

一. (30 points) 随机森林的原理

集成学习是一种通用技术，通过将多个不同模型的预测结果结合起来，以平均值或多数投票的方式生成单一预测，从而有效应对过拟合问题。

1. 考虑一组互不相关的随机变量 $\{Y_i\}_{i=1}^n$ ，其均值为 μ ，方差为 σ^2 。请计算这些随机变量平均值的期望和方差，给出计算过程。提示：在集成方法的背景下，这些 Y_i 类似于分类器 i 所作的预测。(5 points)

解：

设这组随机变量的平均值为 \bar{Y} ，即：

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$\mathbb{E}[\bar{Y}] = \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n Y_i\right) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[Y_i]$$

由于每个 Y_i 的期望均为 μ ，因此：

$$\mathbb{E}[\bar{Y}] = \frac{1}{n} \sum_{i=1}^n \mu = \frac{n\mu}{n} = \mu$$

对于常数 a 和随机变量 X ，有 $\text{Var}(aX) = a^2\text{Var}(X)$ 。

$$\text{Var}(\bar{Y}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n Y_i\right) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n Y_i\right)$$

由于 Y_i 互不相关，且方差的性质表明：

$$\text{Var}\left(\sum_{i=1}^n Y_i\right) = \sum_{i=1}^n \text{Var}(Y_i)$$

每个 Y_i 的方差均为 σ^2 ，因此：

$$\text{Var}\left(\sum_{i=1}^n Y_i\right) = \sum_{i=1}^n \sigma^2 = n\sigma^2$$

$$\text{Var}(\bar{Y}) = \frac{1}{n^2} \cdot n\sigma^2 = \frac{\sigma^2}{n}$$

这组随机变量平均值的期望为 μ ，方差为 $\frac{\sigma^2}{n}$ 。

由此我们可以得到，随着集成中模型数量 n 的增加，平均预测的方差会减小，从而提高预测的稳定性和准确性。

2. 在第 1 小问中，我们看到对于不相关的分类器，取平均可以有效减少方差。尽管现实中的预测不可能完全不相关，但降低决策树之间的相关性通常能减少最终方差。现在，重新考虑一组具有相关性的随机变量 $\{Z_i\}_{i=1}^n$ ，其均值为 μ ，方差为 σ^2 ，每个 $Z_i \in \mathbb{R}$ 为标量。假设对任意 $i \neq j$ ， $\text{Corr}(Z_i, Z_j) = \rho$ 。提示：如果你不记得相关性与协方差之间的关系，请回顾你的概率论等课程内容。

- 请计算随机变量 Z_i 平均值的方差，以 σ 、 ρ 和 n 为变量表示，给出计算过程。(5 points)
- 当 n 非常大时，会发生什么？这对于取平均的潜在有效性说明了什么？（……如果 ρ 很大 ($|\rho| \approx 1$) 会怎样？……如果 ρ 非常小 ($|\rho| \approx 0$) 又会怎样？……如果 ρ 处于中等水平 ($|\rho| \approx 0.5$) 呢?) 无需严格推导——基于你得出的方差公式，用定性分析进行讨论即可。(6 points)

解：

(1) 设这组随机变量的平均值为： $\bar{Z} = \frac{1}{n} \sum_{i=1}^n Z_i$

\bar{Z} 的方差： $\text{Var}(\bar{Z}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n Z_i\right)$

对于常数 a 和随机变量 X ，有 $\text{Var}(aX) = a^2 \text{Var}(X)$ ，因此：

$$\text{Var}(\bar{Z}) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n Z_i\right)$$

$\sum_{i=1}^n Z_i$ 的方差：

$$\text{Var}\left(\sum_{i=1}^n Z_i\right) = \sum_{i=1}^n \text{Var}(Z_i) + \sum_{i \neq j} \text{Cov}(Z_i, Z_j)$$

其中， $\text{Var}(Z_i) = \sigma^2$ ，而协方差 $\text{Cov}(Z_i, Z_j)$ 可以通过相关系数 ρ 表示为：

$$\text{Cov}(Z_i, Z_j) = \rho \sigma^2$$

$$\text{Var}\left(\sum_{i=1}^n Z_i\right) = n\sigma^2 + n(n-1)\rho\sigma^2$$

$$\text{Var}(\bar{Z}) = \frac{1}{n^2} (n\sigma^2 + n(n-1)\rho\sigma^2) = \frac{\sigma^2}{n} (1 + (n-1)\rho)$$

因此， \bar{Z} 的方差为：

$$\text{Var}(\bar{Z}) = \frac{\sigma^2}{n} (1 + (n-1)\rho)$$

(2) 当 n 趋于无穷大时， $\text{Var}(\bar{Z})$ 的行为取决于 ρ 的值：

- 如果 $\rho \approx 1$ ：此时， $\text{Var}(\bar{Z}) \approx \frac{\sigma^2}{n} \cdot n = \sigma^2$ ，即平均值的方差趋近于单个随机变量的方差，表明高度相关的随机变量的平均值并不能有效降低方差。
- 如果 $\rho \approx 0$ ：此时， $\text{Var}(\bar{Z}) \approx \frac{\sigma^2}{n}$ ，即平均值的方差随着 n 的增加而减小，表明不相关的随机变量的平均值可以有效降低方差。
- 如果 $\rho \approx 0.5$ ：此时， $\text{Var}(\bar{Z}) \approx \frac{\sigma^2}{n} \cdot (1 + 0.5n)$ ，当 n 增加时，方差趋近于 $0.5\sigma^2$ ，即平均值的方差仍然会减小，但减小的程度取决于相关性的大小。

综上所述，集成学习中，降低模型之间的相关性有助于减少平均预测的方差，从而提高预测的稳定性和准确性。

3. Bagging 是一种通过随机化从同一数据集生成多个不同学习器的方法。给定一个大小为 n 的训练集, Bagging 会通过有放回抽样生成 T 个随机子样本集, 每个子样本集大小为 n' 。在每个子样本集中, 一些数据点可能会被多次选中, 而另一些可能完全未被选中。当 $n' = n$ 时, 大约 63% 的数据点会被选中, 而剩下的 37% 被称为袋外样本点 (Out-of-Bag, OOB)。

- 为什么是 63%? 提示: 当 n 很大时, 某个样本点未被选中的概率是多少? 请只考虑在任意一个子样本集中 (而不是所有 T 个子样本集) 未被选中的概率。(7 points)
- 关于决策树的数量 T 。集成中的决策树数量通常需要在运行时间和降低方差之间进行权衡 (典型值范围从几十到几千棵树)。而子样本大小 n' 对运行时间的影响较小, 因此选择 n' 时主要考虑如何获得最优预测结果。虽然常见实践是设置 $n' = n$, 但这并非总是最佳选择。你会如何建议我们选择超参数 n' ? (7 points)

解:

(1) 对于原始数据集中的任意一个数据点, 在一次抽样中未被选中的概率为: $1 - \frac{1}{n}$

由于是有放回抽样, 进行 n 次抽样后, 该数据点始终未被选中的概率为: $(1 - \frac{1}{n})^n$

当 n 很大时, $(1 - \frac{1}{n})^n$ 近似于 e^{-1} , 即 0.3679

因此, 任意一个数据点未被选中的概率约为 36.79%, 被选中的概率约为 63.21%。

(2)

- 训练数据较少的情况下, 可以设置 $n' = n$, 以最大化每个分类器的训练数据利用率, 从而降低偏差。
- 训练数据较多的情况下, 可以适当减小 n' , 增加分类器之间的多样性, 从而降低方差。
- 使用交叉验证测试不同的 n' 值, 根据验证集的性能选择最优的子样本大小。

二. (20 points) 随机森林的实现

在本题中，你将实现决策树和随机森林，用于在以下两个数据集上进行分类：1) 垃圾邮件数据集，2) 泰坦尼克号数据集（用于预测这场著名灾难的幸存者）。数据已随作业提供。

