

1. 数据集

采用全部数据进行聚类

```
# 导入数据集
dataset = pd.read_csv('iris.csv')
```

2. 手写K-means算法

kmeans算法手写源码

```
class KMeans:
    # 初始化
    def __init__(self, k=3, max_iters=300):
        self.k = k
        self.max_iters = max_iters
        # 质心
        self.centroids = None
        # 标签
        self.labels = None
        # 分类
        self.clusters = None
        # 损失值
        self.loss_value = []
        # SSE
        self.SSE = 0

    def _init_centroids(self, X):
        # 随机初始化质心
        return X[np.random.choice(X.shape[0], self.k, replace=False)]

    def _get_labels(self, X):
        # 计算距离
        distances = np.sqrt(np.sum((X - self.centroids[:, np.newaxis])**2,
axis=2))
        # 返回最近的质心索引
        return np.argmin(distances, axis=0)

    def _get_centroids(self, X):
        # 更新质心
        centroids = np.zeros((self.k, X.shape[1]))
        for i in range(self.k):
            centroids[i] = np.mean(X[self.labels == i], axis=0)
        return centroids

    def _get_clusters(self, X):
        # 更新SSE
        clusters = [[] for _ in range(self.k)]
        for i in range(self.k):
            clusters[i] = X[self.labels == i]
```

```

        return clusters

    def plot(self):
        # 可视化
        ax = plt.axes(projection='3d')
        for i in range(self.k):
            ax.scatter3D(self.clusters[i][:, 0], self.clusters[i][:, 1],
self.clusters[i][:, 2], label='cluster {}'.format(i))
            ax.scatter3D(self.centroids[:, 0], self.centroids[:, 1],
self.centroids[:, 2], marker='*', c='black', label='centroids')
            ax.set_xlabel('sepal length')
            ax.set_ylabel('sepal width')
            ax.set_zlabel('petal length')

# 训练
    def fit(self, X):
        # 初始化质心
        self.centroids = self._init_centroids(X)
        # 初始化标签
        self.labels = np.zeros((X.shape[0], 1))
        # 初始化分类
        self.clusters = [[] for _ in range(self.k)]
        # 迭代
        for _ in range(self.max_iters):
            # 更新标签
            self.labels = self._get_labels(X)
            # 更新质心
            new_centroids = self._get_centroids(X)
            if np.all(self.centroids == new_centroids):
                break
            self.centroids = self._get_centroids(X)
            # 更新SSE
            self.clusters = self._get_clusters(X)
            # 将损失值加入列表
            self.loss_value.append([_, self.loss(X=X)])

# 计算SSE
        for i in range(self.k):
            self.SSE += np.sum((self.clusters[i] - self.centroids[i])**2)
        # 设置字体为黑体
        plt.rcParams['font.sans-serif'] = ['SimHei']
        # 设置标题为k的值
        self.plot()
        plt.title('k = {}时的聚类结果'.format(self.k))
        # 展示损失值的折线图
        plt.figure()
        plt.plot([i[0] for i in self.loss_value], [i[1] for i in
self.loss_value])
        plt.title('k = {}时的损失值折线图'.format(self.k))

# 预测
    def predict(self, X):
        return self._get_labels(X)

# 损失值

```

```
def loss(self, X):
    loss = 0
    # 代价函数为畸变函数
    for i in range(X.shape[0]):
        # 找到最近的质心
        n_centroid = self.centroids[self.labels[i]]
        # 计算距离
        loss += np.sum((X[i] - n_centroid)**2)
    loss = loss / X.shape[0]
    return loss
```

华为云服务器运行结果

使用aechoterm ssh工具连接到华为云服务器

```
root@ml-routhleek:~/ml_exp3# python3 kmeans.py
k = 1, SSE = 594.8006666666668
k = 2, SSE = 133.46431822602608
k = 3, SSE = 69.42973924466338
k = 4, SSE = 49.43781558441559
k = 5, SSE = 40.325266666666664
root@ml-routhleek:~/ml_exp3#
```

输出的图像文件:

远程 隐藏文件 刷新

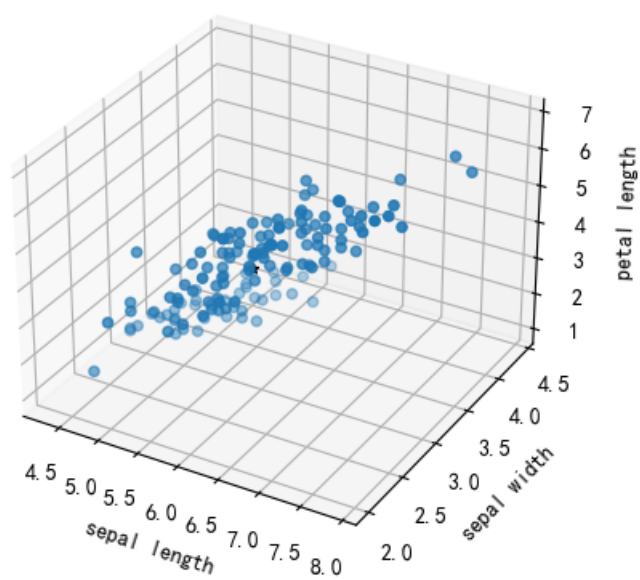
/ root ml_exp3

名称	大小	类型	修改时间
..			
iris.csv	4.86K	文件	2022-12-18 10:44:16
k=1_loss.png	14.51K	文件	2022-12-18 11:00:07
k=1_result.png	78.32K	文件	2022-12-18 11:00:06
k=2_loss.png	18.32K	文件	2022-12-18 11:00:07
k=2_result.png	89.21K	文件	2022-12-18 11:00:07
k=3_loss.png	25.21K	文件	2022-12-18 11:00:07
k=3_result.png	88.34K	文件	2022-12-18 11:00:07
k=4_loss.png	20.51K	文件	2022-12-18 11:00:08
k=4_result.png	94.65K	文件	2022-12-18 11:00:07
k=5_loss.png	21.67K	文件	2022-12-18 11:00:08
k=5_result.png	92.48K	文件	2022-12-18 11:00:08
kmeans.py	4.90K	文件	2022-12-18 10:48:38
SSE.png	18.26K	文件	2022-12-18 11:00:08

