

Fisher线性判别分析

降维作为一种减少特征冗余的方法，也可以应用在线性分类当中。在K分类问题中，Fisher线性判别分析通过最大化类间方差和最小化类内方差，将数据映射到K-1维空间进行分类。本文将着重讨论推导多分类的情况。并在最后给出了降维之后如何对新样本进行分类的方法和建议。

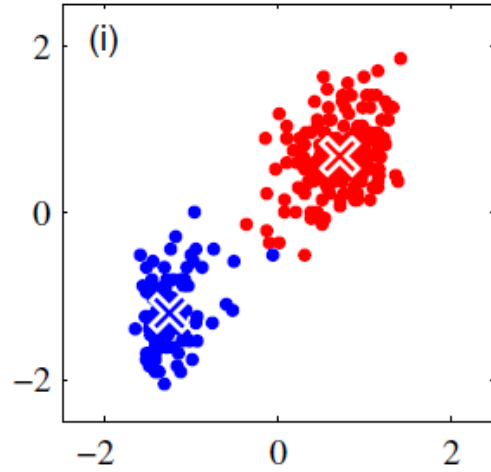
- Fisher线性判别分析
 - 1. 符号标识
 - 2. 散度矩阵(Scatter Matrices)
 - 3. 二分类求解
 - 4. 多分类求解
 - 5. 分类方法
 - 6. 代码实现
 - 后记
 - 参考文献

1. 符号标识

符号	意义
N_k	属于第K类的样本数量
N	样本总数
K	类别总数
$\mathbf{x} \in \mathbb{R}^D$	D维样本向量
$\mathbf{X} \in \mathbb{R}^{N \times D}$	样本矩阵
$\mathbf{S}_W \in \mathbb{R}^{D \times D}$	类内散度矩阵
$\mathbf{S}_B \in \mathbb{R}^{D \times D}$	类间散度矩阵
$\mathbf{W} \in \mathbb{R}^{D \times K-1}$	投影矩阵
$\mathbf{y} \in \mathbb{R}^{K-1}$	投影后样本向量
$\mathbf{u} \in \mathbb{R}^{K-1}$	投影后样本均值
$\mathbf{P}_W \in \mathbb{R}^{K-1 \times K-1}$	投影后类内散度矩阵
$\mathbf{P}_B \in \mathbb{R}^{K-1 \times K-1}$	投影后类间散度矩阵
$Tr(\cdot)$	矩阵的迹

2. 散度矩阵(Scatter Matrices)

考虑一个问题，假设有2个类别，并且不同类别的样本相互杂糅，而我们需要将这两类样本非常清晰的区分开来，需要怎么做？可以使用一种线性变换将纠缠交错的样本给分离开来，那么如何去度量变换之后是否达到了我们的预期呢？我们可以从下面这张图得到一些灵感。



我们可以发现，这两个类的中心距离相差较大，并且在每个类别中，各样本的分布非常紧凑。于是我们可以想到这么一个度量方法。1.使得类与类之间的距离最大化，这可以通过两个类直接的均值进行度量。2.使类内的方差最小化，即同类样本分布紧凑。

上述是以2分类为例的，那么拓展到 K 分类呢($K > 2$)? 如果我们仍采用度量类与类之间的距离的方法，那么将需要最大化 $\frac{K(K-1)}{2}$ 个距离，这样的复杂度是无法接受的。实际上，我们从上图也可以观察到，类中心之间还可以产生一个中心，就是所有样本的均值。那么度量类之间的距离也可以转化为度量个均值之间的距离，那么需要最大化 K 个距离，这依然是一个比较复杂的问题。于是，参考对类内方差的最小化，我们可以用一个等价的方法去最大化类与类之间的距离：最大化类之间的方差。为了实现上述方法，我们定义散度矩阵：

类内散度矩阵

$$\mathbf{S}_W = \sum_{k=1}^K \mathbf{\Sigma}_k \dots\dots (2.1)$$

$$\mathbf{\Sigma}_k = \sum_{n \in C_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T \dots\dots (2.2)$$

其中

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}_n \dots\dots (2.3)$$

类间散度矩阵

$$\mathbf{S}_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T \dots\dots (2.4)$$

其中

$$\mathbf{m}_k = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} \sum_{k=1}^K N_k \mathbf{m}_k \dots\dots (2.5)$$

可得混合散度矩阵(the mixture scatter matrix)

$$\mathbf{S}_M = \mathbf{S}_W + \mathbf{S}_B \dots\dots (2.6)$$

3. 二分类求解

$$\text{设均值向量 } \mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in C_1} \mathbf{x}_n, \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in C_2} \mathbf{x}_n, \text{投影方向为 } \mathbf{w}, y_n = \mathbf{w}^T \mathbf{x}_n.$$

投影后均值为 $u_1 = \mathbf{w}^T \mathbf{m}_1, u_2 = \mathbf{w}^T \mathbf{m}_2$, 类间距离 $(u_1 - u_2)^2$.

