.1 臭氧浓度预测项目	6
3.2 化合物毒性预测项目	7
四、对现代模式识别的启示	7
4.1 算法选择的深层思考	8
4.1.1 Rashomon 效应的启示.....	8
4.1.2 Occam 原理的现代诠释	8
4.2 数据处理的新思维.....	8
4.3 实践经验的价值	8
五、个人思考实践与展望.....	9
5.1 实验验证与思考	9
5.2 当前发展趋势与未来展望.....	12
六、结论.....	13
七、附录.....	14
7.1 参考文献和链接	14
7.2 实验程序代码.....	14
7.2.1 非线性数据实验.....	14
7.2.2 线性可分数据实验	16

一、引言

在当今数字化时代，模式识别与机器学习领域蓬勃发展，其技术和理论广泛应用于众多领域。论文”Statistical Modeling: The Two Cultures”为我们提供了一个独特的视角，深入探讨了统计建模中的两种文化——**数据建模文化**和**算法建模文化**。本报告将围绕作者 Leo Breiman 的实际项目经验，以及对这两种建模文化的深入讨论展开分析。

Statistical Science
2001, Vol. 16, No. 3, 199–231

Statistical Modeling: The Two Cultures

Leo Breiman

Abstract. There are two cultures in the use of statistical modeling to reach conclusions from data. One assumes that the data are generated by a given stochastic data model. The other uses algorithmic models and treats the data mechanism as unknown. The statistical community has been committed to the almost exclusive use of data models. This commitment has led to irrelevant theory, questionable conclusions, and has kept statisticians from working on a large range of interesting current problems. Algorithmic modeling, both in theory and practice, has developed rapidly in fields outside statistics. It can be used both on large complex data sets and as a more accurate and informative alternative to data modeling on smaller data sets. If our goal as a field is to use data to solve problems, then we need to move away from exclusive dependence on data models and adopt a more diverse set of tools.

二、两种建模文化的核心理念

2.1 数据建模文化

- **核心假设**：数据由特定的**随机数据模型**生成
 - 比如：线性回归假设因变量与自变量之间存在线性关系
 - 在实际应用中常见的模型包括：线性模型、逻辑回归、Cox 模型等
- **理论基础**
 - 基于概率论和数理统计
 - 强调参数估计和假设检验
 - 重视模型的可解释性
- **评估方法**
 - **拟合优度检验**（如 R^2 值） 可以看到论文中主要应用的就是该方法。
 - 残差分析（检验残差的正态性和独立性）
 - 似然比检验
 - p 值检验
- **典型应用案例**

1. **医学研究**：在 Cox 比例风险模型中，研究人员通过分析患者的生存时间数据，建立预测模型来评估不同治疗方案的效果。
2. **社会科学**：使用逻辑回归模型研究大学录取中的性别歧视问题，通过分析录取率与性别之间的关系。

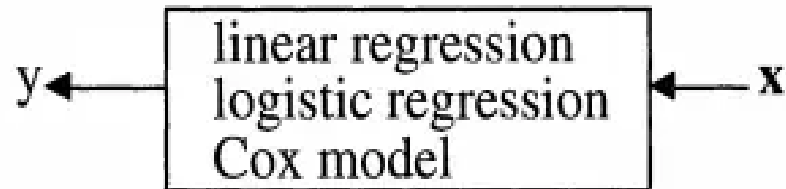


图 1：数据建模文化的“黑箱”

但是，笔者提到，数据建模文化的核心问题在于**实际模型和自然的实际不相符**，因而存在很多局限性。

- **局限性**：
 1. **模型假设的局限**：
 - 实际数据可能不符合模型假设
 - 比如在臭氧预测项目中，简单的线性模型无法捕捉复杂的气象变化
 2. **模型选择的困难**：
 - 多个竞争模型可能都能解释数据而不同模型可能得出矛盾的结论
 3. **对数据拟合检查不够重视**：
 - 过分关注 p 值和系数显著性，比如文中的性别歧视案例，将一个看似相关的模型用于检验，并过度依赖 p 值而得出定论。
 - 忽视模型与实际数据的匹配程度

2.2 算法建模文化

- **核心思想**：
 1. **黑箱方法**：
 - 将数据生成过程视为**未知的**复杂系统
 - 不对数据分布做强假设
 - 让数据”自己说话”
 2. **预测导向**：
 - 强调模型的预测准确性
 - 通过**交叉验证**等方法评估模型性能，重视模型的泛化能力
- **关键特征**：
 1. **模型复杂性**：