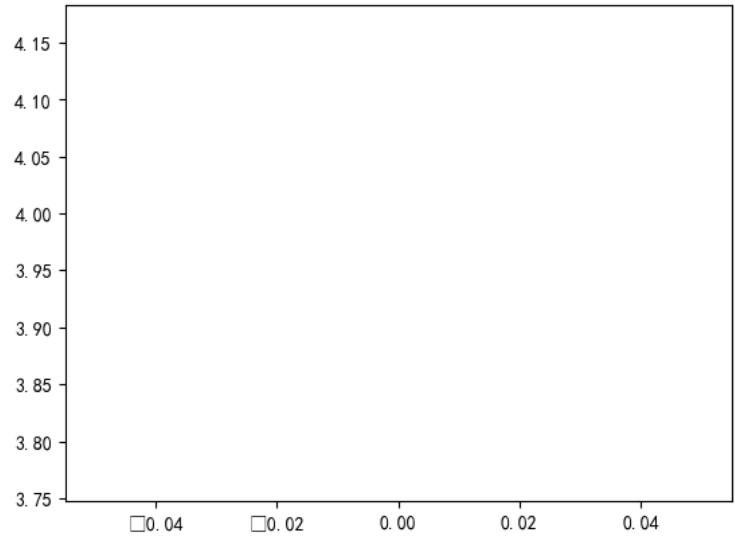
3. 实验结果可视化

k=1

k = 1时的聚类结果

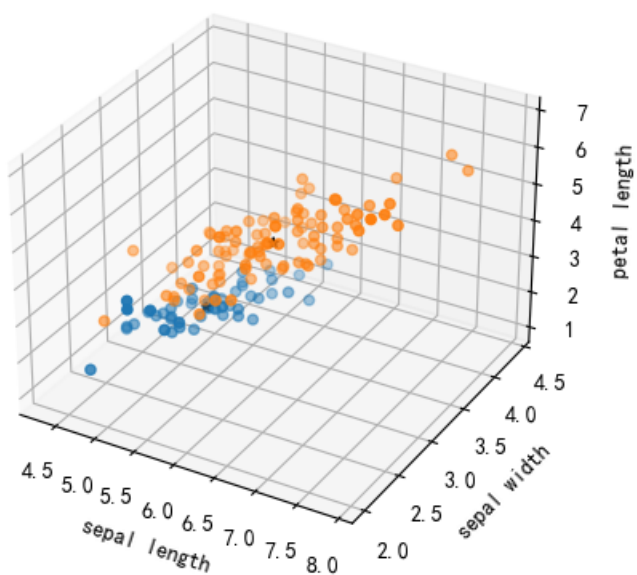


k = 1时的损失值折线图

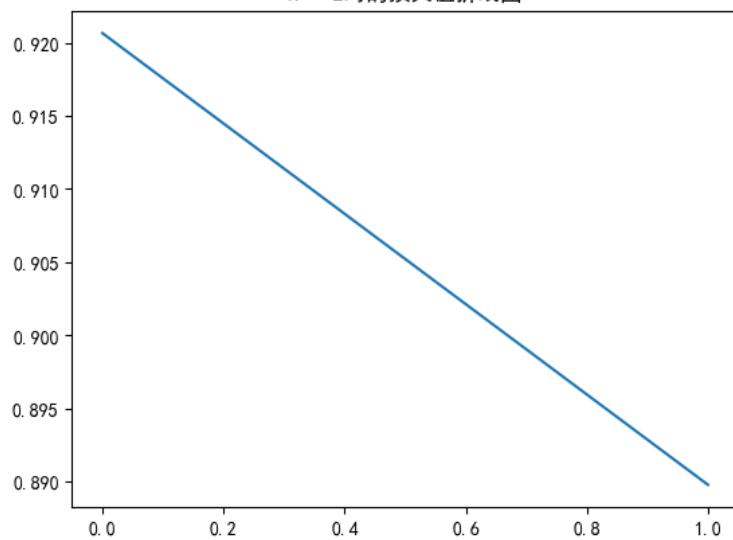


k=2

k = 2时的聚类结果

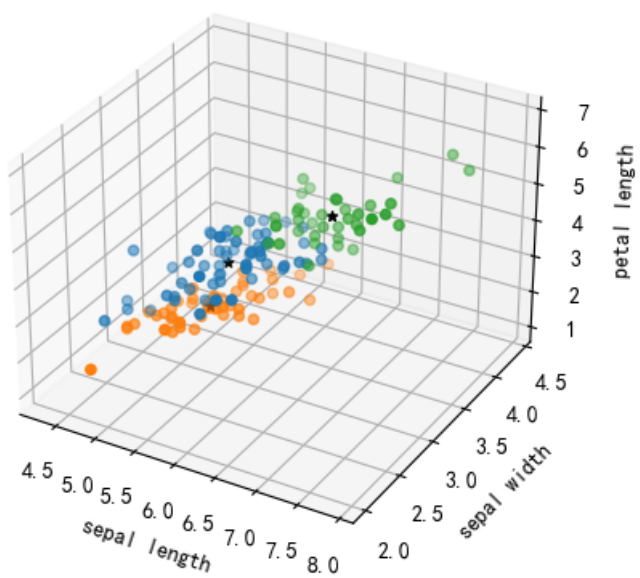