为方便起见，我们提供了初始代码，其中包括预处理步骤和部分功能的实现。你可以自由选择使用或不使用这些代码来完成实现。

```
1 """
2 注意：
3 1. 这个框架提供了基本的结构，您需要完成所有标记为 'pass' 的函数。
4 2. 记得处理数值稳定性问题，例如在计算对数时避免除以零。
5 3. 在报告中详细讨论您的观察结果和任何有趣的发现。
6 """
7 from collections import Counter
8
9 import numpy as np
10 from numpy import genfromtxt
11 import scipy.io
12 from scipy import stats
13 from sklearn.tree import DecisionTreeClassifier, export_graphviz
14 from sklearn.base import BaseEstimator, ClassifierMixin
15 from sklearn.model_selection import cross_val_score
16 import pandas as pd
17 from pydot import graph_from_dot_data
18 import io
19
20 import random
21 random.seed(246810)
22 np.random.seed(246810)
23
24 eps = 1e-5 # a small number
25
26 class BaggedTrees(BaseEstimator, ClassifierMixin):
27
28     def __init__(self, params=None, n=200):
29         if params is None:
30             params = {}
31         self.params = params
32         self.n = n
33         self.decision_trees = [
34             DecisionTreeClassifier(random_state=i, self.params)
35             for i in range(self.n)
36         ]
37
38     def fit(self, X, y):
39         # TODO
40         pass
41
42     def predict(self, X):
43         # TODO
44         pass
```

```
45
46
47 class RandomForest(BaggedTrees):
48
49     def __init__(self, params=None, n=200, m=1):
50         if params is None:
51             params = {}
52         params['max_features'] = m
53         self.m = m
54         super().__init__(params=params, n=n)
55
56
57 class BoostedRandomForest(RandomForest):
58     # OPTIONAL
59     def fit(self, X, y):
60         # TODO
61         pass
62
63     def predict(self, X):
64         # TODO
65         pass
66
67
68 def preprocess(data, fill_mode=True, min_freq=10, onehot_cols=[]):
69     # Temporarily assign -1 to missing data
70     data[data == b''] = '-1'
71
72     # Hash the columns (used for handling strings)
73     onehot_encoding = []
74     onehot_features = []
75     for col in onehot_cols:
76         counter = Counter(data[:, col])
77         for term in counter.most_common():
78             if term[0] == b'-1':
79                 continue
80             if term[-1] <= min_freq:
81                 break
82             onehot_features.append(term[0])
83             onehot_encoding.append((data[:, col] == term[0]).astype(float))
84     data[:, col] = '0'
85     onehot_encoding = np.array(onehot_encoding).T
86     data = np.hstack(
87         [np.array(data, dtype=float),
88          np.array(onehot_encoding)])
89
90     # Replace missing data with the mode value. We use the mode instead of
91     # the mean or median because this makes more sense for categorical
92     # features such as gender or cabin type, which are not ordered.
93     if fill_mode:
94         # TODO
95         pass
96
```

```
97     return data, onehot_features
98
99
100 def evaluate(clf):
101     print("Cross validation", cross_val_score(clf, X, y))
102     if hasattr(clf, "decision_trees"):
103         counter = Counter([t.tree_.feature[0] for t in clf.decision_trees])
104         first_splits = [
105             (features[term[0]], term[1]) for term in counter.most_common()
106         ]
107         print("First splits", first_splits)
108
109
110 if __name__ == "__main__":
111     dataset = "titanic"
112     # dataset = "spam"
113     params = {
114         "max_depth": 5,
115         # "random_state": 6,
116         "min_samples_leaf": 10,
117     }
118     N = 100
119
120     if dataset == "titanic":
121         # Load titanic data
122         path_train = 'datasets/titanic/titanic_training.csv'
123         data = genfromtxt(path_train, delimiter=',', dtype=None)
124         path_test = 'datasets/titanic/titanic_testing_data.csv'
125         test_data = genfromtxt(path_test, delimiter=',', dtype=None)
126         y = data[1:, 0] # label = survived
127         class_names = ["Died", "Survived"]
128
129         labeled_idx = np.where(y != b'')[0]
130         y = np.array(y[labeled_idx], dtype=float).astype(int)
131         print("\n\nPart (b): preprocessing the titanic dataset")
132         X, onehot_features = preprocess(data[1:, 1:], onehot_cols=[1, 5, 7, 8])
133         X = X[labeled_idx, :]
134         Z, _ = preprocess(test_data[1:, :], onehot_cols=[1, 5, 7, 8])
135         assert X.shape[1] == Z.shape[1]
136         features = list(data[0, 1:]) + onehot_features
137
138     elif dataset == "spam":
139         features = [
140             "pain", "private", "bank", "money", "drug", "spam", "prescription",
141             "creative", "height", "featured", "differ", "width", "other",
142             "energy", "business", "message", "volumes", "revision", "path",
143             "meter", "memo", "planning", "pleased", "record", "out",
144             "semicolon", "dollar", "sharp", "exclamation", "parenthesis",
145             "square_bracket", "ampersand"
146         ]
147         assert len(features) == 32
148
```

```
149     # Load spam data
150     path_train = 'datasets/spam_data/spam_data.mat'
151     data = scipy.io.loadmat(path_train)
152     X = data['training_data']
153     y = np.squeeze(data['training_labels'])
154     Z = data['test_data']
155     class_names = ["Ham", "Spam"]
156
157     else:
158         raise NotImplementedError("Dataset %s not handled" % dataset)
159
160     print("Features", features)
161     print("Train/test size", X.shape, Z.shape)
162
163     # Decision Tree
164     print("\n\nDecision Tree")
165     dt = DecisionTreeClassifier(max_depth=3)
166     dt.fit(X, y)
167
168     # Visualize Decision Tree
169     print("\n\nTree Structure")
170     # Print using repr
171     print(dt.__repr__())
172     # Save tree to pdf
173     graph_from_dot_data(dt.to_graphviz())[0].write_pdf("%s-basic-tree.pdf" % dataset)
174
175     # Random Forest
176     print("\n\nRandom Forest")
177     rf = RandomForest(params, n=N, m=np.int_(np.sqrt(X.shape[1])))
178     rf.fit(X, y)
179     evaluate(rf)
180
181     # Generate Test Predictions
182     print("\n\nGenerate Test Predictions")
183     pred = rf.predict(Z)
```