类内总方差 $s_1^2 = \sum_{n \in C_1} (y_n - u_1)^2, s_2^2 = \sum_{n \in C_2} (y_n - u_2)^2$. 最大化类间距离, 最小化类内方差。

$$J = \frac{(u_1 - u_2)^2}{s_1 + s_2} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

其中,

$$\mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T, \mathbf{S}_W = \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in C_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T$$

J 对 \mathbf{w} 求导可得

$$\frac{\partial J}{\partial \mathbf{w}} = \frac{2(\mathbf{S}_B \mathbf{w} \mathbf{w}^T \mathbf{S}_W \mathbf{w} - \mathbf{w}^T \mathbf{S}_B \mathbf{w} \mathbf{S}_W \mathbf{w})}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} = 0$$

$$\mathbf{S}_W \mathbf{w} \propto \mathbf{S}_B \mathbf{w} = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} \propto (\mathbf{m}_1 - \mathbf{m}_2)$$

$\rightarrow \mathbf{w} \propto \mathbf{S}_W^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$, 我们只关注投影的方向, 也可以给其范数加上约束。

4. 多分类求解

由于我们有 K 个类别, 根据贝叶斯分类器对此类问题的处理, 是得到 K 个后验概率 $p_1(\mathbf{x}) \dots p_K(\mathbf{x})$, 然而我们知道 $\sum_i p_i = 1$, 因此, 只有 $K - 1$ 个是线性无关的。那么我们讲 D 维样本空间映射到 $K - 1$ 维空间是没有分类信息的损失的。

于是, 有线性映射

$$\mathbf{y} = \mathbf{W}^T \mathbf{x} \dots\dots (4.1)$$

$$\mathbf{P}_W = \mathbf{W}^T \mathbf{P}_W \mathbf{W} \dots\dots (4.2)$$

$$\mathbf{P}_B = \mathbf{W}^T \mathbf{P}_B \mathbf{W} \dots\dots (4.3)$$

在二分类时的思想是最大化类间方差, 最小化类内方差, 于是可得二分类时的损失函数

$$J(\mathbf{W}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

与之不同的是，多分类情况下分子分母都是矩阵而不是标量，且矩阵没有除法，因此需要采用另一种判别准则。

判别准则有多种，我们这里使用其中一种。可以先从直觉上理解，具体是为什么等我明白了再补充吧。

$$J(\mathbf{W}) = \text{Tr}(\mathbf{P}_W^{-1} \mathbf{P}_B) = \text{Tr}((\mathbf{W}^T \mathbf{S}_W \mathbf{W})^{-1} \mathbf{W}^T \mathbf{S}_B \mathbf{W}) \dots \dots (4.4)$$

对其求微分可得

$$\begin{aligned} d(J) &= \text{Tr}[d(\mathbf{W}^T \mathbf{S}_W \mathbf{W})^{-1} \mathbf{W}^T \mathbf{S}_B \mathbf{W}] + \text{Tr}[2(\mathbf{W}^T \mathbf{S}_W \mathbf{W})^{-1} \mathbf{W}^T \mathbf{S}_B d(\mathbf{W})] \\ &= \text{Tr}[-(\mathbf{W}^T \mathbf{S}_W \mathbf{W})^{-1} d(\mathbf{W}^T \mathbf{S}_W \mathbf{W})(\mathbf{W}^T \mathbf{S}_W \mathbf{W})^{-1} \mathbf{W}^T \mathbf{S}_B \mathbf{W}] + \text{Tr}[2(\mathbf{W}^T \mathbf{S}_W \mathbf{W})^{-1} \mathbf{W}^T \mathbf{S}_B d(\mathbf{W})] \\ &= \text{Tr}[-2\mathbf{P}_W^{-1} \mathbf{W}^T \mathbf{S}_W d(\mathbf{W}) \mathbf{P}_W^{-1} \mathbf{P}_B] + \text{Tr}[2\mathbf{P}_W^{-1} \mathbf{W}^T \mathbf{S}_B d(\mathbf{W})] \\ &= \text{Tr}[(-2\mathbf{P}_W^{-1} \mathbf{P}_B \mathbf{P}_W^{-1} \mathbf{W}^T \mathbf{S}_W + 2\mathbf{P}_W^{-1} \mathbf{W}^T \mathbf{S}_B) d(\mathbf{W})] \end{aligned} \quad (4.5)$$

