

# 基于 K-means 算法的鸢尾花数据集聚类分析

## 数据挖掘实验报告

2024 年 12 月 18 日

### 摘要

本实验基于经典的鸢尾花（Iris）数据集，实现并分析了 K-means 聚类算法。通过对算法的实现、参数优化和结果可视化，深入研究了 K-means 算法在实际数据集上的应用效果。实验结果表明 K-means 算法能够有效地对鸢尾花数据进行聚类分析，并取得了良好的聚类效果。

## 目录

## 1 引言

### 1.1 研究背景

聚类分析是数据挖掘和机器学习中的重要任务之一，其目的是将相似的数据对象划分到同一个簇中，而将不相似的对象划分到不同簇中。K-means 算法作为最经典的聚类算法之一，因其简单、高效的特点而被广泛应用。

### 1.2 研究目的

本实验旨在：

- 实现 K-means 聚类算法
- 分析算法在鸢尾花数据集上的表现
- 通过可视化手段理解聚类结果
- 评估聚类效果并进行参数优化

## 2 理论基础

### 2.1 K-means 算法原理

K-means 算法的基本思想是通过迭代方式寻找 K 个簇的一种划分方案，使得聚类结果对应的代价函数最小。其主要步骤如下：

1. 随机选择 K 个点作为初始聚类中心
2. 计算每个数据点到各个聚类中心的距离，将其划分到最近的聚类中心所对应的簇
3. 重新计算每个簇的中心点（计算簇内所有点的均值）
4. 重复步骤 2 和 3，直到聚类中心不再发生变化或达到最大迭代次数

### 2.2 评估指标

本实验使用以下指标评估聚类效果：

- SSE（簇内误差平方和）：评估簇内的紧密度
- 轮廓系数：评估簇的分离度和紧密度
- 聚类准确率：与真实标签比较的正确率

## 3 实验设计

### 3.1 数据集说明

鸢尾花数据集包含 150 个样本，每个样本有 4 个特征：

- 萼片长度（Sepal Length）
- 萼片宽度（Sepal Width）
- 花瓣长度（Petal Length）
- 花瓣宽度（Petal Width）

数据集包含三个品种的鸢尾花：Setosa、Versicolor 和 Virginica，每类 50 个样本。

### 3.2 实验环境

- 编程语言：Python 3.6+
- 主要库：NumPy, Pandas, Matplotlib, Scikit-learn
- 开发环境：Visual Studio Code

### 3.3 实验流程

#### 1. 数据预处理

- 数据加载和清洗
- 特征标准化
- 数据可视化分析

#### 2. 算法实现

- 实现 K-means 核心算法
- 实现评估指标计算
- 实现结果可视化

#### 3. 参数优化

- 使用肘部法则确定最优 K 值
- 分析不同 K 值对聚类效果的影响

## 4 实验结果与分析

### 4.1 数据分布分析

通过特征分布图和散点矩阵，我们可以观察到：

- 不同特征的分布特征
- 特征之间的相关性
- 数据的可分性

### 4.2 聚类结果分析

#### 4.2.1 最优 K 值选择

通过肘部曲线和轮廓系数分析，我们发现：

- K=3 时，SSE 下降趋势变缓
- K=3 时，轮廓系数达到较高值
- 这与数据集中实际的三个品种数量相符

### 4.2.2 聚类效果评估

在  $K=3$  的情况下：

- 聚类准确率达到较高水平
- 簇内紧密度良好
- 簇间分离度明显

## 4.3 算法性能分析

- 收敛速度快，一般在 10 次迭代内收敛
- 计算复杂度适中，适合中等规模数据集
- 对初始中心点的选择较为敏感

# 5 结论与展望

## 5.1 主要结论

- K-means 算法在鸢尾花数据集上表现 好
- 算法能够有效识别数据集中的自然聚类
- 通过参数优化可以进一步提升聚类效果

## 5.2 改进方向

- 优化初始中心点的选择策略
- 考虑使用其他距离度量方式
- 结合其他算法进行集成学习

# 6 参考文献

1. MacQueen, J. (1967). Some methods for classification and analysis of multivariate observations.
2. Fisher, R. A. (1936). The use of multiple measurements in taxonomic problems.
3. Scikit-learn: Machine Learning in Python, Pedregosa et al., JMLR 12, pp. 2825-2830, 2011.