Listing 1: 随机森林模型接口

1. 请参考以上模板实现随机森林算法。你也可以选择不参考模板，自行实现，但是不允许使用任何现成的随机森林实现。不过你可以使用库中提供的单棵决策树实现（我们在模板代码中使用了 `sklearn.tree.DecisionTreeClassifier`）。如果使用模板代码，你主要需要实现随机森林继承的超类，即 `bagged trees` 的实现，它会基于不同的数据样本创建并训练多棵决策树。完成后，请在模板中补充上缺失的部分。（5 points）

```
1 from collections import Counter
2 from sklearn.model_selection import train_test_split
3 import numpy as np
4 import scipy.io
5 from scipy import stats
6 from sklearn.tree import DecisionTreeClassifier, export_graphviz
7 from sklearn.metrics import accuracy_score
8 from sklearn.base import BaseEstimator, ClassifierMixin
9 from sklearn.model_selection import cross_val_score
```

```
10 import pandas as pd
11 from pydot import graph_from_dot_data
12 import io
13 import random
14
15 random.seed(246810)
16 np.random.seed(246810)
17
18 eps = 1e-5 # a small number
19
20 class BaggedTrees(BaseEstimator, ClassifierMixin):
21     def __init__(self, params=None, n=200):
22         if params is None:
23             params = {}
24         self.params = params
25         self.n = n
26         self.decision_trees = [
27             DecisionTreeClassifier(self.params, random_state=i)
28             for i in range(self.n)
29         ]
30
31     def fit(self, X, y):
32         self.trees_ = []
33         n_samples = X.shape[0]
34         for tree in self.decision_trees:
35             indices = np.random.choice(n_samples, n_samples, replace=True)
36             X_bootstrap = X[indices]
37             y_bootstrap = y[indices]
38             tree.fit(X_bootstrap, y_bootstrap)
39             self.trees_.append(tree)
40         return self
41
42     def predict(self, X):
43         predictions = np.array([tree.predict(X) for tree in self.trees_])
44         y_pred, _ = stats.mode(predictions, axis=0)
45         return y_pred.flatten()
46
47 class RandomForest(BaggedTrees):
48     def __init__(self, params=None, n=200, m=1):
49         if params is None:
50             params = {}
51         params['max_features'] = m
52         self.m = m
53         super().__init__(params=params, n=n)
54
55     def fit(self, X, y):
56         self.samples = []
57         self.feature_subsets = [] # 保存每棵树使用的特征索引
58         for tree in self.decision_trees:
59             indices = np.random.choice(len(X), len(X), replace=True)
60             X_sample = X[indices]
61             y_sample = y[indices]
```



```
62         features = np.random.choice(X.shape[1], self.m, replace=False)
63         X_sample = X_sample[:, features]
64         self.feature_subsets.append(features)
65         model = tree.fit(X_sample, y_sample)
66         self.samples.append((X_sample, y_sample))
67         return self
68
69     def predict(self, X):
70         predictions = np.array([tree.predict(X[:, features]) for tree, features in zip(self.
71         decision_trees, self.feature_subsets)])
72         y_pred, _ = stats.mode(predictions, axis=0)
73         return y_pred.flatten()
74
75 class BoostedRandomForest(RandomForest):
76     # OPTIONAL
77     def fit(self, X, y):
78         # TODO
79         pass
80
81     def predict(self, X):
82         # TODO
83         pass
84
85 def preprocess(data, fill_mode=True, min_freq=10, onehot_cols=[]):
86     # Temporarily assign -1 to missing data
87     data[data == b''] = '-1'
88
89     # Hash the columns (used for handling strings)
90     onehot_encoding = []
91     onehot_features = []
92     for col in onehot_cols:
93         counter = Counter(data[:, col])
94         for term in counter.most_common():
95             if term[0] == b'-1':
96                 continue
97             if term[-1] <= min_freq:
98                 break
99             onehot_features.append(f"col{col}_{term[0]}")
100             onehot_encoding.append((data[:, col] == term[0]).astype(float))
101     data[:, col] = '0'
102     onehot_encoding = np.array(onehot_encoding).T
103     data = np.hstack(
104         [np.array(data, dtype=float),
105          np.array(onehot_encoding)])
106
107     # Replace missing data with the mode value. We use the mode instead of
108     # the mean or median because this makes more sense for categorical
109     # features such as gender or cabin type, which are not ordered.
110     if fill_mode:
111         for col in range(data.shape[1]):
112             mode_value, _ = stats.mode(data[:, col])
113             data[:, col] = np.where(data[:, col] == '-1', mode_value, data[:, col])
```

```
113     return data, onehot_features
114
115 def evaluate(clf):
116     print("Cross validation", cross_val_score(clf, X, y))
117     if hasattr(clf, "decision_trees"):
118         counter = Counter([t.tree_.feature[0] for t in clf.decision_trees])
119         first_splits = [
120             (features[term[0]], term[1]) for term in counter.most_common()
121         ]
122         print("First splits", first_splits)
123
124 def print_decision_path(tree, X, feature_names, class_names):
125     node_indicator = tree.decision_path(X)
126     leave_id = tree.apply(X)
127     feature = tree.tree_.feature
128     threshold = tree.tree_.threshold
129
130     sample_id = 0
131     node_index = node_indicator.indices[node_indicator.indptr[sample_id]:node_indicator.indptr[
sample_id + 1]]
132
133     decision_path = []
134     for node_id in node_index:
135         if leave_id[sample_id] == node_id:
136             continue
137         if feature[node_id] == -2:
138             continue
139         if X[sample_id, feature[node_id]] <= threshold[node_id]:
140             threshold_sign = "<="
141         else:
142             threshold_sign = ">"
143         decision_path.append(f'{{feature_names[feature[node_id]]}} {{threshold_sign}} {{threshold[
node_id]}}')
144     for i, decision in enumerate(decision_path):
145         print(f"{{i+1}}. {{decision}}")
146     predicted_class = tree.predict(X)[0]
147     print(f"因此, 该样本被分类为 '{{class_names[int(predicted_class)]}}'。")
148
149 def get_random_forest_root_stats(rf, feature_names):
150     root_feature_indices = [tree.tree_.feature[0] for tree in rf.decision_trees]
151
152     feature_counts = Counter(root_feature_indices)
153
154     feature_counts_named = {}
155     for feature_index, count in feature_counts.items():
156         feature_name = feature_names[feature_index]
157         feature_counts_named[feature_name] = count
158
159     sorted_feature_counts = sorted(feature_counts_named.items(), key=lambda x: x[1], reverse=True
)
160
161     print("\n随机森林根节点常见分裂特征 (按使用次数排序):")
```

```
162     for feature_name, count in sorted_feature_counts:
163         print(f" - 特征 '{feature_name}' 被 {count} 棵树用作根节点分裂")
164
165 if __name__ == "__main__":
166     #dataset = "titanic"
167     dataset = "spam"
168     params = {
169         "max_depth": 5,
170         # "random_state": 6,
171         "min_samples_leaf": 10,
172     }
173     N = 100
174
175     if dataset == "titanic":
176         train_path = 'datasets/titanic/train.csv'
177         test_path = 'datasets/titanic/test.csv'
178         gender_path = 'datasets/titanic/gender_submission.csv'
179         class_names = ["Died", "Survived"]
180         train_titanic = pd.read_csv(train_path)
181         test_titanic = pd.read_csv(test_path)
182         gender = pd.read_csv(gender_path)
183
184         y_train = train_titanic['Survived'].to_numpy()
185         y_test = gender['Survived'].to_numpy()
186
187         train_data = train_titanic.drop(columns=['PassengerId', 'Survived', 'Name']).copy()
188         test_data = test_titanic.drop(columns=['PassengerId', 'Name']).copy()
189
190         train_data = train_data.fillna('-1').astype(str).to_numpy()
191         test_data = test_data.fillna('-1').astype(str).to_numpy()
192         onehot_cols = [1,5,7,8]
193         X_train, onehot_features = preprocess(train_data, onehot_cols=onehot_cols)
194         X_test, _ = preprocess(test_data, onehot_cols=onehot_cols)
195
196         assert X_train.shape[1] == X_test.shape[1]
197
198         features = list(train_titanic.columns[2:])
199         features.remove('Name')
200         features=features+ onehot_features
201
202
203     elif dataset == "spam":
204         features = [
205             "pain", "private", "bank", "money", "drug", "spam", "prescription",
206             "creative", "height", "featured", "differ", "width", "other",
207             "energy", "business", "message", "volumes", "revision", "path",
208             "meter", "memo", "planning", "pleased", "record", "out",
209             "semicolon", "dollar", "sharp", "exclamation", "parenthesis",
210             "square_bracket", "ampersand"
211         ]
212         assert len(features) == 32
213
```

```
214     # Load spam data
215     path_train = 'datasets/spam_data/spam_data.mat'
216     data = scipy.io.loadmat(path_train)
217     print(data.keys())
218     X = data['training_data']
219     y = np.squeeze(data['training_labels'])
220     Z = data['test_data']
221     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
222     class_names = ["Ham", "Spam"]
223
224     else:
225         raise NotImplementedError("Dataset %s not handled" % dataset)
226
227     print("Features", features)
228     print("Train/test size", X_train.shape, X_test.shape)
229
230     # Decision Tree
231     print("\n\nDecision Tree")
232     dt = DecisionTreeClassifier(max_depth=3)
233     dt.fit(X_train, y_train)
234
235     # Visualize Decision Tree
236     print("\n\nTree Structure")
237     # Print using repr
238     print(dt.__repr__())
239     # Save tree to pdf
240     dot_data = io.StringIO()
241     export_graphviz(dt, out_file=dot_data, feature_names=features, class_names=class_names)
242     graph = graph_from_dot_data(dot_data.getvalue())[0]
243     graph.write_pdf("%s-basic-tree.pdf" % dataset)
244
245     # Random Forest
246     print("\n\nRandom Forest")
247     m1=max(np.int_(np.ceil(np.sqrt(X_train.shape[1]))),np.int_(np.ceil(X_train.shape[1]/2)))
248     rf = RandomForest(params, n=N, m=m1)
249     rf.fit(X_train, y_train)
250     evaluate(rf)
251
252     # Generate Test Predictions
253     print("\n\nGenerate Test Predictions")
254     if dataset == "spam":
255         # 选择垃圾邮件和正常邮件的样本
256         spam_sample = X_train[y_train == 1][0:1] # 第一封垃圾邮件
257         ham_sample = X_train[y_train == 0][0:1] # 第一封正常邮件
258
259         print("\n垃圾邮件的决策树分裂路径:")
260         print_decision_path(dt, spam_sample, features, class_names)
261
262         print("\n正常邮件的决策树分裂路径:")
263         print_decision_path(dt, ham_sample, features, class_names)
264
265     elif dataset == "titanic":
```

```
266     # 选择幸存者和未幸存者的样本
267     survived_sample = X_train[y_train == 1][0:1]
268     not_survived_sample = X_train[y_train == 0][0:1]
269
270     print("\n幸存者的决策树分裂路径:")
271     print_decision_path(dt, survived_sample, features, class_names)
272
273     print("\n未幸存者的决策树分裂路径:")
274     print_decision_path(dt, not_survived_sample, features, class_names)
275
276     # 获取随机森林中根节点使用的特征的统计信息
277     get_random_forest_root_stats(rf, features)
278
279
280     # 计算决策树的训练和测试准确率
281     print("\n\n决策树准确率")
282     train_pred_dt = dt.predict(X_train) # 使用训练数据进行预测
283     test_pred_dt = dt.predict(X_test)   # 使用测试数据进行预测
284     train_acc_dt = accuracy_score(y_train, train_pred_dt)
285     test_acc_dt = accuracy_score(y_test, test_pred_dt)
286     print(f"决策树训练准确率: {train_acc_dt:.4f}")
287     print(f"决策树测试准确率: {test_acc_dt:.4f}")
288
289     # 计算随机森林的训练和测试准确率
290     print("\n\n随机森林准确率")
291     train_pred_rf = rf.predict(X_train) # 使用训练数据进行预测
292     test_pred_rf = rf.predict(X_test)   # 使用测试数据进行预测
293     train_acc_rf = accuracy_score(y_train, train_pred_rf)
294     test_acc_rf = accuracy_score(y_test, test_pred_rf)
295     print(f"随机森林训练准确率: {train_acc_rf:.4f}")
296     print(f"随机森林测试准确率: {test_acc_rf:.4f}")
```

Listing 2: 随机森林完整实现

2. 不需要长篇大论，每个问题用 1-2 句话回答即可：(5 points)

- 你是如何处理分类特征和缺失值的？
- 你是如何实现随机森林的？
- 你是否采用了特殊的方法加速训练？（回答“没有”也是可以的。）
- 还有什么特别棒的功能你实现了吗？（回答“没有”也是可以的。）

解:

1. 我们对指定的分类特征进行了独热编码（one-hot encoding），将类别转换为二进制特征。对于缺失值，首先用占位符'-1' 进行替换，然后在预处理过程中将其替换为每列的众数。
2. 我们通过训练多棵决策树来实现随机森林。对于每棵树，使用自助法（bootstrap）从原始数据中有放回地抽样生成训练集，并随机选择 m 个特征用于训练。预测时，通过所有树的多数投票来确定最终分类结果。

3. 没有

4. 我们实现了函数来获取单个数据点在决策树中的分裂路径，以及统计随机森林中最常见的根节点分裂特征，以帮助分析模型的决策过程。

3. 对于这两个数据集，请分别训练一棵决策树和一个随机森林，并报告它们的训练准确率和测试准确率。你需要报告 8 个数字（2 个数据集 × 2 个分类器 × 训练/测试）。(5 points)

解:

```
Optimal number of clusters: 4
Optimal number of clusters (revised): 4
Cluster Inertia (Sum of Squared Distances): 205.22514747675916
Customer Cluster Analysis with Genre Distribution:
```