可得

$$\frac{\partial J}{\partial \mathbf{W}} = -2\mathbf{S}_W \mathbf{W} \mathbf{P}_W^{-1} \mathbf{P}_B \mathbf{P}_W^{-1} + 2\mathbf{S}_B \mathbf{W} \mathbf{P}_W^{-1} = 0 \dots \dots (4.6)$$

$$\mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{W} = \mathbf{W} \mathbf{P}_W^{-1} \mathbf{P}_B \dots \dots (4.7)$$

式4.7的形式容易与矩阵的特征值联系起来。式中的散度矩阵 \mathbf{S}_B 是不满秩的，它是由 K 个秩为1的矩阵相加得到的，而在式2.3的约束下，只有 $K - 1$ 个矩阵是线性无关的，因此它的秩最多为 $K - 1$ 。而 \mathbf{S}_W 是满秩的，则 $\mathbf{S}_W^{-1} \mathbf{S}_B$ 只有 $K - 1$ 个非零特征值。

命题1：存在一个线性变换 $\mathbf{Q} \in \mathbb{R}^{K-1 \times K-1}$ 且 \mathbf{Q}^{-1} 存在，使得

$$\begin{aligned} \mathbf{Q}^T \mathbf{P}_W \mathbf{Q} &= \mathbf{I}, \mathbf{Q}^T \mathbf{P}_B \mathbf{Q} = \mathbf{\Lambda} \dots \dots (4.8) \\ \text{其中, } \mathbf{I} &\text{为} K - 1 \text{阶单位矩阵, } \mathbf{\Lambda} \text{为} K - 1 \text{阶对角矩阵。} \end{aligned}$$

证明：

$\because \mathbf{P}_W$ 正定，存在 \mathbf{C} 使得 $\mathbf{C}^T \mathbf{P}_W \mathbf{C} = \mathbf{I}$ ，而 $\mathbf{C}^T \mathbf{P}_B \mathbf{C}$ 是实对称矩阵，从而存在正交变换

$$\begin{aligned} \mathbf{D}^T (\mathbf{C}^T \mathbf{P}_B \mathbf{C}) \mathbf{D} &= \text{diag}(\lambda_1, \dots, \lambda_{K-1}) \\ \text{令 } \mathbf{Q} &= \mathbf{C} \mathbf{D}, \text{ 则 } \mathbf{Q}^T \mathbf{P}_W \mathbf{Q} = \mathbf{D}^T (\mathbf{C}^T \mathbf{P}_W \mathbf{C}) \mathbf{D} = \mathbf{D}^T \mathbf{D} = \mathbf{I}, \text{ 得证。} \end{aligned}$$

将式4.8代入式4.7可得

$$\mathbf{S}_W^{-1} \mathbf{S}_B (\mathbf{W} \mathbf{Q}) = (\mathbf{W} \mathbf{Q}) \mathbf{\Lambda} \dots \dots (4.9)$$

可以发现， $\mathbf{\Lambda}$ 不仅是 \mathbf{P}_B 的特征值矩阵，还是 $\mathbf{S}_W^{-1} \mathbf{S}_B$ 的特征值矩阵。则有，

$$\begin{aligned} J(\mathbf{W}) &= \text{Tr}(\mathbf{P}_W^{-1} \mathbf{P}_B) = \sum_i^{K-1} \lambda_i \\ \text{Tr}(\mathbf{S}_W^{-1} \mathbf{S}_B) &= \sum_i^D \mu_i \dots \dots (4.10) \end{aligned}$$