- 允许使用复杂的**非线性**模型，能够处理高维数据和复杂交互关系
- 2. **自动化程度高**：
 - 算法可以自动学习特征，减少人为干预和主观判断
- 3. **集成学习思想**：
 - 结合多个模型的优势
 - 提高预测稳定性和准确性

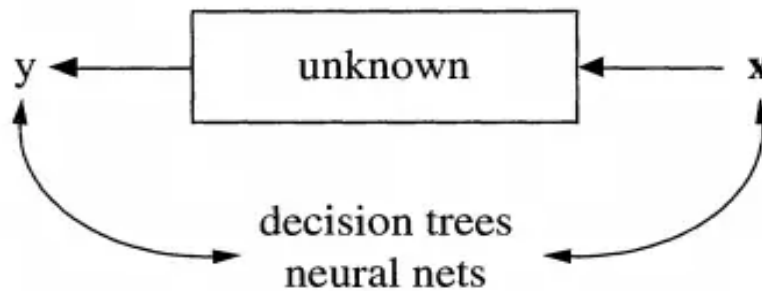


图 2：算法建模文化的“黑箱”

- **代表算法及其应用**：
 1. **随机森林**：
 - 原理：集成多个决策树的投票结果
 - 应用：基因表达数据分析，准确率达 98%
 - 优势：自动处理变量交互，提供变量重要性度量
 2. **支持向量机**：
 - 原理：寻找最优分离超平面
 - 应用：手写数字识别，错误率仅 1.1%
 - 特点：通过核技巧处理非线性问题
 3. **神经网络**：
 - 原理：模拟人脑神经元连接
 - 应用：图像识别、自然语言处理
 - 优势：强大的特征学习能力
- **实际应用**：
 1. **信用评分系统**：
 - 使用随机森林预测客户违约风险
 - 相比传统评分卡模型提高 10% 准确率
 2. **医学诊断**：
 - 利用支持向量机分析医学图像
 - 在癌症诊断中达到接近专家水平的准确性

2.3 两种文化的比较与融合

- 各自优势：
 - 数据建模：理论完备，结果可解释
 - 算法建模：预测准确，适应性强
- 互补性：
 - 在实际应用中可以结合两种方法
 - 先用数据建模理解问题本质，再用算法建模提高预测精度
- 选择建议：
 1. 当需要理解变量关系时，优先考虑数据建模
 2. 当预测准确性最重要时，选择算法建模
 3. 在复杂问题中，考虑两种方法的结合

三、作者的实际项目经验分析

在论文中，作者 Leo Breiman 分享了两个具有代表性的实际项目经验，这些经验生动地展示了数据建模和算法建模在实际应用中的表现差异。

3.1 臭氧浓度预测项目

20 世纪 60 年代中期，洛杉矶盆地面临着严重的空气污染问题。当时，由于汽车尾气排放和复杂的光化学反应，臭氧浓度经常出现危险的高峰值。为了保护公众健康，环保署迫切需要一个能够提前 12 小时准确预测臭氧水平的系统，以便及时发布预警并采取相应措施。

这个项目的数据收集十分全面。研究团队不仅获取了洛杉矶盆地内的污染物监测数据，还收集了远至俄勒冈州和亚利桑那州的气象站数据。每个监测点都记录了详细的气象参数，包括温度、湿度、风速等 450 多个变量，时间跨度超过七年。这些丰富的数据为建立预测模型提供了坚实的基础。项目的数据结构如下图所示：