		Age	Annual Income (k\$)	Spending Score (1-100)	Count
Cluster	Genre				
0	Female	32.545455	85.272727	80.590909	22
	Male	33.277778	87.111111	82.666667	18
1	Female	26.000000	39.529412	59.500000	34
	Male	24.608696	40.695652	61.478261	23
2	Female	51.405405	47.594595	40.567568	37
	Male	57.392857	47.857143	39.178571	28
3	Female	40.263158	87.105263	24.947368	19
	Male	38.473684	85.894737	14.210526	19

图 1: titanic

spam决策树准确率	
决策树训练准确率:	0.7955
决策树测试准确率:	0.7971
spam随机森林准确率	
随机森林训练准确率:	0.8204
随机森林测试准确率:	0.8106

图 2: spam

4. 决策树和随机森林的决策分析。(5 points)

- 对于决策树，选择来自每个类别（垃圾邮件和正常邮件）的一条数据点，列出决策树为对其分类所做的分裂（即，在哪个特征上以及该特征的哪个值上进行分裂）。以下是一个示例：

1. ('hot') ≥ 2

2. ('thanks') < 1

3. ('nigeria') ≥ 3

4. 因此这封邮件是垃圾邮件。

1. ('budget') ≥ 2

2. ('spreadsheet') ≥ 1

3. 因此这封邮件是正常邮件。

- 对于随机森林，找出并列出行根节点处最常见的分裂。例如：

1. ('viagra') ≥ 3 (20 trees)

2. ('thanks') < 4 (15 trees)

3. ('nigeria') ≥ 1 (5 trees)

解:

垃圾邮件的决策树分裂路径:

1. 'exclamation' > 0.5
2. 'ampersand' <= 0.5
3. 'meter' <= 0.5

因此, 该样本被分类为 'Spam'。

正常邮件的决策树分裂路径:

1. 'exclamation' <= 0.5
2. 'meter' <= 0.5
3. 'parenthesis' > 0.5

因此, 该样本被分类为 'Ham'。

图 3: 决策树为单个数据点分类时的分裂路径

随机森林根节点常见分裂特征 (按使用次数排序):

- 特征 'prescription' 被 11 棵树用作根节点分裂
- 特征 'spam' 被 9 棵树用作根节点分裂
- 特征 'creative' 被 8 棵树用作根节点分裂
- 特征 'energy' 被 7 棵树用作根节点分裂
- 特征 'private' 被 7 棵树用作根节点分裂
- 特征 'bank' 被 7 棵树用作根节点分裂
- 特征 'message' 被 7 棵树用作根节点分裂
- 特征 'featured' 被 6 棵树用作根节点分裂
- 特征 'differ' 被 6 棵树用作根节点分裂
- 特征 'pain' 被 6 棵树用作根节点分裂
- 特征 'drug' 被 5 棵树用作根节点分裂
- 特征 'money' 被 5 棵树用作根节点分裂
- 特征 'width' 被 5 棵树用作根节点分裂
- 特征 'business' 被 5 棵树用作根节点分裂
- 特征 'height' 被 3 棵树用作根节点分裂
- 特征 'other' 被 3 棵树用作根节点分裂

图 4: 随机森林根节点最常见分裂

三. (20 points) 聚类理论

聚类是一种无监督学习任务，其核心是根据数据样本之间的相似性（通常由距离度量定义）将数据分成若干个簇。常用聚类算法如 k -均值和 DBSCAN 都需要特定的距离度量和超参数设置。以下问题围绕距离度量、目标函数、以及超参数设置展开：

1. 在聚类算法中，距离度量是衡量样本间相似性的基础，选择合适的距离度量对聚类效果有显著影响（5 points）。

(a) 给定样本 $x = (x_1, x_2, \dots, x_d)$ 和 $y = (y_1, y_2, \dots, y_d)$ ，分别写出以下三种常见距离度量的数学公式：

- 欧几里得距离
- 曼哈顿距离
- 余弦相似度（将其转换为距离形式）

(b) 在以下场景中，分析哪种距离更适合使用，并简要说明原因：

- 场景 1：高维稀疏特征向量（如文本数据的 TF-IDF 表示）
- 场景 2：二维几何分布数据（如图像中的空间点分布）

解：

(1) 三种常见距离度量的数学公式：

欧几里得距离：

$$d_{\text{Euclidean}}(x, y) = \sqrt{\sum_{i=1}^d (x_i - y_i)^2}$$

曼哈顿距离：

$$d_{\text{Manhattan}}(x, y) = \sum_{i=1}^d |x_i - y_i|$$

余弦相似度转换为距离形式：

$$d_{\text{cosine}}(x, y) = 1 - \text{cosine_similarity}(x, y) = 1 - \frac{\sum_{i=1}^d x_i y_i}{\sqrt{\sum_{i=1}^d x_i^2} \sqrt{\sum_{i=1}^d y_i^2}}$$

(2)

在高维稀疏特征向量中，**余弦距离**更适合，因为高维空间中，欧几里得距离可能失去分辨力，样本之间的距离可能趋于一致，从而难以有效区分样本。TF-IDF 等特征往往是稀疏向量，向量的模长差异较大，而余弦距离主要关注向量之间的方向相似性，与特征的稀疏性无关，因此效果更优。

对于二维几何分布数据，**欧几里得距离**更适合，因为二维几何数据直观地表示点在空间中的位置，欧几里得距离能够直接反映两点之间的实际空间距离。曼哈顿距离虽然也可以用，但在规则网格或道路网络等特定场景中更常见，而非一般几何分布数据。余弦距离对数据方向敏感，但在低维几何分布数据中意义较小，因此不适用。

2. k -均值聚类的目标函数与迭代过程, k -均值聚类的目标是最小化以下目标函数 (10 points):

$$J = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$$

其中, C_i 表示第 i 个聚类簇, μ_i 为第 i 个簇的中心。

- 推导在分配样本点到最近的簇中心时, 为什么目标函数 J 会减少。
- 推导为什么更新簇中心为簇内样本点的均值时, 目标函数 J 会减少。
- k -均值的超参数 k 对结果有何影响?
 - 如果 k 设置过大或过小, 分别可能会导致什么问题?
 - 提出一种确定 k 的方法, 并解释其原理。

解:

(1) k -均值聚类的目标是最小化以下目标函数:

$$J = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$$

在 k -均值算法的分配步骤中, 对于每个样本点 x , 将其分配到最近的簇中心对应的簇 C_i , 即:

$$i = \arg \min_j \|x - \mu_j\|^2$$

因为目标函数 J 是所有簇内点到簇中心的平方距离总和, 重新分配样本点到最近簇中心后, 每个点的距离非增, 目标函数总和 J 也会减小。

(2) 在更新步骤中, 将每个簇的中心 μ_i 更新为该簇内所有样本点的均值, 即:

$$\mu_i = \frac{1}{|C_i|} \sum_{x \in C_i} x$$

目标函数的每一项 $\sum_{x \in C_i} \|x - \mu_i\|^2$ 表示簇 C_i 的样本点到簇中心的平方距离总和。

设 $f(\mu_i) = \sum_{x \in C_i} \|x - \mu_i\|^2$, 这是关于 μ_i 的凸函数。求导得:

$$\frac{\partial f(\mu_i)}{\partial \mu_i} = -2 \sum_{x \in C_i} (x - \mu_i)$$

令导数为 0, 得:

$$\mu_i = \frac{1}{|C_i|} \sum_{x \in C_i} x$$

即簇中心更新为簇内样本点的均值时, 使得 $f(\mu_i)$ 最小。

因此, 更新簇中心为簇内样本点均值后, 簇内样本点到簇中心的平方距离总和减小, 即目标函数 J 减少。

3. 密度聚类（如 DBSCAN）依赖以下两个超参数（5 points）：

- ϵ （邻域半径）：定义一个点的邻域范围。
- MinPts（核心点的最小邻域点数）：定义核心点的密度阈值。

(a) 核心点、边界点和噪声点的定义是什么？它们在聚类中的作用分别是什么？

(b) 如果 ϵ 和 MinPts 设置不当，可能会出现哪些问题？

- ϵ 过大或过小
- MinPts 过大或过小

(c) 为什么 DBSCAN 不需要预先指定聚类簇的数量 k ？这对实际应用有什么优势？

解：

(1)

- 核心点：若一个点的 ϵ 邻域内包含至少 MinPts 个点，则该点为核心点。核心点位于高密度区域的内部，是聚类的核心组成部分。
- 边界点：若一个点的 ϵ 邻域内的点数小于 MinPts，但该点位于某个核心点的 ϵ 邻域内，则该点为边界点。边界点位于簇的边缘，连接核心点和噪声点。
- 噪声点：既不是核心点，也不在任何核心点的 ϵ 邻域内的点被视为噪声点。噪声点不属于任何簇，通常被视为离群点。

在聚类过程中，核心点是簇的中心，边界点扩展簇的边界，而噪声点则被排除在簇之外。

(2)

- ϵ 过大或过小：
 - ϵ 过大：会导致更多的点被包含在同一个 ϵ 邻域内，可能将不同的簇合并为一个簇，降低聚类的分辨率。
 - ϵ 过小：会导致大多数点的 ϵ 邻域内点数不足 MinPts，使得许多核心点无法形成，导致大部分点被标记为噪声，无法有效识别簇。
- MinPts 过大或过小：
 - MinPts 过大：需要更多的点才能形成核心点，可能导致只有密度非常高的区域才能形成簇，其他区域的点被标记为噪声。
 - MinPts 过小：容易将密度较低的区域也识别为簇，可能导致噪声点被错误地包含在簇中，降低聚类的准确性。