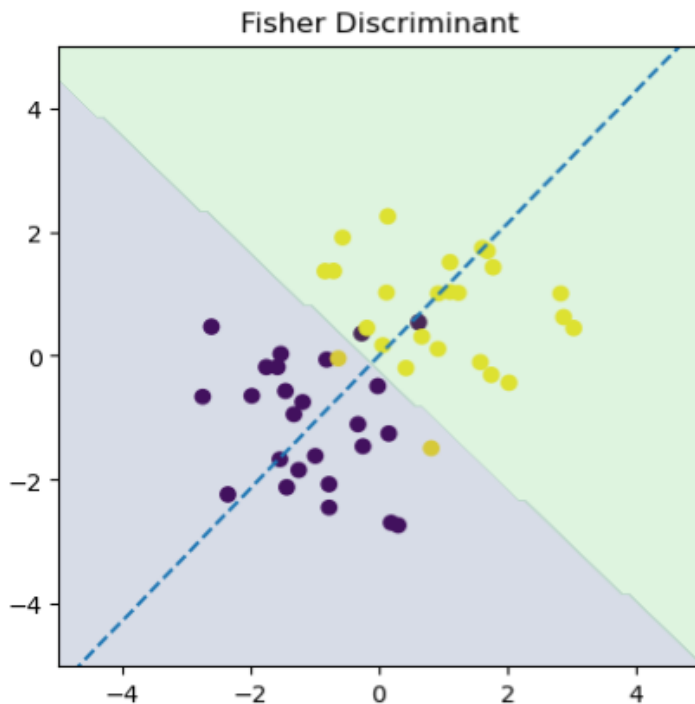
注意，这里 $\mathbf{S}_W^{-1}\mathbf{S}_B$ 是我们可以通过观测到的样本计算出来的，所以特征值是确定的 $\mu_i, i = 1, \dots, D$,式4.10给出了与目标函数之间的关系，并且由正交变换的不变性，我们可得知 \mathbf{W} 就是由 $\mathbf{S}_W^{-1}\mathbf{S}_B$ 最大的 $K - 1$ 个特征值对应的特征向量构成的。

5. 分类方法

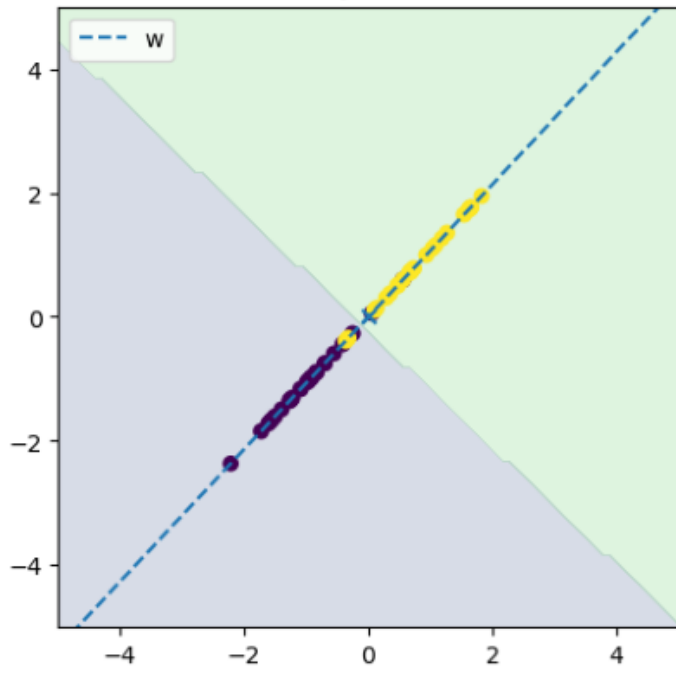
当我们得到投影方向之后，我们可以线性映射到 $K - 1$ 维空间。那么映射之后，该如何对未知点进行预测呢？可以用以下方法。

1. 最简单的方法，是降维之后，计算点到各个类别中心的距离，并将点分类为最近中心所属的类别，这与K-means算法的思想是相似的。但实际上，这与K-means一样都是不考虑分布的方差的，只考虑每个类别的均值，因为这等价于K个方差趋近0的高斯分布的混合。
2. 在第1种方法的基础上进行改进，把分布的方差也考虑进去，那么我们可以对投影后的数据点用高斯分布进行建模，但是我们注意到高斯分布是单峰的，而数据在簇中的分布不一定是单峰的，而real-world data往往是复杂的呈现出多峰的分布，因此用单一的单峰的分布建模效果会受到影响。
3. 承接上述两点，我们可以用高斯混合模型来解决分布是多峰的问题，但实际上我们很难去判断这个数据的分布到底有几个峰值，这决定了我们使用几个高斯分量去构建混合模型，那么高斯分量的数量将会是一个超参数。