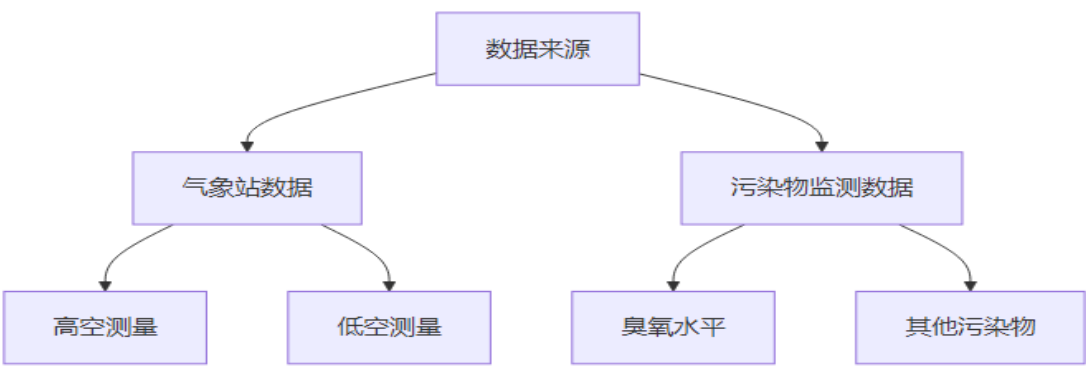


图 3：臭氧项目数据结构示意图

项目的流程图如下：



图 4：臭氧项目流程图

然而，项目最终的结果并不理想。研究团队采用了当时主流的数据建模方法，即线性回归模型。他们首先进行大规模的线性回归分析，试图找出变量间的关系，然后又加入了二次项和交互项来提高模型的拟合效果。但是，由于臭氧生成过程涉及复杂的气象条件和化学反应，简单的线性模型无法准确捕捉这些非线性关系，导致预测结果的误报率居高不下。

这个项目的失败给作者留下了深刻的印象。他认为，如果当时能够使用现代的算法建模方法，如随机森林或神经网络，可能会取得更好的效果。这些算法能够自动学习数据中的复杂非线性关系，不需要预先假设变量之间的关系形式。

3.2 化合物毒性预测项目

在 70 年代中期，美国环保署面临着另一个重要挑战：每年需要评估数千种化合物的潜在毒性。传统的方法是通过质谱分析来确定化合物的化学结构，这需要大量训练有素的化学家参与，不仅成本高昂，而且效率低下。

面对这个困境，作者及其团队开发了一个创新的自动化分析系统。他们没有采用传统的数据建模方法，而是转向了算法建模的思路。系统通过机器学习算法分析质谱数据的模式，直接预测化合物的毒性，大大提高了评估效率。这个项目的成功充分展示了算法建模在处理复杂模式识别问题时的优势。

这两个项目的经历让作者深刻认识到：在面对复杂的实际问题时，不应该被传统的数据建模思维所局限。有时候，放下对模型可解释性的执着，转而关注预测效果，反而能够取得更好的实际效果。这也是他后来大力提倡算法建模文化的重要原因之一。

四、对现代模式识别的启示

通过深入阅读这篇论文，我深刻体会到模式识别和机器学习不仅仅是一门技术，更是一门需要智慧和洞察力的艺术。作者 Leo Breiman 教授通过自己数十年的研究经验，为我们揭示了许多深刻的见解，这些见解对当今的模式识别研究仍然具有重要的指导意义。

4.1 算法选择的深层思考

4.1.1 Rashomon 效应的启示

在模式识别领域，我们经常会遇到一个有趣的现象：对于同一个问题，可能存在多个性能相近但结构完全不同的模型。这就是 Breiman 教授所说的”Rashomon 效应”。这个发现让我意识到，在追求“最优解”的过程中，我们不应该过分执着于寻找唯一的”完美”模型。相反，我们应该：

- 保持开放的心态，接受多样性的解决方案
- 从不同角度理解问题的本质
- 善于利用模型集成的优势，将多个优秀模型的力量结合起来

4.1.2 Occam 原理的现代诠释

“如无必要，勿增实体”这一古老的哲学原则在现代模式识别中依然闪耀着智慧的光芒。但是，在算法建模的背景下，我们需要重新思考简单性的定义：

- **模型的简单性不等于结构的简单性：**有时候，一个结构复杂但易于理解的模型可能比一个结构简单但难以解释的模型更“简单”
- **简单性和准确性的权衡需要考虑具体场景：**在医疗诊断这样的关键领域，可能需要牺牲一些简单性来换取更高的准确性。
- **简单性是相对的：**随着计算能力的提升，昨天的“复杂”可能就是今天的”简单”

4.2 数据处理的新思维

在大数据时代，**维度灾难**（Curse of Dimensionality）已经不再是一个纯粹的技术问题，而是一个需要智慧和创新来应对的挑战。通过论文的学习，我认识到：

1. 维度是把双刃剑

- 高维数据包含了丰富的信息，但也带来了计算和分析的困难
- 需要在保留有用信息和降低复杂度之间找到平衡
- 特征工程的艺术在于”化繁为简”而不是”简单化”