k = 2时的损失值折线图

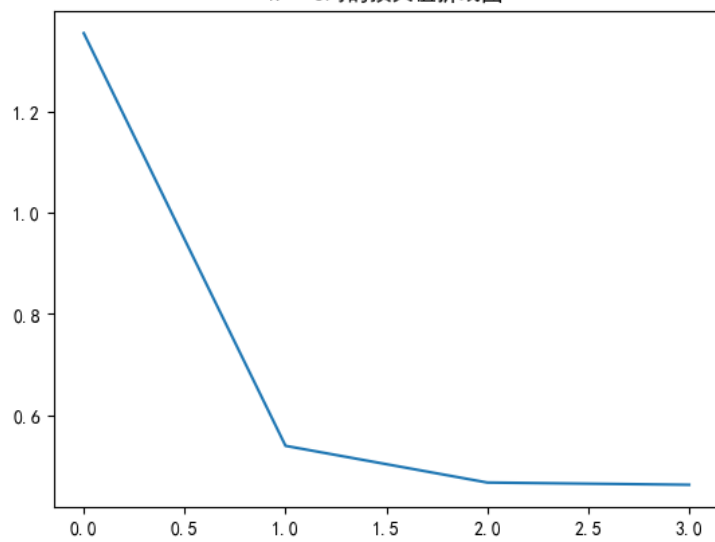


k=3

k = 3时的聚类结果

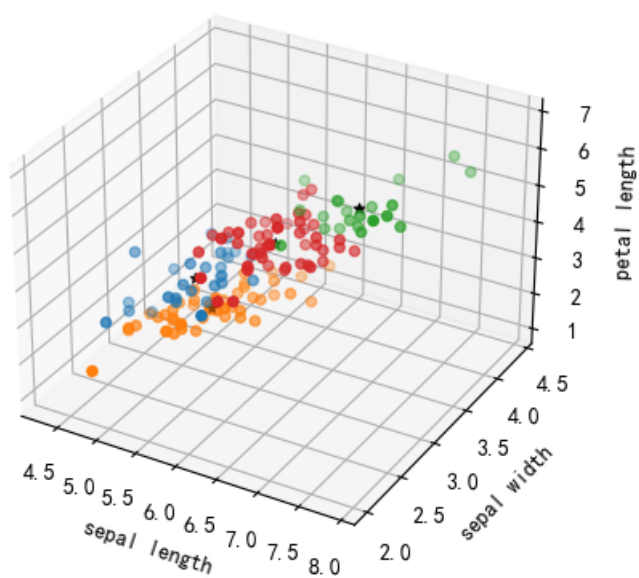


k = 3时的损失值折线图

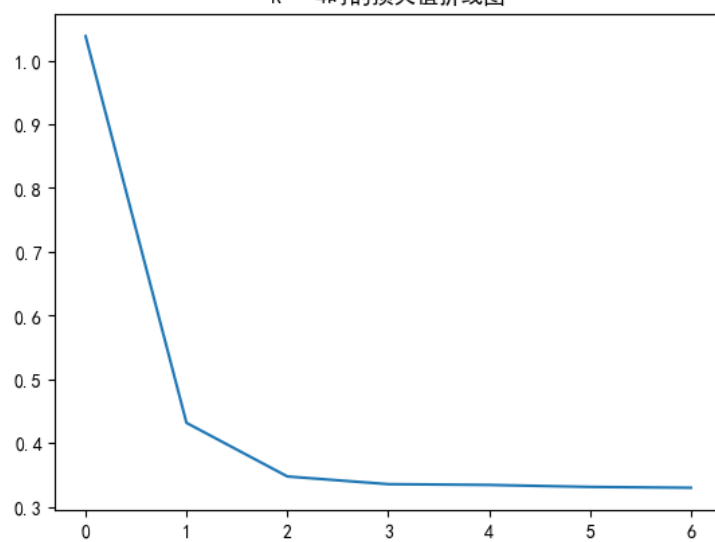


k=4

k = 4时的聚类结果

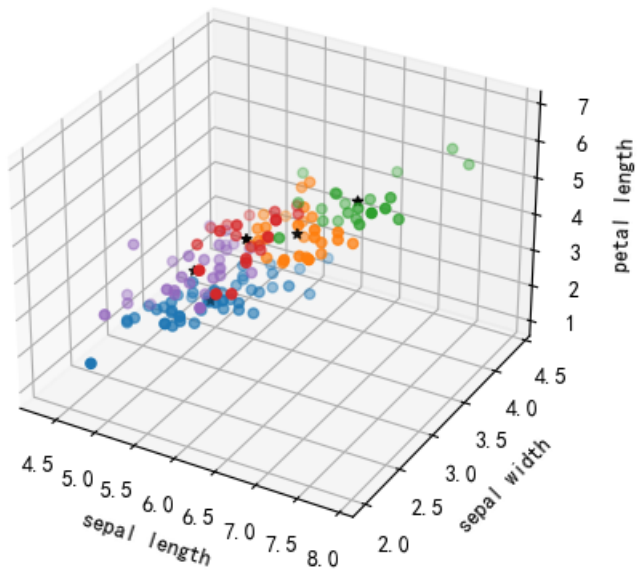