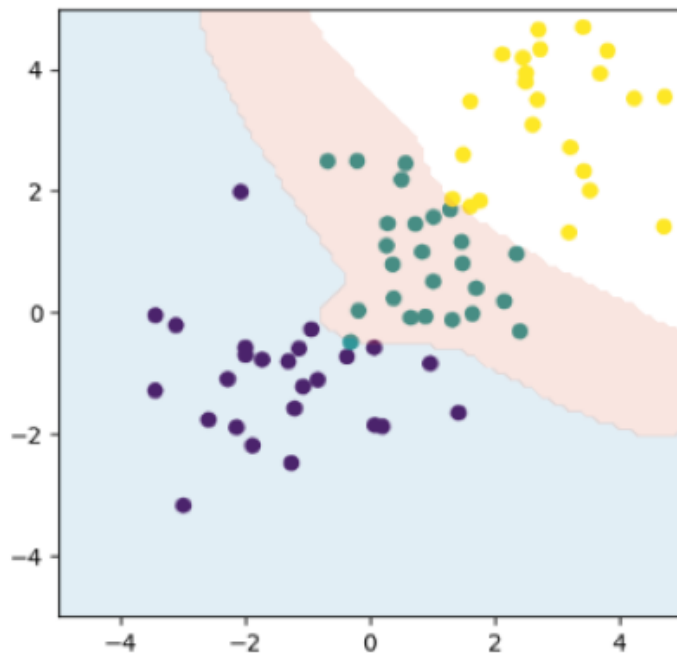
6. 代码实现



Projection



Fisher Discriminant



二分类

```
class FisherLinearDiscriminant:
    """
    Only for 2 classes
    """
    def __init__(self, w=None, threshold=None):
        self.w = w
        self.threshold = threshold

    def fit(self, x_train: np.ndarray, y_train: np.ndarray):
        x0 = x_train[y_train == 0]
        x1 = x_train[y_train == 1]
        u1 = np.mean(x0, axis=0)
        u2 = np.mean(x1, axis=0)
        cov = np.cov(x0, rowvar=False) + np.cov(x1, rowvar=False)
        w = np.linalg.inv(cov) @ (u2 - u1)
        self.w = w / np.linalg.norm(w)
        g0 = Gaussian()
        g0.fit(x0 @ self.w)
        g1 = Gaussian()
        g1.fit(x1 @ self.w)
        x = np.roots([g1.var - g0.var,
                      2*(g1.mean*g0.var - g0.mean*g1.var),
                      g1.var * g0.mean ** 2 - g0.var * g1.mean ** 2
                      - g1.var * g0.var * np.log(g1.var / g0.var)
                      ])
        if g0.mean < x[0] < g1.mean or g1.mean < x[0] < g0.mean:
            self.threshold = x[0]
        else:
            self.threshold = x[1]

    def project(self, x: np.ndarray):
        return x @ self.w

    def classify(self, x: np.ndarray):
        return (x @ self.w > self.threshold).astype(int)

class MultiFisherLinearDiscriminant:
    """
    For K >= 2
    """
    def __init__(self, W=None, threshold=None, n_classes=3, peaks: Iterable = None):
        self.W = W
        self.threshold = threshold
        self.n_classes = n_classes
        self.peaks = peaks or [2] * n_classes
        assert len(self.peaks) == self.n_classes, "peaks shape error"

    def fit(self, x_train: np.ndarray, y_train: np.ndarray):
        cov_b = [] # between
        cov_w = [] # within
        mean = []
        mu = x_train.mean(0, keepdims=True) # 1 D
        for k in range(self.n_classes):
            x_k = x_train[y_train == k] # N_k D
            mean_k = np.mean(x_k, axis=0, keepdims=True) # 1 D
            mean.append(mean_k)
            dist = x_k - mean_k # N_K D
```

```

        cov_k = dist.T @ dist
        cov_w.append(cov_k)
        dist = mean_k - mu
        cov_k = (y_train == k).sum() * dist.T @ dist
        cov_b.append(cov_k)
    self.mean = np.concatenate(mean, axis=0)
    cov_b = np.sum(cov_b, 0) # D D
    cov_w = np.sum(cov_w, 0)
    A = np.linalg.inv(cov_w) @ cov_b
    _, vectors = np.linalg.eig(A)
    self.W = vectors[:, -(self.n_classes - 1):] # D, K-1
    x_prj = x_train @ self.W
    self.__getDistributions(x_prj, y_train)

def __getDistributions(self, x_train, y_train):
    distributions = []
    for k in range(self.n_classes):
        x_k = x_train[y_train == k] # N_k D
        gmm = GaussianMixture(classes=self.peaks[k]) if self.peaks[k] > 1 else Gaussian()
        gmm.fit(x_k)
        distributions.append(gmm)
    self.distributions = distributions

def project(self, x: np.ndarray):
    assert x.ndim <= 2, "ndim should be less than 3"
    if x.ndim == 1:
        x = x[None, :]
    return x @ self.W # N, K-1

def classify(self, x: np.ndarray):
    assert x.ndim <= 2, "ndim should be less than 3"
    probs = []
    if x.ndim == 1:
        x = x[None, :]
    x = self.project(x)
    for i in range(x.shape[0]):
        probs.append(np.concatenate([gmm._pdf(x[i]) for gmm in self.distributions]))
    classes = np.argmax(probs, axis=-1)
    return classes