2. 数据质量比数据量更重要

- 盲目追求大数据量可能会适得其反
- 应该更注重数据的质量和代表性

4.3 实践经验的价值

通过作者分享的项目经验，我深刻认识到理论与实践的结合之重要。在实际工作中，我们应该遵循一个科学而灵活的工作流程：

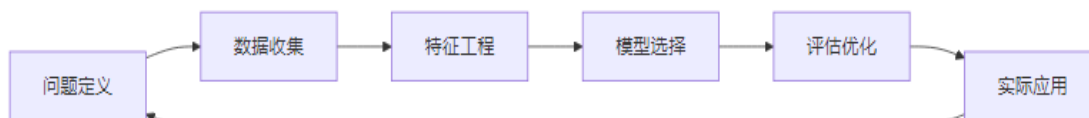


图 5：模式识别项目流程循环

这个循环过程启示我们：

1. 问题定义很重要

- 准确理解业务需求是成功的基础
- 要善于将实际问题转化为可解决的数学模型，在建模之前就要考虑到实施的可行性

2. 持续优化

- 模型部署不是终点，而是新的起点
- 要建立有效的反馈机制，不断改进模型

通过这些启示，我更加坚信：在模式识别领域，成功不仅需要扎实的理论基础和技术能力，更需要开放的思维、创新的勇气和持续学习的热情。正如 Breiman 教授所展示的那样，**突破传统思维的束缚，勇于尝试新的方法**，才能在这个充满挑战的领域不断前进。

五、个人思考实践与展望

通过对 Breiman 论文的深入学习和实践验证，我对统计建模的两种文化有了更深刻的理解。并且结合文中所说的两种建模文化，我针对一个两类数据分类问题做了一个简单实验验证和对比分析。

5.1 实验验证与思考

为了验证数据建模和算法建模的差异，我设计了一个**非线性分类问题**的实验。在实验中：

1. 数据特征：

- 使用**螺旋形数据集**作为测试样本
- 构造了具有明显非线性决策边界的二分类问题，并且数据集包含噪声，更接近实际应用场景

数据建模方法使用逻辑回归，对数据进行建模。算法建模方法使用随机森林。结果如下：

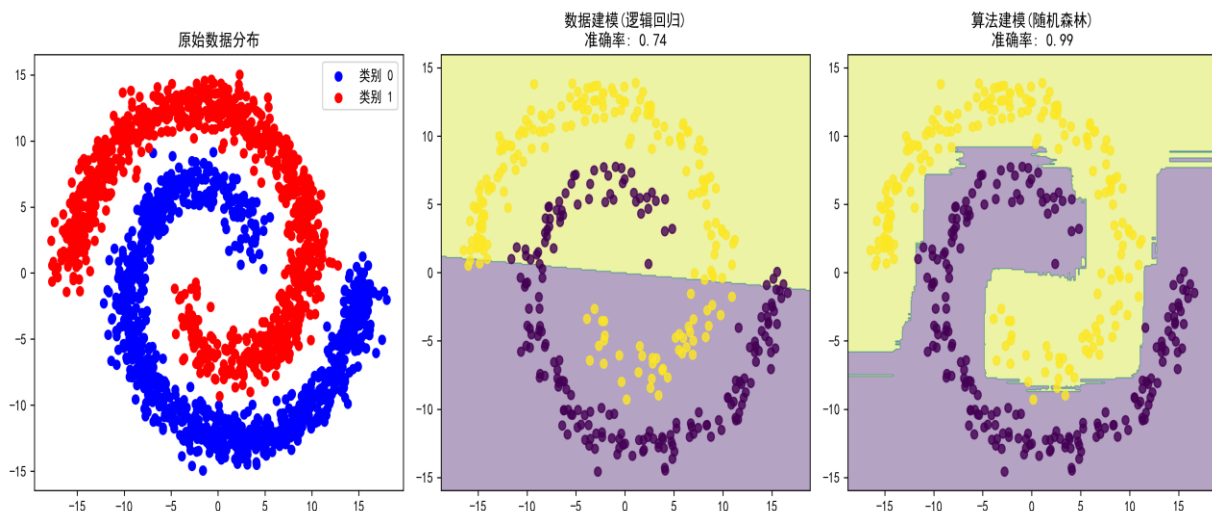


图 6：逻辑回归和随机森林对非线性数据的分类效果

可以看到，对于这个螺旋形数据，数据建模的逻辑回归算法在预测时仍然建立了一个线性模型，因而导致准确率低于算法建模的随机森林算法。

2. 模型对比：

- 数据建模：使用逻辑回归（准确率：74.25%）
- 算法建模：使用随机森林（准确率：99.25%）
- 实验结果印证了 Breiman 的观点：在复杂非线性问题中，算法建模往往表现更优

