

## 实验二 KNN 分类任务

### 1 问题描述

通过 KNN 算法对 Iris 鸢尾花数据集中的测试集进行分类，分别使用欧式距离以及马氏距离作为度量实现，同时实现基于马氏距离的度量学习算法 NCA

### 2 实现步骤与流程

#### KNN 算法

KNN 算法选取训练集中与未见样本的度量最小的  $k$  个样本中类别众数作为预测值

#### NCA 算法

样本的散布矩阵

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T$$

马氏距离

$$d_M(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{S}^{-1} (\mathbf{x}_i - \mathbf{x}_j)}$$

可学习的马氏距离（伪马氏距离）

$$d_M(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{M} (\mathbf{x}_i - \mathbf{x}_j)}$$

其中， $\mathbf{M}$  是半正定矩阵，重写以上距离数学形式

$$d_M(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{A}^T \mathbf{A} (\mathbf{x}_i - \mathbf{x}_j)}$$

其中变换矩阵  $\mathbf{A} \in \mathbb{R}^{k \times d}$ ，并且  $k \geq \text{rank}(\mathbf{M})$ 。设样本  $\mathbf{x}_i$  的近邻分布为

$$p_{ij} = \exp\{-d_M^2(\mathbf{x}_i, \mathbf{x}_j)\} / Z_i$$

$$Z_i = \sum_{k \neq i} \exp\{-d_M^2(\mathbf{x}_i, \mathbf{x}_k)\}$$

代表样本  $\mathbf{x}_j$  与样本  $\mathbf{x}_i$  属于相同类别的概率。针对变换矩阵  $\mathbf{A}$  最大化同类样本相似程度

$$\max_{\mathbf{A}} f(\mathbf{A}) = \sum_{i=1}^n \sum_{j \in \Omega_i} p_{ij}$$

其中  $\Omega_i$  表示和样本  $\mathbf{x}_i$  属于同一类别的其它样本的集合。令

$$d_{ij} = d_M^2(\mathbf{x}_i, \mathbf{x}_j)$$

准则函数对变换矩阵  $\mathbf{A}$  求偏导数得

$$\frac{\partial p_{ij}}{\partial d_{ik}} = \begin{cases} -p_{ij}(1 - p_{ij}) & k = j \\ p_{ij}p_{ik} & k \neq j \end{cases}$$

$$\frac{\partial d_{ik}}{\partial \mathbf{A}} = 2\mathbf{A}(\mathbf{x}_i - \mathbf{x}_k)(\mathbf{x}_i - \mathbf{x}_k)^T$$

$$\frac{\partial p_{ij}}{\partial \mathbf{A}} = 2p_{ij}\mathbf{A} \left[ \sum_{k \neq i} p_{ik}(\mathbf{x}_i - \mathbf{x}_k)(\mathbf{x}_i - \mathbf{x}_k)^T - (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T \right]$$

$$\frac{\partial f}{\partial \mathbf{A}} = \sum_{i=1}^n \sum_{j \in \Omega_i} \frac{\partial p_{ij}}{\partial \mathbf{A}}$$

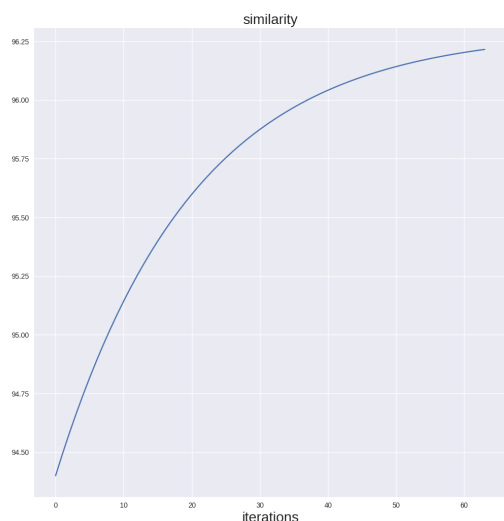
利用梯度下降法即可学习到一个较优的变换矩阵，使得变换后的同类样本距离，异类样本远离。但注意到该优化问题非凸，意味着最终收敛得到的解可能是局部最优解，并且与迭代初始值密切相关，实际实验中需要多次尝试才能得到满足要求的解。

### 3 实验结果与分析

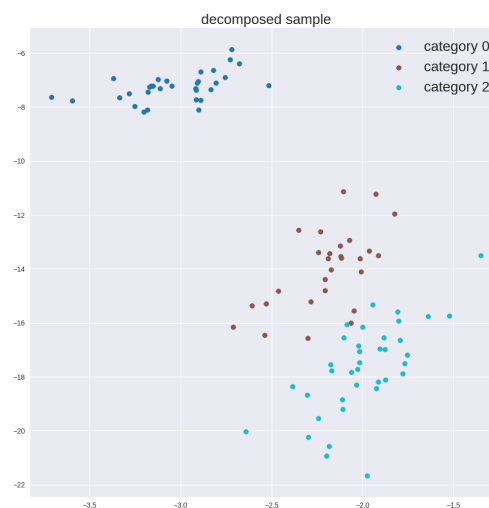
由于实验要求的 MindSpore 框架相关的资料较为匮乏，因此本次实验采取开源机器学习算法库 sklearn 与实验中的算法进行对比。sklearn 是一个开源的基于 python 语言的机器学习工具包。它通过 numpy, scipy 等 python 数值计算的库实现高效的算法应用，并且涵盖了几乎所有主流机器学习算法。

#### NCA 算法

以下是 NCA 算法的度量学习历史以及最终收敛的模型的变换矩阵对样本进行降维后的分布图



(a) 度量学习历史



(b) 降维样本分布

图 1: NCA 算法

观察可得，随着模型的迭代，同类样本的相似度逐渐提高，这点也可以从降维样本上看出

## KNN 算法

以下是实验 KNN 算法（欧式距离、马氏距离）与 sklearn 的 KNN 算法的在验证集上的准确率对比

```
● (python) wzh@wzh:~/workspace/course/pattern recognition/实验/实验2$
euclidean accuracy: 1.0
100%|████████████████████████████████████████████████████████████████████████████████
mahalanobis accuracy: 1.0
```

(a) KNN (实验)

```
● (python) wzh@wzh:~/workspace/course/pattern recognition/实验/实验2$
euclidean accuracy: 1.0
```

(b) KNN (sklearn)

图 2: KNN 算法准确率

经过实验，采用 5 近邻在验证集上可以达到最高的准确率 100%，但是由于样本数量较少，因此结果在一定程度上作为参考

## 4 MindSpore 学习使用心得体会

由于本次实验并未采用 MindSpore，此部分用 sklearn 的学习使用心得体会来代替。

本次实验中调用了 sklearn 中的以下算法接口

- KNN (neighbors KNeighborsClassifier)

作为实验算法的对照算法，sklearn 提供的算法接口方便调用，并且执行效率较高，能够高效地实现实验的需求。

## 5 代码附录