```


#测试

```
def create_data(size=50, add_outlier=False, add_class=False):
    assert size % 2 == 0
    x0 = np.random.normal(size=size).reshape(-1, 2) - 1
    x1 = np.random.normal(size=size).reshape(-1, 2) + 1
    if add_outlier:
        x = np.random.normal(size=10).reshape(-1, 2) + np.array([5, 10])
        return np.concatenate([x0, x1, x]), np.concatenate([np.zeros(size//2), np.ones(size//2 + 5)])
    if add_class:
        x = np.random.normal(size=size).reshape(-1, 2) + 3
        return np.concatenate([x0, x1, x]), np.concatenate([np.zeros(size//2), np.ones(size//2), 2*np.ones(size//2)])
    return np.concatenate([x0, x1]), np.concatenate([np.zeros(size//2), np.ones(size//2)])
```

```
model = FisherLinearDiscriminant()
```

```
model.fit(x_train, y_train)
```

```
plt.scatter(x_train[:, 0], x_train[:, 1], c=y_train)
x1_test, x2_test = np.meshgrid(np.linspace(-5, 5, 100), np.linspace(-5, 5, 100))
x_test = np.concatenate([x1_test, x2_test]).reshape(2, -1).T
y_pred = model.classify(x_test)
x = np.linspace(-5, 5, 20)
plt.contourf(x1_test, x2_test, y_pred.reshape(100, -1), alpha=0.2, levels=np.linspace(0,1,3))
plt.plot(x, x * model.w[1]/model.w[0], label='w', linestyle='--')
plt.title('Fisher Discriminant')
plt.gca().set_aspect('equal', adjustable='box')
plt.xlim(-5, 5)
plt.ylim(-5, 5)
plt.show()
```

```
plt.plot(x, x * model.w[1]/model.w[0], label='w', linestyle='--')
w = model.w
rollmat = np.zeros((2,2))
div = np.sqrt(w[0]**2 + w[1]**2)
rollmat[0,0] = w[0]/div
rollmat[0,1] = w[1]/div
rollmat[1,0] = -w[1]/div
rollmat[1,1] = w[0]/div
x_proj = x_train@w
x_proj = np.concatenate([x_proj[:,None], np.zeros_like(x_proj[:,None])],axis=-1).reshape(-1, 2)
#plt.scatter(x_proj[:,0], x_proj[:,1]-5, c=y_train)
x_roll = x_proj @ rollmat
plt.contourf(x1_test, x2_test, y_pred.reshape(100, -1), alpha=0.2, levels=np.linspace(0,1,3))
plt.scatter(x_roll[:, 0], x_roll[:,1], c=y_train)
plt.scatter(0, 0, marker='x', alpha=1)
plt.title('Projection')
plt.gca().set_aspect('equal', adjustable='box')
plt.xlim(-5, 5)
plt.ylim(-5, 5)
plt.legend()
plt.show()
```

后记

有些地方还没整明白，明白了再回来补充。

参考文献

- [1] Fukunaga, K. (1990). Introduction to Statistical Pattern Recognition (Second ed.). Academic Press. 441-454.
- [2] Christopher M. Bishop.(2007). Pattern Recognition and Machine Learning. 187-192.