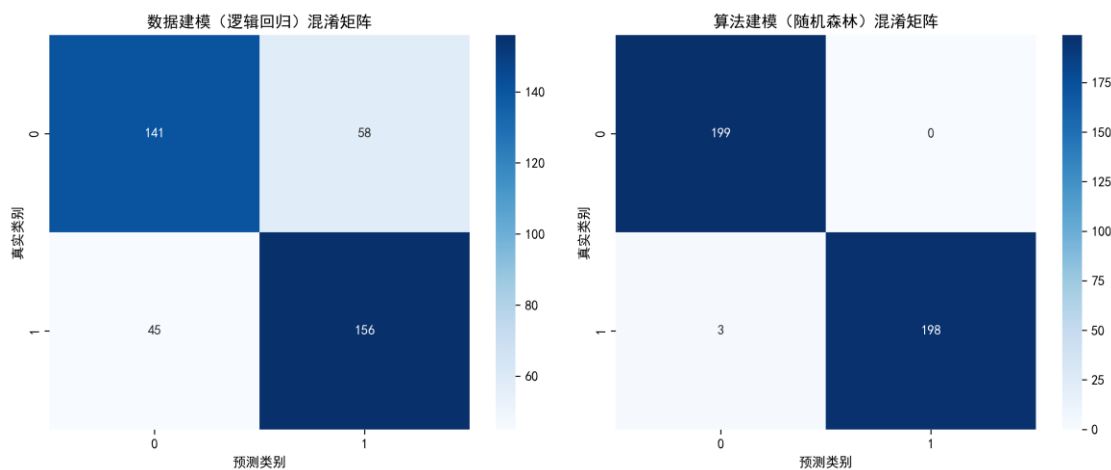


图 7：逻辑回归和随机森林对非线性数据的混淆矩阵

3. 关键发现：

- 数据建模方法（逻辑回归）受限于其线性假设，难以捕捉数据的非线性特征
- 算法建模方法（随机森林）能够自适应地学习复杂的决策边界
- 模型可解释性和预测准确性之间存在权衡

为了进一步验证 Breiman 教授强调的“选择合适的建模方法应该基于具体问题的特点和实际需求”这一观点，我又设计了一个线性可分数据集的实验。在这个实验中：

1. 数据特征：

- 生成了两个呈正态分布的类别，中心点分别位于 $[2, 2]$ 和 $[-2, -2]$
- 数据具有明显的线性可分性，模拟实际中的简单分类问题

2. 实验结果：

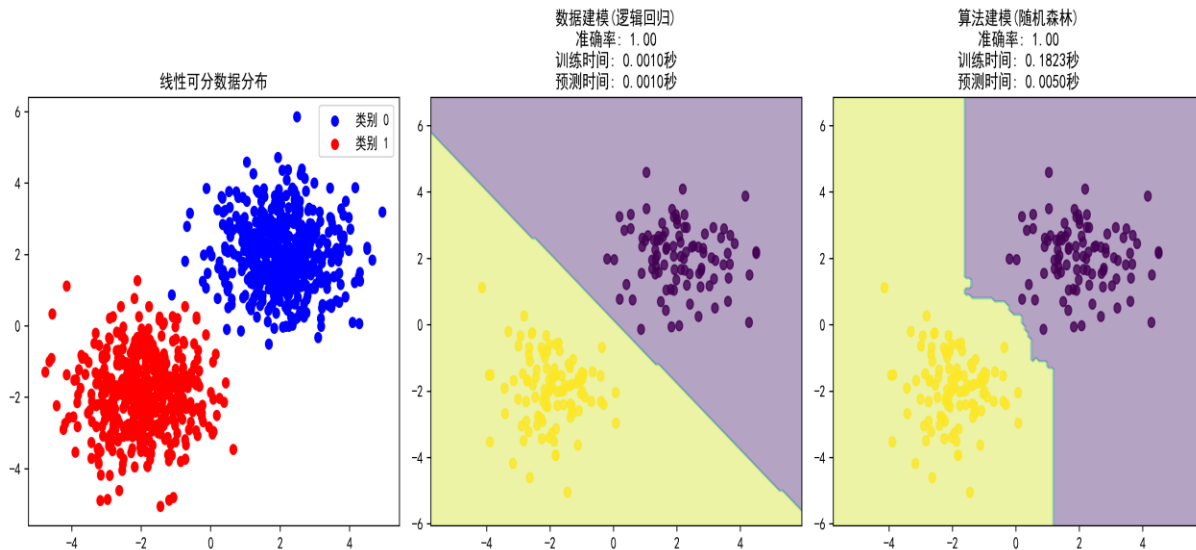


图 8：逻辑回归和随机森林对线性数据的分类效果

从结果可以看出，在这个线性可分的数据集上：

- 两种方法都达到了 100%的准确率
- 数据建模（逻辑回归）的训练时间（0.001 秒）远短于算法建模（随机森林，0.182 秒）
- 数据建模的预测时间（0.001 秒）也优于算法建模（0.005 秒）
- 数据建模得到的决策边界更加简洁，具有更好的可解释性