## A 代码实现

### A.1 K-means 核心实现

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.metrics import silhouette_score
4 from data_visualization import load_data
5
6 class KMeans:
7     def __init__(self, n_clusters=3, max_iters=300):
8         self.n_clusters = n_clusters
9         self.max_iters = max_iters
10        self.centroids = None
11        self.labels = None
12
13    def fit(self, X):
14        # 随机初始化聚类中心
15        idx = np.random.choice(len(X), self.n_clusters, replace=False)
16        self.centroids = X[idx]
17
18        for _ in range(self.max_iters):
19            old_centroids = self.centroids.copy()
20
21            # 计算每个样本到各个聚类中心的距离
22            distances = np.sqrt(((X - self.centroids[:, np.newaxis])
23                                **2).sum(axis=2))
24            self.labels = np.argmin(distances, axis=0)
25
26            # 更新聚类中心
27            for k in range(self.n_clusters):
28                if sum(self.labels == k) > 0:
29                    self.centroids[k] = X[self.labels == k].mean(axis=0)
30
31            # 检查是否收敛
32            if np.all(old_centroids == self.centroids):
33                break
```

```

34         return self
35
36     def predict(self, X):
37         distances = np.sqrt(((X - self.centroids[:, np.newaxis])**2).
38                               sum(axis=2))
39         return np.argmin(distances, axis=0)
40
41 def calculate_sse(X, kmeans):
42     """计算SSE（簇内误差平方和）"""
43     sse = 0
44     for k in range(kmeans.n_clusters):
45         cluster_points = X[kmeans.labels == k]
46         centroid = kmeans.centroids[k]
47         sse += np.sum((cluster_points - centroid) ** 2)
48     return sse
49
50 def plot_elbow_curve(X):
51     """绘制肘部曲线"""
52     sse_values = []
53     silhouette_values = []
54     k_range = range(2, 8)
55
56     for k in k_range:
57         kmeans = KMeans(n_clusters=k)
58         kmeans.fit(X)
59         sse = calculate_sse(X, kmeans)
60         silhouette = silhouette_score(X, kmeans.labels)
61         sse_values.append(sse)
62         silhouette_values.append(silhouette)
63
64     # 绘制SSE曲线
65     plt.figure(figsize=(12, 5))
66     plt.subplot(1, 2, 1)
67     plt.plot(k_range, sse_values, 'bo-')
68     plt.xlabel('聚类数量 (k)')
69     plt.ylabel('SSE')
70     plt.title('肘部曲线')
71
72     # 绘制轮廓系数曲线
73     plt.subplot(1, 2, 2)

```

```

73 plt.plot(k_range, silhouette_values, 'ro-')
74 plt.xlabel('聚类数量 (k)')
75 plt.ylabel('轮廓系数')
76 plt.title('轮廓系数曲线')
77
78 plt.tight_layout()
79 plt.savefig('kmeans_evaluation.png')
80 plt.close()
81
82 def plot_clusters(X, kmeans, title='K-means 聚类结果'):
83     """绘制聚类结果"""
84     plt.figure(figsize=(10, 8))
85     scatter = plt.scatter(X[:, 0], X[:, 1], c=kmeans.labels, cmap='
        viridis')
86     plt.scatter(kmeans.centroids[:, 0], kmeans.centroids[:, 1],
87                 c='red', marker='x', s=200, linewidths=3, label='聚类
        中心')
88     plt.colorbar(scatter)
89     plt.title(title)
90     plt.xlabel('Sepal Length')
91     plt.ylabel('Petal Length')
92     plt.legend()
93     plt.savefig('kmeans_clusters.png')
94     plt.close()
95
96 if __name__ == '__main__':
97     # 加载数据
98     X_scaled, y, _ = load_data()
99
100    # 绘制肘部曲线和轮廓系数
101    plot_elbow_curve(X_scaled)
102
103    # 使用最优k值进行聚类
104    kmeans = KMeans(n_clusters=3)
105    kmeans.fit(X_scaled)
106
107    # 绘制聚类结果
108    plot_clusters(X_scaled, kmeans)
109
110    # 计算聚类准确率

```

```

111     correct = sum(kmeans.labels == y)
112     accuracy = correct / len(y) * 100
113     print(f"聚类准确率: {accuracy:.2f}%")

```

Listing 1: K-means 算法实现

## A.2 数据可视化实现