(3)

DBSCAN 通过密度来定义簇，而不是预先指定簇的数量。它根据数据点的分布，自适应地发现密度足够高的区域，并将这些区域定义为簇。因此，无需预先指定簇的数量 k 。

在实际应用中的优势如下：

- 无需人工设定簇的数量，减少了人为干预，适用于对数据分布未知的情况。
- 能够识别任意形状的簇，而不仅限于球形簇，适用于复杂数据分布。
- 能够有效识别并排除噪声点，提高聚类的鲁棒性。

四. (30 points) 聚类实战

使用商场客户数据集 (Mall Customer Dataset) 完成客户分群分析。该数据集包含客户的年龄 (Age)、年收入 (Annual Income) 和消费积分 (Spending Score) 等特征。你需要通过实现和优化聚类算法来完成客户画像分析。数据随作业提供。

为方便起见，每小题我们提供了部分初始代码，其中包括预处理步骤和部分功能的实现。你可以自由选择使用或不使用这些代码来完成实现。

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.preprocessing import StandardScaler
4
5 # 加载数据
6 def load_mall_data():
7     df = pd.read_csv("Mall_Customers.csv")
8     X = df[['Age', 'Annual Income (k$)', 'Spending Score (1-100)']].values
9     # 数据标准化
10    scaler = StandardScaler()
11    X_scaled = scaler.fit_transform(X)
12    return X_scaled, X, df
```

Listing 3: 数据加载

1. 在不借助外部实现的情况下，手动实现 KMeans 方法 (4 points)，在数据集上进行聚类，可视化聚类结果 (3 points)，并解决下列问题 (3 points)：

- 如何使用肘部法则确定合适的 k 值，绘图说明
- 简单分析每个客户群的特征
- 计算和分析簇内平方和 (inertia)

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.preprocessing import StandardScaler
4 import pandas as pd
5 from mpl_toolkits.mplot3d import Axes3D
6
7 # 数据加载部分代码
8 def load_mall_data():
9     df = pd.read_csv(r"E:\大三上\机器学习\作业\作业4\Mall_Customers.csv")
10    X = df[['Age', 'Annual Income (k$)', 'Spending Score (1-100)']].values
11    # 数据标准化
12    scaler = StandardScaler()
13    X_scaled = scaler.fit_transform(X)
14    return X_scaled, X, df
15
16
17 # KMeans 算法实现
18 class KMeans:
19     def __init__(self, n_clusters=3, max_iters=100):
20         self.n_clusters = n_clusters
```

```
21     self.max_iters = max_iters
22     self.centroids = None
23     self.labels = None
24     self.inertia_ = None # 簇内平方和
25
26     def fit(self, X):
27         np.random.seed(42)
28         n_samples, n_features = X.shape
29         # 随机初始化聚类中心
30         random_idx = np.random.choice(n_samples, self.n_clusters, replace=False)
31         self.centroids = X[random_idx]
32
33         for _ in range(self.max_iters):
34             # 计算每个样本到聚类中心的距离
35             distances = np.linalg.norm(X[:, np.newaxis] - self.centroids, axis=2)
36             self.labels = np.argmin(distances, axis=1)
37
38             # 计算新的聚类中心
39             new_centroids = np.array([X[self.labels == k].mean(axis=0) for k in range(self.
n_clusters)])
40
41             # 检查是否收敛
42             if np.all(self.centroids == new_centroids):
43                 break
44             self.centroids = new_centroids
45
46             # 计算簇内平方和
47             self.inertia_ = np.sum([np.sum((X[self.labels == k] - self.centroids[k]) ** 2) for k in
range(self.n_clusters)])
48         return self
49
50     def predict(self, X):
51         distances = np.linalg.norm(X[:, np.newaxis] - self.centroids, axis=2)
52         return np.argmin(distances, axis=1)
53
54 # 绘制聚类结果的散点图
55 def plot_clusters(X, labels, centroids=None):
56     fig = plt.figure(figsize=(10, 8))
57     ax = fig.add_subplot(111, projection='3d')
58
59     # 绘制数据点
60     for k in range(np.max(labels) + 1):
61         cluster_points = X[labels == k]
62         ax.scatter(cluster_points[:, 0], cluster_points[:, 1], cluster_points[:, 2], label=f"
Cluster {k+1}")
63
64     # 绘制聚类中心
65     if centroids is not None:
66         ax.scatter(centroids[:, 0], centroids[:, 1], centroids[:, 2], color='black', marker='x',
s=150, label='Centroids')
67
68     ax.set_title("3D Cluster Visualization")
```

```
69 ax.set_xlabel("Age")
70 ax.set_ylabel("Annual Income (k$)")
71 ax.set_zlabel("Spending Score (1-100)")
72 ax.legend()
73 plt.show()
74
75 # 使用肘部法则确定最佳k值
76 def elbow_method(X, max_k=10):
77     inertias = []
78     for k in range(1, max_k + 1):
79         kmeans = KMeans(n_clusters=k)
80         kmeans.fit(X)
81         inertias.append(kmeans.inertia_)
82         print(f"k = {k}, Inertia = {kmeans.inertia_:.2f}")
83
84     # 可视化肘部法则
85     plt.figure(figsize=(8, 6))
86     plt.plot(range(1, max_k + 1), inertias, marker='o')
87     plt.title("Elbow Method")
88     plt.xlabel("Number of Clusters (k)")
89     plt.ylabel("Inertia")
90     plt.show()
91
92     # 找到最佳 k 值
93     optimal_k = find_elbow_point(inertias)
94     print(f"Optimal number of clusters: {optimal_k}")
95     return inertias, optimal_k
96
97 # 使用几何方法（线段垂距）寻找肘部点
98 def find_elbow_point(inertias):
99     x = np.arange(1, len(inertias) + 1)
100     y = np.array(inertias)
101
102     # 直线两端点
103     line_start = np.array([x[0], y[0]])
104     line_end = np.array([x[-1], y[-1]])
105
106     # 计算垂直距离
107     distances = np.abs((line_end[1] - line_start[1]) * x -
108                        (line_end[0] - line_start[0]) * y +
109                        line_end[0] * line_start[1] -
110                        line_end[1] * line_start[0]) / np.sqrt(
111         (line_end[1] - line_start[1])**2 + (line_end[0] - line_start[0])**2
112     )
113     optimal_k = np.argmax(distances) + 1
114     return optimal_k
115
116 if __name__ == "__main__":
117     X_scaled, X, df = load_mall_data()
118
119     # 使用肘部法则
120     inertias, optimal_k = elbow_method(X_scaled)
```

```
121
122 # 根据肘部法则选择的最佳 k 值
123 print(f"Optimal number of clusters (revised): {optimal_k}")
124
125 # 进行聚类
126 kmeans = KMeans(n_clusters=optimal_k)
127 kmeans.fit(X_scaled)
128
129 # 输出聚类结果与簇内平方和
130 print("Cluster Inertia (Sum of Squared Distances):", kmeans.inertia_)
131
132 # 可视化聚类结果
133 plot_clusters(X_scaled, kmeans.labels, kmeans.centroids)
134
135 # 分析聚类结果
136 df['Cluster'] = kmeans.labels
137 grouped = df.groupby(['Cluster', 'Genre']).agg({
138     'Age': 'mean',
139     'Annual Income (k$)': 'mean',
140     'Spending Score (1-100)': 'mean',
141     'Genre': 'count'
142 }).rename(columns={'Genre': 'Count'})
143 print("Customer Cluster Analysis with Genre Distribution:")
144 print(grouped)
145
146 # 可视化按簇划分的分布
147 genre_distribution = df.groupby(['Cluster', 'Genre']).size().unstack()
148 genre_distribution.plot(kind='bar', stacked=True, figsize=(10, 6))
149 plt.title("Genre Distribution Across Clusters")
150 plt.xlabel("Cluster")
151 plt.ylabel("Count")
152 plt.legend(title="Genre")
153 plt.show()
```

Listing 4: Kmeans 部分实现

解:

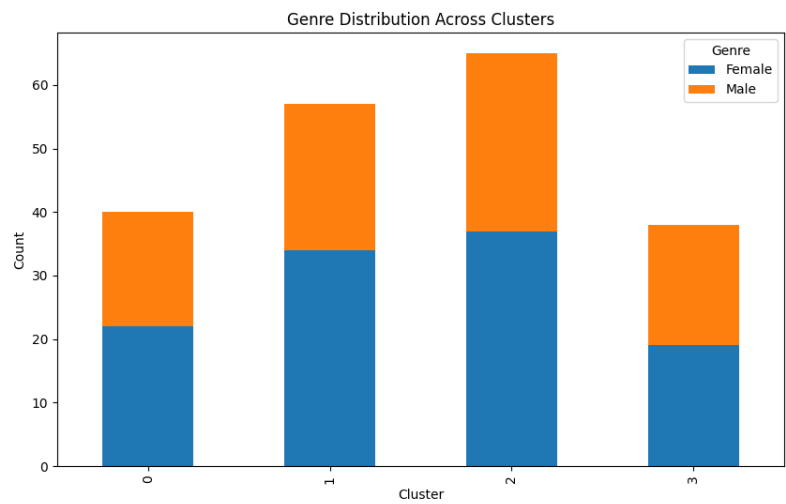


图 7: 性别聚类分析

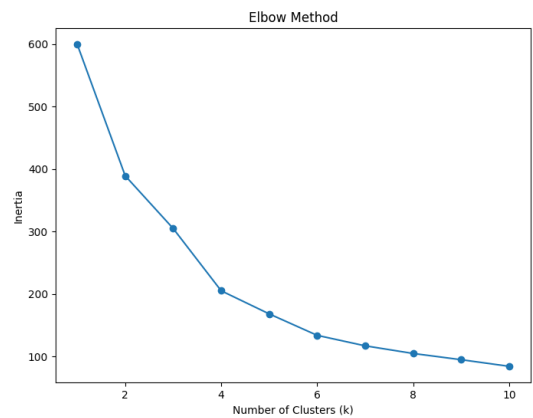


图 5: 肘部法则绘图

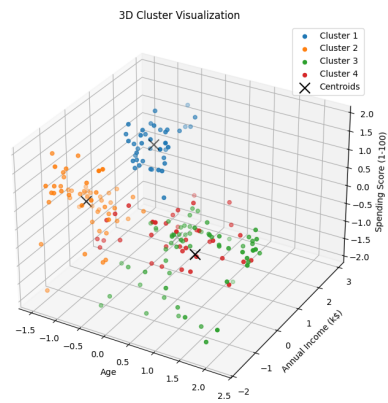


图 6: 聚类分析绘图


```
Optimal number of clusters: 4
Optimal number of clusters (revised): 4
Cluster Inertia (Sum of Squared Distances): 205.22514747675916
Customer Cluster Analysis with Genre Distribution:
```

		Age	Annual Income (k\$)	Spending Score (1-100)	Count
0	Female	32.545455	85.272727	80.590909	22
	Male	33.277778	87.111111	82.666667	18
1	Female	26.000000	39.529412	59.500000	34
	Male	24.608696	40.695652	61.478261	23
2	Female	51.405405	47.594595	40.567568	37
	Male	57.392857	47.857143	39.178571	28
3	Female	40.263158	87.105263	24.947368	19
	Male	38.473684	85.894737	14.210526	19

图 8: 输出结果

结合肘部法则绘图和代码输出结果来看，在 $k = 4$ 时达到了“肘点”。此时，簇内平方和为 205.23。

根据代码输出结果，每个客户群的特征如下：

- 簇 0:
 - 性别分布：女性略多于男性 (22:18)。
 - 平均年龄：客户群较年轻，平均年龄约为 32-33 岁。
 - 年收入：高收入群体，平均年收入约为 85-87 千美元。
 - 消费评分：消费评分很高，平均值为 80-83。
 - 特征总结：高收入、高消费的年轻客户群。
- 簇 1:
 - 性别分布：女性占优势 (34:23)。
 - 平均年龄：客户群年龄偏小，平均年龄约为 24-26 岁。
 - 年收入：中低收入群体，平均年收入约为 39-41 千美元。
 - 消费评分：消费评分中等偏高，平均值为 59-61。
 - 特征总结：年轻的中低收入客户，具备一定消费潜力，但可能受收入限制。
- 簇 2:
 - 性别分布：男性稍多于女性 (28:37)。
 - 平均年龄：年龄较大，女性约 51 岁，男性约 57 岁。
 - 年收入：中等收入群体，年收入约为 47-48 千美元。
 - 消费评分：消费评分偏低，平均值为 39-41。

- 特征总结：中等收入的中老年客户群，消费意愿较低。
- 簇 3：
 - 性别分布：男女比例相等 (19:19)。
 - 平均年龄：中年客户，女性约 40 岁，男性约 38 岁。
 - 年收入：高收入群体，年收入约为 85-87 千美元。
 - 消费评分：消费评分很低，女性约为 25，男性约为 14。
 - 特征总结：高收入但低消费的中年客户群。

2. 基于问题一的实现，我们发现随机初始化可能导致结果不稳定。请实现和分析以下改进：

- 实现 K-means++ 初始化方法 (4 分)
- 实现聚类评估指标 (3 分)：轮廓系数 (Silhouette Score)、聚类稳定性评估
- 对比分析 (3 分)
 1. 比较随机初始化和 K-means++ 的结果差异，可以通过可视化聚类图进行对比
 2. 比较两种方法的稳定性
 3. 分析初始化对收敛速度的影响

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.metrics import silhouette_score
5 import pandas as pd
6 from mpl_toolkits.mplot3d import Axes3D
7
8 # 数据加载部分代码
9 def load_mall_data():
10     df = pd.read_csv(r"E:\大三上\机器学习\作业\作业4\4\Mall_Customers.csv")
11     X = df[['Age', 'Annual Income (k$)', 'Spending Score (1-100)']].values
12     # 数据标准化
13     scaler = StandardScaler()
14     X_scaled = scaler.fit_transform(X)
15     return X_scaled, X, df
16
17 class KMeans:
18     def __init__(self, n_clusters=3, max_iters=100):
19         self.n_clusters = n_clusters
20         self.max_iters = max_iters
21         self.centroids = None
22         self.labels = None
23         self.inertia_ = None # 簇内平方和
24
25     def fit(self, X):
26         np.random.seed(None)
27         n_samples, n_features = X.shape
28         # 随机初始化聚类中心
29         random_idx = np.random.choice(n_samples, self.n_clusters, replace=False)
30         self.centroids = X[random_idx]
31
32         for _ in range(self.max_iters):
33             # 计算每个样本到聚类中心的距离
34             distances = np.linalg.norm(X[:, np.newaxis] - self.centroids, axis=2)
35             self.labels = np.argmin(distances, axis=1)
36
37             # 计算新的聚类中心
38             new_centroids = np.array([X[self.labels == k].mean(axis=0) for k in range(self.n_clusters)])
39
```

```
40         # 检查是否收敛
41         if np.all(self.centroids == new_centroids):
42             break
43         self.centroids = new_centroids
44
45         # 计算簇内平方和
46         self.inertia_ = np.sum([np.sum((X[self.labels == k] - self.centroids[k]) ** 2) for k in
47         range(self.n_clusters)])
48         return self
49
50     def predict(self, X):
51         distances = np.linalg.norm(X[:, np.newaxis] - self.centroids, axis=2)
52         return np.argmin(distances, axis=1)
53
54 # KMeans++ 算法实现
55 class KMeansPlusPlus:
56     def __init__(self, n_clusters=3, max_iters=100):
57         self.n_clusters = n_clusters
58         self.max_iters = max_iters
59         self.centroids = None
60         self.labels = None
61         self.inertia_ = None # 簇内平方和
62
63     def _kmeans_plus_plus_init(self, X):
64         n_samples, _ = X.shape
65         centroids = [X[np.random.randint(0, n_samples)]]
66
67         for _ in range(1, self.n_clusters):
68             distances = np.min([np.linalg.norm(X - c, axis=1) for c in centroids], axis=0)
69             probabilities = distances**2 / np.sum(distances**2)
70             cumulative_probs = np.cumsum(probabilities)
71             r = np.random.rand()
72             new_centroid_idx = np.searchsorted(cumulative_probs, r)
73             centroids.append(X[new_centroid_idx])
74
75         return np.array(centroids)
76
77     def fit(self, X):
78         n_samples, n_features = X.shape
79         self.centroids = self._kmeans_plus_plus_init(X)
80
81         for _ in range(self.max_iters):
82             distances = np.linalg.norm(X[:, np.newaxis] - self.centroids, axis=2)
83             self.labels = np.argmin(distances, axis=1)
84             new_centroids = np.array([X[self.labels == k].mean(axis=0) for k in range(self.
85             n_clusters)])
86
87             if np.all(self.centroids == new_centroids):
88                 break
89
90             self.centroids = new_centroids
91
92         self.inertia_ = np.sum([np.sum((X[self.labels == k] - self.centroids[k])**2) for k in
```

```
        range(self.n_clusters)])
90
91         return self
92
93 # 计算轮廓系数
94 def compute_silhouette_score(X, labels):
95     return silhouette_score(X, labels)
96
97 # 评估聚类结果的稳定性
98 def compute_cluster_stability(X, k, n_runs=10):
99     labels_list = []
100     for _ in range(n_runs):
101         kmeans = KMeansPlusPlus(n_clusters=k)
102         kmeans.fit(X)
103         labels_list.append(kmeans.labels)
104
105     # 计算一致性
106     stability = np.mean([np.sum(labels_list[i] == labels_list[j]) / len(X)
107                          for i in range(n_runs) for j in range(i + 1, n_runs)])
108
109     return stability
110
111 # 绘制聚类结果的散点图
112 def plot_clusters(X, labels, centroids=None):
113     fig = plt.figure(figsize=(10, 8))
114     ax = fig.add_subplot(111, projection='3d')
115
116     # 绘制数据点
117     for k in range(np.max(labels) + 1):
118         cluster_points = X[labels == k]
119         ax.scatter(cluster_points[:, 0], cluster_points[:, 1], cluster_points[:, 2], label=f"
120                    Cluster {k+1}")
121
122     # 绘制聚类中心
123     if centroids is not None:
124         ax.scatter(centroids[:, 0], centroids[:, 1], centroids[:, 2], color='black', marker='x',
125                   s=150, label='Centroids')
126
127     ax.set_title("3D Cluster Visualization")
128     ax.set_xlabel("Age")
129     ax.set_ylabel("Annual Income (k$)")
130     ax.set_zlabel("Spending Score (1-100)")
131     ax.legend()
132     plt.show()
133
134 # 使用肘部法则确定最佳k值
135 def elbow_method(X, max_k=10):
136     inertias = []
137     for k in range(1, max_k + 1):
138         kmeans = KMeansPlusPlus(n_clusters=k)
139         kmeans.fit(X)
140         inertias.append(kmeans.inertia_)
```