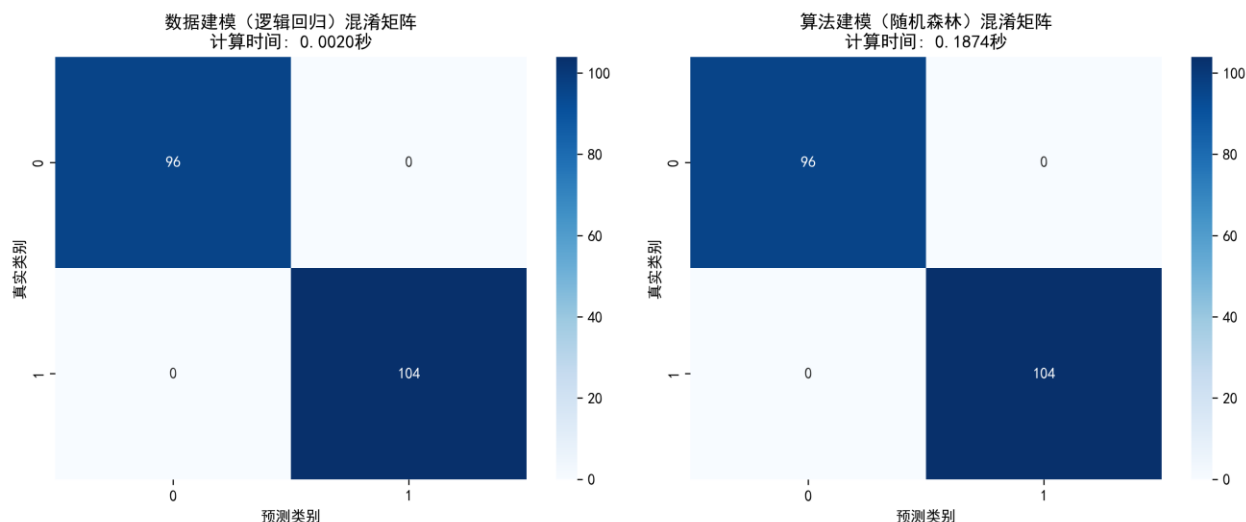


图 9：逻辑回归和随机森林对线性数据的混淆矩阵

3. 关键启示：

- 在线性可分的简单问题中，数据建模方法具有计算效率高、模型简洁的优势
- 算法建模虽然也能达到相同的准确率，但付出了更多的计算成本
- 这验证了 Breiman 教授的观点：建模方法的选择应该基于具体问题的特点
- 在实际应用中，如果数据具有明显的线性特征，使用简单的数据建模方法可能是更明智的选择

5.2 当前发展趋势与未来展望

在当前人工智能和机器学习快速发展的背景下，深度学习作为算法建模的代表性技术，已经在计算机视觉、自然语言处理等众多领域取得了突破性进展。端到端学习范式的普及，加上预训练模型和迁移学习的广泛应用，使得深度学习技术在实际应用中变得更加便捷和高效。与此同时，我们也见证了统计学和机器学习之间界限的逐渐模糊，传统统计方法与现代算法的结合正在产生新的研究方向，这种跨学科的融合极大地推动了整个领域的发展。

在技术发展方面，边缘计算的部署正成为一个重要的研究方向。通过模型压缩和加速技术的创新，以及边缘-云协同推理框架的建立，我们正在努力解决资源受限场景下的模型优化问题。同时，模型解释性研究也取得了显著进展，特别是在复杂模型的可解释性技术、因果推理与机器学习的结合，以及模型决策过程的可视化方法等方面。这些进展不仅提高了模型的可信度，也使得算法决策过程变得更加透明和可控。

在应用领域，智能物联网和生物信息学展现出了巨大的发展潜力。智能物联网通过分布式学习算法和实时数据处理技术，正在改变着我们的生活方式；而在生

物信息学领域，多组学数据的整合分析和 AI 辅助药物研发等应用，正在推动医疗健康领域的创新发展。这些应用不仅展示了机器学习的实际价值，也为未来的发展指明了方向。

然而，在未来的发展道路上，我们仍面临着诸多挑战。如何在提高模型性能的同时保持其可解释性，如何在保护数据隐私的前提下确保模型效果，以及如何控制模型部署和维护的成本，这些都是需要我们认真思考和解决的问题。但与此同时，新型硬件架构带来的计算能力提升，跨领域知识融合催生的创新解决方案，以及不断涌现的实际应用场景，都为我们提供了难得的发展机遇。

