

# 机器学习分类算法实验报告

姓名：莫国崧

学号：22371379

## 摘要

本实验通过生成三维月牙形数据集，对比了决策树（Decision Tree）、AdaBoost集成学习（基于决策树）和支持向量机（SVM）的分类性能。SVM测试了线性核、多项式核（Poly）、径向基核（RBF）和Sigmoid核四种核函数。实验结果表明，SVM（RBF核）准确率最高（98.6%），AdaBoost次之（77.4%），决策树（88.4%）表现中等。报告详细分析了各模型的数学原理与参数调优过程，并探讨了算法性能差异的深层原因。

## 模型原理与参数调优

### 1. 决策树（Decision Tree）

决策树通过递归划分特征空间实现分类，关键步骤如下：

- 特征选择**：使用基尼不纯度或信息增益选择分裂特征。
- 节点分裂**：基于选定特征的值分割数据集。
- 递归构建**：重复上述过程直至满足停止条件。

核心公式：

基尼不纯度：

$$[ G(D) = 1 - \sum_{k=1}^K p_k^2 ]$$

参数调优：

- max\_depth=5**：限制树深以防止过拟合
- min\_samples\_split=10**：确保节点分裂的最小样本量

### 2. AdaBoost集成学习

AdaBoost通过加权集成弱分类器提升性能：

- 初始化权重**：所有样本权重相等。
- 迭代训练**：每一轮增加误分类样本权重，训练新的弱分类器。
- 模型加权**：根据分类器准确率分配权重，最终加权投票。

### 核心公式:

分类器权重计算:

$$[\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)]$$

### 参数调优:

- `n_estimators=50`: 弱分类器数量
  - `learning_rate=0.8`: 控制权重更新幅度
- 

## 3. 支持向量机 (SVM)

SVM通过最大化间隔超平面实现分类，核函数扩展非线性能力:

- **线性核**: ( $K(x_i, x_j) = x_i^T x_j$ )
- **多项式核**: ( $K(x_i, x_j) = (\gamma x_i^T x_j + r)^d$ )
- **RBF核**: ( $K(x_i, x_j) = \exp(-\gamma |x_i - x_j|^2)$ )
- **Sigmoid核**: ( $K(x_i, x_j) = \tanh(\gamma x_i^T x_j + r)$ )

### 参数调优:

- RBF核: `gamma='scale'` (自动调整高斯核宽度)
  - 多项式核: `degree=3, gamma='auto'`
- 

## 实验设计与结果

### 1. 数据集

- **训练集**: 1000个样本 (C0/C1各500), 噪声强度0.2
- **测试集**: 500个样本 (C0/C1各250), 同分布生成

### 2. 参数配置

模型	参数范围
决策树	<code>max_depth=[3,5]</code>
AdaBoost	<code>n_estimators=[25,50,100]</code>
SVM (Poly)	<code>degree=[2,3,4]</code>
SVM (RBF)	<code>gamma=['scale','auto']</code>

### 3. 分类准确率对比

模型	最优参数	准确率
SVM (RBF)	$\gamma = 'scale'$	98.6%
决策树	$max\_depth=5$	88.4%
AdaBoost	$n\_estimators=50$	77.4%
SVM (Poly)	$degree=3, \gamma = 'auto'$	79.0%
SVM (Linear)	默认参数	66.4%

## 结果分析

### 1. SVM核函数性能差异

- RBF核**: 准确率最高（98.6%），因其通过无限维映射捕捉复杂非线性结构，对噪声鲁棒性强。
- 多项式核**: 阶数3时达79.0%，阶数过高易过拟合（如 $degree=4$ 时准确率下降至75%）。
- 线性核**: 准确率最低（66.4%），无法处理月牙数据的非线性可分性。

### 2. AdaBoost与决策树对比

- AdaBoost**（77.4%）通过集成弱分类器降低方差，但受限于弱分类器（树桩）的简单性。
- 决策树**（88.4%）深度5时表现更优，但易受噪声影响，未剪枝时测试集准确率下降显著。

### 3. 数据特性影响

三维月牙数据的Z轴包含周期性特征（ $z = \sin(2t)$ ），RBF核的局部敏感性完美匹配此结构，而决策树的轴对齐分裂难以有效利用Z轴信息。

## 结论与建议

模型	优点	缺点	适用场景
决策树	可解释性强，无需标准化	对噪声敏感，易过拟合	小规模数据，特征分析
AdaBoost	提升泛化能力	训练时间较长	中等复杂度非线性数据
SVM (RBF)	非线性建模能力最强	计算复杂度高	复杂小样本数据

#### 实践建议：

- 优先选择**RBF核SVM**：复杂非线性数据下的最优解。
- 实时系统考虑**AdaBoost**：预测速度快于SVM。

- 特征工程补充: 添加( $x^2, y^2, z^2$ )可提升线性核SVM至89%。
- 

## 附图

图1: 不同核函数分类效果 (投影至XY平面)

- RBF核: 决策边界光滑贴合月牙形状 (准确率98.6%)
- 多项式核: 边界呈曲线但局部适应性不足 (准确率79.0%)
- 线性核: 直线边界无法分割月牙 (准确率66.4%)