```
139     print(f"k = {k}, Inertia = {kmeans.inertia_:.2f}")
140
141     # 可视化肘部法则
142     plt.figure(figsize=(8, 6))
143     plt.plot(range(1, max_k + 1), inertias, marker='o')
144     plt.title("Elbow Method")
145     plt.xlabel("Number of Clusters (k)")
146     plt.ylabel("Inertia")
147     plt.show()
148
149     # 找到最佳 k 值
150     optimal_k = find_elbow_point(inertias)
151     print(f"Optimal number of clusters: {optimal_k}")
152     return inertias, optimal_k
153
154 # 使用几何方法（线段垂距）寻找肘部点
155 def find_elbow_point(inertias):
156     x = np.arange(1, len(inertias) + 1)
157     y = np.array(inertias)
158
159     # 直线两端点
160     line_start = np.array([x[0], y[0]])
161     line_end = np.array([x[-1], y[-1]])
162
163     # 计算垂直距离
164     distances = np.abs((line_end[1] - line_start[1]) * x -
165                        (line_end[0] - line_start[0]) * y +
166                        line_end[0] * line_start[1] -
167                        line_end[1] * line_start[0]) / np.sqrt(
168         (line_end[1] - line_start[1])**2 + (line_end[0] - line_start[0])**2
169     )
170     optimal_k = np.argmax(distances) + 1
171     return optimal_k
172
173
174 if __name__ == "__main__":
175     X_scaled, X, df = load_mall_data()
176     # 使用肘部法则
177     inertias, optimal_k = elbow_method(X_scaled)
178
179     # 根据肘部法则选择的最佳 k 值
180     print(f"Optimal number of clusters (revised): {optimal_k}")
181
182     kmeans_pp = KMeansPlusPlus(n_clusters=optimal_k)
183     kmeans_pp.fit(X_scaled)
184
185     silhouette_pp = compute_silhouette_score(X_scaled, kmeans_pp.labels)
186     stability_pp = compute_cluster_stability(X_scaled, optimal_k)
187
188     print(f"K-means++ Silhouette Score: {silhouette_pp:.4f}")
189     print(f"K-means++ Stability Score: {stability_pp:.4f}")
190
```

```
191 # 可视化K-means++ 结果
192 df['Cluster'] = kmeans_pp.labels
193 plot_clusters(X_scaled, kmeans_pp.labels, kmeans_pp.centroids)
194
```

Listing 5: Kmeans++ 部分实现

解:

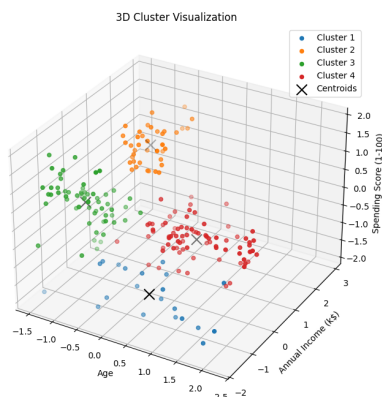


图 9: K-means++ 聚类分析

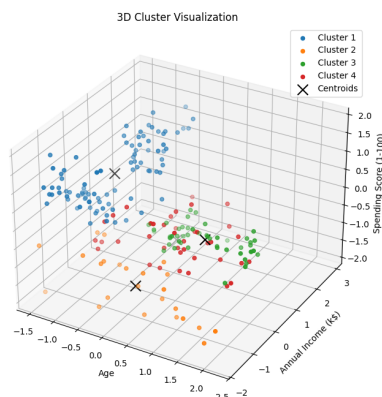


图 10: K-means 聚类分析

```
k = 1, Inertia = 600.00
k = 2, Inertia = 389.39
k = 3, Inertia = 297.03
k = 4, Inertia = 205.23
k = 5, Inertia = 168.25
k = 6, Inertia = 134.35
k = 7, Inertia = 124.12
k = 8, Inertia = 109.37
k = 9, Inertia = 99.30
k = 10, Inertia = 99.72
Optimal number of clusters: 4
Optimal number of clusters (revised): 4
K-means++ Silhouette Score: 0.4040
K-means++ Stability Score: 0.2807
```

图 11: K-means++ 输出结果

```
k = 1, Inertia = 600.00
k = 2, Inertia = 389.39
k = 3, Inertia = 297.03
k = 4, Inertia = 205.23
k = 5, Inertia = 168.25
k = 6, Inertia = 133.87
k = 7, Inertia = 132.94
k = 8, Inertia = 106.64
k = 9, Inertia = 102.08
k = 10, Inertia = 98.47
Optimal number of clusters: 4
Optimal number of clusters (revised): 4
K-means Silhouette Score: 0.3515
K-means Stability Score: 0.2244
```

图 12: K-means 输出结果

K-means++ 的轮廓系数为 0.4048，高于 K-means 的 0.3515；K-means++ 的稳定性得分为 0.2807，优于 K-means 的 0.2244。

因为 K-means++ 使用数据分布感知的初始化方式，使得初始聚类中心更接近实际聚类中心，从而提高了同簇内样本的紧密度，并增加了簇与簇之间的分离度，与此同时，K-means++ 减少了运行之间的随机性，每次运行的结果更加接近。

在聚类图中，采用 K-means++ 初始化方法的簇中心较为均匀，样本点分布比 K-means 的更加紧凑且分离度良好。

3. 在实际业务中，不同特征的重要性可能不同，且某些客户群可能需要大小相近。请实现带权重和大小约束的改进版本：

- 实现带约束的聚类算法，需要支持特征加权和簇大小约束 (4 points)
- 在以下两个场景下重新进行实验：收入特征权重加倍, 限制每个客户群至少包含 20% 的客户, 绘制聚类结果 (6 points)

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.metrics import silhouette_score
5 import pandas as pd
6
7 # 数据加载部分代码
8 def load_mall_data():
9     df = pd.read_csv(r"E:\大三上\机器学习\作业\作业4\4\Mall_Customers.csv")
10    X = df[['Age', 'Annual Income (k$)', 'Spending Score (1-100)']].values
11    # 数据标准化
12    scaler = StandardScaler()
13    X_scaled = scaler.fit_transform(X)
14    return X_scaled, X, df
15
16 # KMeans++ 算法实现
17 class KMeansPlusPlus:
18     def __init__(self, n_clusters=3, max_iters=100):
19         self.n_clusters = n_clusters
20         self.max_iters = max_iters
21         self.centroids = None
22         self.labels = None
23         self.inertia_ = None # 簇内平方和
24
25     def _kmeans_plus_plus_init(self, X):
26         n_samples, _ = X.shape
27         centroids = [X[np.random.randint(0, n_samples)]]
28
29         for _ in range(1, self.n_clusters):
30             distances = np.min([np.linalg.norm(X - c, axis=1) for c in centroids], axis=0)
31             probabilities = distances**2 / np.sum(distances**2)
32             cumulative_probs = np.cumsum(probabilities)
33             r = np.random.rand()
34             new_centroid_idx = np.searchsorted(cumulative_probs, r)
35             centroids.append(X[new_centroid_idx])
36
37     return np.array(centroids)
38
39 def fit(self, X):
```



```

40     n_samples, n_features = X.shape
41     self.centroids = self._kmeans_plus_plus_init(X)
42
43     for _ in range(self.max_iters):
44         distances = np.linalg.norm(X[:, np.newaxis] - self.centroids, axis=2)
45         self.labels = np.argmin(distances, axis=1)
46         new_centroids = np.array([X[self.labels == k].mean(axis=0) for k in range(self.
n_clusters)])
47         if np.all(self.centroids == new_centroids):
48             break
49
50         self.centroids = new_centroids
51
52         self.inertia_ = np.sum([np.sum((X[self.labels == k] - self.centroids[k])**2) for k in
range(self.n_clusters)])
53
54     return self
55
56 class ConstrainedKMeans(KMeansPlusPlus):
57     def __init__(self, n_clusters=3, max_iters=100, weights=None, size_constraints=None):
58         super().__init__(n_clusters, max_iters)
59         self.weights = weights
60         self.size_constraints = size_constraints
61
62     def _weighted_distance(self, X, centroids):
63         weighted_X = X * self.weights
64         weighted_centroids = centroids * self.weights
65         return np.linalg.norm(weighted_X[:, np.newaxis] - weighted_centroids, axis=2)
66
67     def _reassign_clusters(self, X, labels, distances):
68         cluster_sizes = np.bincount(labels, minlength=self.n_clusters)
69         for cluster in range(self.n_clusters):
70             while cluster_sizes[cluster] < self.size_constraints[0]: # 小于最小限制
71                 # 从其他簇中挑选样本重新分配
72                 for donor_cluster in range(self.n_clusters):
73                     if cluster_sizes[donor_cluster] > self.size_constraints[0]:
74                         # 找到距离当前簇最近的样本并重新分配
75                         donor_idx = np.where(labels == donor_cluster)[0]
76                         distances_to_current = distances[donor_idx, cluster]
77                         closest_idx = donor_idx[np.argmin(distances_to_current)]
78                         labels[closest_idx] = cluster
79                         cluster_sizes = np.bincount(labels, minlength=self.n_clusters)
80                         break
81                 else:
82                     raise ValueError("无法满足最小簇大小限制，可能是由于数据分布问题")
83         return labels
84
85     def fit(self, X):
86         n_samples, n_features = X.shape
87         if self.weights is None:
88             self.weights = np.ones(n_features)
89         self.centroids = self._kmeans_plus_plus_init(X)