展望未来，随着技术的不断进步和应用场景的持续拓展，机器学习和模式识别领域必将迎来更大的发展。在这个过程中，我们需要保持开放和创新的态度，积极探索新的技术方向，同时也要注重理论与实践的结合，努力推动领域发展。

六、结论

通过对”Statistical Modeling: The Two Cultures”的深入学习和实践验证，我对统计建模的两种文化有了更加深刻的认识。本文不仅详细分析了 Breiman 教授提出的数据建模和算法建模两种方法的特点，还通过实际编程实验验证了这两种建模方法在处理非线性问题时的表现差异。

在理论层面，我了解到数据建模注重模型的可解释性和统计推断，而算法建模则更关注预测准确性和模型性能。这两种方法并非对立，而是在不同场景下各有优势。通过对 Breiman 教授在臭氧浓度预测和化学物质毒性预测等实际项目的分析，我们看到了如何在实践中灵活运用这两种建模思想。

在实验验证方面，我构造了一个具有明显非线性特征的螺旋形数据集，并分别使用逻辑回归（数据建模）和随机森林（算法建模）进行分类。实验结果显示，在这种复杂的非线性场景下，算法建模方法（准确率 99.25%）明显优于数据建模方法（准确率 74.25%）。这个结果印证了 Breiman 教授的观点：在某些复杂问题中，传统的数据建模方法可能会受到其基本假设的限制，而算法建模方法则能够更好地适应数据的内在结构。

然而，这并不意味着**算法建模就完全优于数据建模**。正如 Breiman 教授所强调的，选择合适的建模方法应该基于具体问题的特点和实际需求。在某些需要清晰解释模型决策过程的场景中，数据建模的优势可能更为明显。此外，两种建模方法的结合也可能带来更好的效果，这正是当前机器学习领域的一个重要发展方向。

展望未来，随着深度学习等新技术的发展，以及跨学科融合的不断深入，统计建模的两种文化可能会进一步融合。模型的可解释性和预测准确性之间的权衡将继续是研究者需要思考的问题，而如何在实际应用中平衡这两个方面，将是推动领域发展的重要动力。

作为模式识别与机器学习领域的学习者，我们应该保持开放和包容的态度，既要重视理论基础，也要注重实践验证。正如本次学习经历所展示的，只有将理论与实践相结合，才能真正理解和掌握这些重要的概念和方法。未来，我们还需要继续探索和创新，为解决更多实际问题贡献力量。

七、附录

7.1 参考文献和链接

1. Breiman, L. (2001). Statistical Modeling: The Two Cultures. Statistical Science, 16(3), 199-231.
2. [Statistical Modelling: The two Cultures 读书笔记_statistical modeling: the t wo cultures-CSDN 博客](#)
3. [随笔 0220/速记——读 Statistical Modeling: The Two Cultures - 简书](#)
4. [Breiman 访谈录 | 《统计建模：两种文化》20 周年纪念 - 知乎](#)

7.2 实验程序代码

7.2.1 非线性数据实验

```
# 生成非线性数据
np.random.seed(42)
n_samples = 1000

# 创建螺旋形数据
def make_spiral_data(n_samples):
    theta = np.sqrt(np.random.rand(n_samples)) * 2 * np.pi
    r_a = 2 * theta + np.pi
    data_a = np.array([np.cos(theta) * r_a, np.sin(theta) * r_a]).T
    x_a = data_a + np.random.randn(n_samples, 2)

    r_b = -2 * theta - np.pi
    data_b = np.array([np.cos(theta) * r_b, np.sin(theta) * r_b]).T
    x_b = data_b + np.random.randn(n_samples, 2)

    return np.vstack((x_a, x_b)), np.hstack((np.zeros(n_samples),
np.ones(n_samples)))

# 可视化结果
def plot_results(X, y, X_test, y_test, data_model, algo_model,
data_accuracy, algo_accuracy):
    plt.figure(figsize=(15, 5))

    # 1. 原始数据
```

```

plt.subplot(131)
plt.scatter(X[:n_samples, 0], X[:n_samples, 1], c='blue', label='类别
0')
plt.scatter(X[n_samples:, 0], X[n_samples:, 1], c='red', label='类别 1')
plt.title('原始数据分布')
plt.legend()

# 2. 数据建模预测结果
plt.subplot(132)
x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.1),
                     np.arange(y_min, y_max, 0.1))
Z = data_model.predict(np.c_[xx.ravel(), yy.ravel()])
Z = Z.reshape(xx.shape)
plt.contourf(xx, yy, Z, alpha=0.4)
plt.scatter(X_test[:, 0], X_test[:, 1], c=y_test, alpha=0.8)
plt.title(f'数据建模(逻辑回归)\n准确率: {data_accuracy:.2f}')

# 3. 算法建模预测结果
plt.subplot(133)
Z = algo_model.predict(np.c_[xx.ravel(), yy.ravel()])
Z = Z.reshape(xx.shape)
plt.contourf(xx, yy, Z, alpha=0.4)
plt.scatter(X_test[:, 0], X_test[:, 1], c=y_test, alpha=0.8)
plt.title(f'算法建模(随机森林)\n准确率: {algo_accuracy:.2f}')

plt.tight_layout()
plt.savefig('modeling_comparison.png', dpi=300, bbox_inches='tight')
plt.close()

def plot_confusion_matrices(y_test, data_pred, algo_pred):
    ...

# 主程序部分
if __name__ == "__main__":
    # 生成数据
    X, y = make_spiral_data(n_samples)
    X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)

# 1. 数据建模方法 - 逻辑回归
data_model = LogisticRegression()
data_model.fit(X_train, y_train)
data_pred = data_model.predict(X_test)