k = 4时的损失值折线图

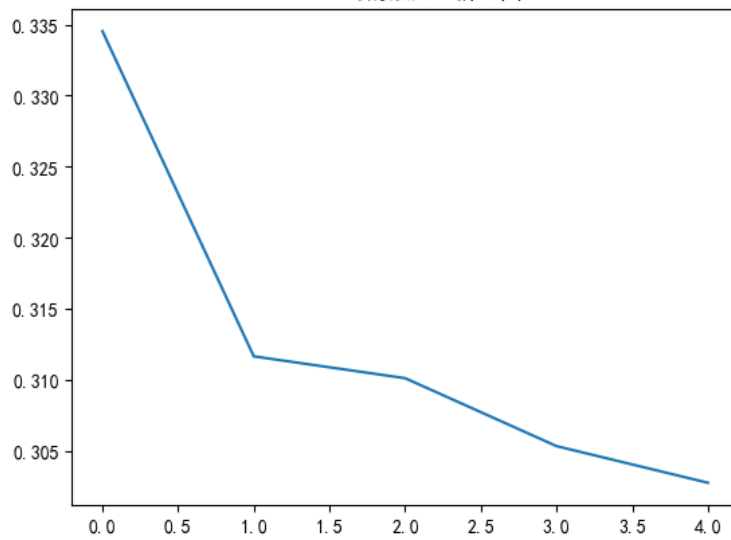


k=5

k = 5时的聚类结果

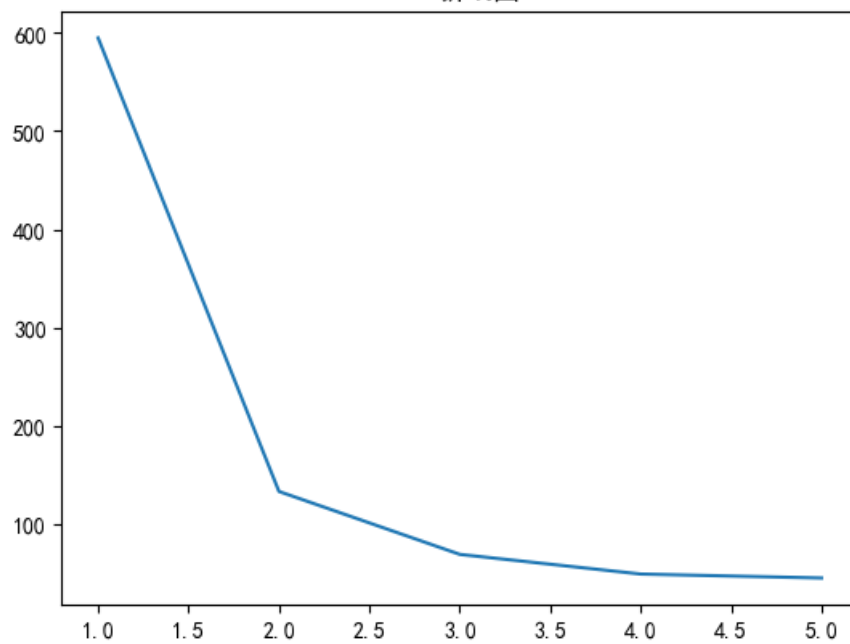


k = 5时的损失值折线图



4. 分析实验结果

SSE折线图



根据SSE折线图我们可以看到，当 $k=3$ 后SSE的下降幅度明显趋向于缓慢，可以看做 $k=3$ 为肘点，所以最佳 k 值为3