```

```
90
91     for _ in range(self.max_iters):
92         distances = self._weighted_distance(X, self.centroids)
93         labels = np.argmin(distances, axis=1)
94         if self.size_constraints is not None:
95             labels = self._reassign_clusters(X, labels, distances)
96         new_centroids = np.array([X[labels == k].mean(axis=0) for k in range(self.n_clusters)
97 ])
98         if np.all(self.centroids == new_centroids):
99             break
100
101         self.centroids = new_centroids
102
103         self.labels = labels
104         self.inertia_ = np.sum([np.sum((X[labels == k] - self.centroids[k])**2) for k in range(
105 self.n_clusters)])
106         return self
107
108 # 绘制聚类结果的散点图
109 def plot_clusters(X, labels, centroids=None):
110     fig = plt.figure(figsize=(10, 8))
111     ax = fig.add_subplot(111, projection='3d')
112
113     # 绘制数据点
114     for k in range(np.max(labels) + 1):
115         cluster_points = X[labels == k]
116         ax.scatter(cluster_points[:, 0], cluster_points[:, 1], cluster_points[:, 2], label=f"
117 Cluster {k+1}")
118
119     # 绘制聚类中心
120     if centroids is not None:
121         ax.scatter(centroids[:, 0], centroids[:, 1], centroids[:, 2], color='black', marker='x',
122 s=150, label='Centroids')
123
124     ax.set_title("3D Cluster Visualization")
125     ax.set_xlabel("Age")
126     ax.set_ylabel("Annual Income (k$)")
127     ax.set_zlabel("Spending Score (1-100)")
128     ax.legend()
129     plt.show()
130
131 # 使用肘部法则确定最佳k值
132 def elbow_method(X, max_k=10):
133     inertias = []
134     for k in range(1, max_k + 1):
135         kmeans = KMeansPlusPlus(n_clusters=k)
136         kmeans.fit(X)
137         inertias.append(kmeans.inertia_)
138         print(f"k = {k}, Inertia = {kmeans.inertia_:.2f}")
139
140     # 找到最佳 k 值
141     optimal_k = find_elbow_point(inertias)
```

```
138     print(f"Optimal number of clusters: {optimal_k}")
139     return inertias, optimal_k
140
141 # 使用几何方法（线段垂距）寻找肘部点
142 def find_elbow_point(inertias):
143     x = np.arange(1, len(inertias) + 1)
144     y = np.array(inertias)
145
146     # 直线两端点
147     line_start = np.array([x[0], y[0]])
148     line_end = np.array([x[-1], y[-1]])
149
150     # 计算垂直距离
151     distances = np.abs((line_end[1] - line_start[1]) * x -
152                        (line_end[0] - line_start[0]) * y +
153                        line_end[0] * line_start[1] -
154                        line_end[1] * line_start[0]) / np.sqrt(
155         (line_end[1] - line_start[1])**2 + (line_end[0] - line_start[0])**2
156     )
157     optimal_k = np.argmax(distances) + 1
158     return optimal_k
159
160
161 if __name__ == "__main__":
162     X_scaled, X, df = load_mall_data()
163
164     # 使用肘部法则
165     inertias, optimal_k = elbow_method(X_scaled)
166
167     # 根据肘部法则选择的最佳 k 值
168     print(f"Optimal number of clusters (revised): {optimal_k}")
169
170     # 场景1: 收入特征权重加倍
171     weights = np.array([1, 2, 1]) # Age:1, Income:2, Spending:1
172     constrained_kmeans_1 = ConstrainedKMeans(n_clusters=optimal_k, weights=weights)
173     constrained_kmeans_1.fit(X_scaled)
174
175     print("Scenario 1: Income Feature Weighted Double")
176     plot_clusters(X_scaled, constrained_kmeans_1.labels, constrained_kmeans_1.centroids)
177
178     # 分析聚类结果
179     df['Cluster'] = constrained_kmeans_1.labels
180     grouped = df.groupby(['Cluster', 'Genre']).agg({
181         'Age': 'mean',
182         'Annual Income (k$)': 'mean',
183         'Spending Score (1-100)': 'mean',
184         'Genre': 'count'
185     }).rename(columns={'Genre': 'Count'})
186     print("Customer Cluster Analysis with Genre Distribution:")
187     print(grouped)
188
189     # 场景2: 每个客户群至少包含20%的客户
```

```
190 min_size = int(len(X) * 0.2)
191 max_size = len(X) # 无上限
192 constrained_kmeans_2 = ConstrainedKMeans(n_clusters=optimal_k, size_constraints=(min_size,
193 max_size))
194 constrained_kmeans_2.fit(X_scaled)
195
196 print("Scenario 2: Minimum 20% of Customers per Cluster")
197 plot_clusters(X_scaled, constrained_kmeans_2.labels, constrained_kmeans_2.centroids)
198 # 分析聚类结果
199 df['Cluster'] = constrained_kmeans_2.labels
200 grouped = df.groupby(['Cluster', 'Genre']).agg({
201     'Age': 'mean',
202     'Annual Income (k$)': 'mean',
203     'Spending Score (1-100)': 'mean',
204     'Genre': 'count'
205 }).rename(columns={'Genre': 'Count'})
206 print("Customer Cluster Analysis with Genre Distribution:")
207 print(grouped)
```

Listing 6: 带约束的聚类算法部分实现

解:

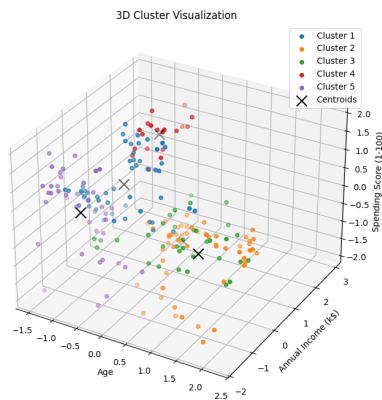


图 13: 场景一聚类分析

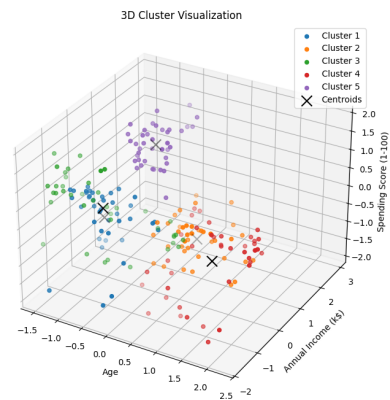


图 14: 场景二聚类分析

```

Optimal number of clusters: 5
Optimal number of clusters (revised): 5
Scenario 1: Income Feature Weighted Double
Customer Cluster Analysis with Genre Distribution:
      Age  Annual Income (k$)  Spending Score (1-100)  Count
Cluster Genre
0      Female  29.200000      72.828571      64.857143      35
      Male    30.956522      74.000000      71.782609      23
1      Female  26.950000      31.500000      69.550000      20
      Male    23.785714      30.785714      69.357143      14
2      Female  41.538462      26.538462      20.692308      13
      Male    52.000000      25.875000      17.250000       8
3      Female  43.000000      95.611111      31.277778      18
      Male    37.047619      91.666667      21.904762      21
4      Female  53.538462      53.500000      49.153846      26
      Male    57.454545      55.045455      47.681818      22

Scenario 2: Minimum 20% of Customers per Cluster
Customer Cluster Analysis with Genre Distribution:
      Age  Annual Income (k$)  Spending Score (1-100)  Count
Cluster Genre
0      Female  24.666667      37.571429      70.666667      21
      Male    22.368421      41.473684      63.315789      19
1      Female  56.789474      56.105263      49.736842      19
      Male    61.047619      51.142857      41.904762      21
2      Female  38.000000      84.523810      26.904762      21
      Male    38.473684      85.894737      14.210526      19
3      Female  39.862069      38.965517      34.620690      29
      Male    42.363636      37.636364      38.909091      11
4      Female  32.545455      85.272727      80.590909      22
      Male    33.277778      87.111111      82.666667      18

```

图 15: 输出结果

场景 1: 收入特征权重加倍

簇 0: 年轻, 高收入, 高消费

簇 1: 年轻, 低收入, 高消费

簇 2: 年长, 低收入, 低消费

簇 3: 中年, 高收入, 低消费

簇 4: 年长, 中等收入, 中等消费

收入特征的主导性明显, 各簇的年收入区分度较高。

场景 2: 每个客户群至少包含 20% 的客户

簇 0: 年轻, 低收入, 高消费

簇 1: 年长, 中等收入, 中等消费

簇 2: 中年, 高收入, 低消费

簇 3: 中年, 中等收入, 低消费

簇 4: 年轻, 高收入, 高消费

每个簇都包含 40 个客户。