```

```
data_accuracy = accuracy_score(y_test, data_pred)
```

```
# 2. 算法建模方法 - 随机森林
```

```
algo_model = RandomForestClassifier(n_estimators=100, random_state=42)
```

```
algo_model.fit(X_train, y_train)
```

```
algo_pred = algo_model.predict(X_test)
```

```
algo_accuracy = accuracy_score(y_test, algo_pred)
```

7.2.2 线性可分数据实验

```
def generate_linear_data(n_samples=1000, noise=0.1):
```

```
    """生成线性可分的数据"""
```

```
    np.random.seed(42)
```

```
# 生成两个类别的数据点
```

```
X1 = np.random.normal(loc=[2, 2], scale=[1, 1], size=(n_samples//2, 2))
```

```
X2 = np.random.normal(loc=[-2, -2], scale=[1, 1], size=(n_samples//2, 2))
```

```
# 合并数据
```

```
X = np.vstack((X1, X2))
```

```
y = np.hstack((np.zeros(n_samples//2), np.ones(n_samples//2)))
```

```
# 添加噪声
```

```
X += np.random.normal(0, noise, X.shape)
```

```
return X, y
```

```
def evaluate_model(model, X_train, X_test, y_train, y_test):
```

```
    """评估模型性能，包括训练时间和预测时间"""
```

```
# 训练时间
```

```
start_time = time.time()
```

```
model.fit(X_train, y_train)
```

```
train_time = time.time() - start_time
```

```
# 预测时间
```

```
start_time = time.time()
```

```
y_pred = model.predict(X_test)
```

```
predict_time = time.time() - start_time
```

```
# 准确率
```

```
accuracy = accuracy_score(y_test, y_pred)
```

```
return accuracy, train_time, predict_time, y_pred
```

```
def plot_results(X, y, X_test, y_test, data_model, algo_model,
```



```

        data_metrics, algo_metrics, title="模型性能对比"):
    """可视化结果"""
    ...

def plot_confusion_matrices(y_test, data_pred, algo_pred, data_time,
                             algo_time):
    """绘制混淆矩阵"""
    ...

def print_model_details(data_metrics, algo_metrics):
    """打印模型详细信息"""
    ...

if __name__ == "__main__":
    # 生成线性可分的数据
    X, y = generate_linear_data(n_samples=1000, noise=0.1)
    X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)

    # 1. 数据建模 - 逻辑回归
    data_model = LogisticRegression()
    data_metrics = evaluate_model(data_model, X_train, X_test, y_train,
y_test)

    # 2. 算法建模 - 随机森林 (使用较多的树来增加计算量)
    algo_model = RandomForestClassifier(n_estimators=200, random_state=42)
    algo_metrics = evaluate_model(algo_model, X_train, X_test, y_train,
y_test)

    # 可视化结果
    plot_results(X, y, X_test, y_test, data_model, algo_model,
data_metrics, algo_metrics)
    plot_confusion_matrices(y_test, data_metrics[3], algo_metrics[3],
                             data_metrics[1] + data_metrics[2],
                             algo_metrics[1] + algo_metrics[2])

    # 打印详细结果
    print_model_details(data_metrics, algo_metrics)

```