```

1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  from sklearn.preprocessing import StandardScaler
6
7  def load_data():
8      """加载和预处理数据"""
9      try:
10         # 加载数据, 跳过第一行, 使用第二行作为列名
11         data = pd.read_csv('iris_train_shuzhi.csv', skiprows=[0])
12         print("原始数据形状:", data.shape)
13         print("原始数据前几行:\n", data.head())
14
15         # 删除第一列 (索引列)
16         data = data.iloc[:, 1:]
17
18         # 检查是否有缺失值
19         if data.isnull().any().any():
20             print("警告: 数据中存在缺失值, 将进行处理")
21             data = data.dropna()
22
23         # 分离特征和标签
24         X = data.iloc[:, :-1].astype(float) # 特征
25         y = data.iloc[:, -1].astype(int)    # 标签
26
27         print("\n数据统计信息:")
28         print(X.describe())
29
30         return X.values, y.values, X.values
31     except Exception as e:
32         print(f"数据加载错误: {str(e)}")

```



```

33         raise
34
35 def plot_features_distribution(X, y):
36     """绘制特征分布图"""
37     try:
38         features = ['Sepal Length', 'Sepal Width', 'Petal Length', '
39                     Petal Width']
40
41         plt.figure(figsize=(15, 10))
42         for i, feature in enumerate(features):
43             plt.subplot(2, 2, i+1)
44             for class_label in np.unique(y):
45                 mask = y == class_label
46                 plt.hist(X[mask, i],
47                         bins=20,
48                         alpha=0.5,
49                         label=f'Class {class_label}',
50                         density=True)
51             plt.title(f'{feature} Distribution')
52             plt.xlabel(feature)
53             plt.ylabel('Density')
54             plt.legend()
55
56         plt.tight_layout()
57         plt.savefig('feature_distribution.png')
58         print("特征分布图数据范围: ", np.ptp(X, axis=0))
59         plt.close()
60     except Exception as e:
61         print(f"特征分布图绘制错误: {str(e)}")
62         plt.close()
63
64 def plot_scatter_matrix(X, y):
65     """绘制特征散点矩阵"""
66     try:
67         plt.figure(figsize=(10, 8))
68
69         # 使用不同的颜色和标记绘制散点图
70         markers = ['o', 's', '^'] # 不同的标记样式
71         for i, class_label in enumerate(np.unique(y)):
72             mask = y == class_label

```

```

72         plt.scatter(X[mask, 0],
73                     X[mask, 2],
74                     alpha=0.6,
75                     marker=markers[i],
76                     s=50,
77                     label=f'Class {class_label}')
78
79     plt.xlabel('Sepal Length')
80     plt.ylabel('Petal Length')
81     plt.title('Sepal Length vs Petal Length')
82     plt.legend()
83     plt.grid(True, alpha=0.3)
84
85     plt.tight_layout()
86     plt.savefig('scatter_matrix.png')
87     print("散点图数据范围: ",
88           "\nSepal Length:", np.ptp(X[:, 0]),
89           "\nPetal Length:", np.ptp(X[:, 2]))
90     plt.close()
91 except Exception as e:
92     print(f"散点矩阵图绘制错误: {str(e)}")
93     plt.close()
94
95 def print_data_summary(X, y):
96     """打印数据摘要信息"""
97     try:
98         print("\n数据集摘要信息: ")
99         print(f"样本数量: {len(X)}")
100        print(f"特征数量: {X.shape[1]}")
101        print(f"类别数量: {len(np.unique(y))}")
102        print("\n各类别样本数量: ")
103        unique_labels, counts = np.unique(y, return_counts=True)
104        for label, count in zip(unique_labels, counts):
105            print(f"类别 {label}: {count}")
106
107        print("\n特征值范围: ")
108        for i, feature in enumerate(['Sepal Length', 'Sepal Width', 'Petal Length', 'Petal Width']):
109            print(f"{feature}: [{X[:, i].min():.2f}, {X[:, i].max():.2f}]")

```

```

110     except Exception as e:
111         print(f"数据摘要信息生成错误: {str(e)}")
112
113 if __name__ == '__main__':
114     try:
115         # 设置中文显示
116         plt.rcParams['font.sans-serif'] = ['SimHei']
117         plt.rcParams['axes.unicode_minus'] = False
118
119         # 加载数据
120         X_scaled, y, X_original = load_data()
121
122         # 打印数据摘要
123         print_data_summary(X_original, y)
124
125         # 绘制特征分布图
126         plot_features_distribution(X_original, y)
127         print("特征分布图已保存为 'feature_distribution.png'")
128
129         # 绘制散点矩阵
130         plot_scatter_matrix(X_original, y)
131         print("散点矩阵图已保存为 'scatter_matrix.png'")
132
133     except Exception as e:
134         print(f"程序执行出错: {str(e)}")

```

Listing 2: 数据可视化实现