

线性判别分析 (LDA)

简介

线性判别分析(Linear Discriminant Analysis, **LDA**), 于1936年由Ronald Fisher首次提出, 因此, 也叫做Fisher线性判别(Fisher Linear Discriminant, **FLD**), 或者Fisher Discriminant Analysis (**FDA**), 是模式识别的经典算法, 并在1996年由Belhumeur引入模式识别和人工智能领域。

基本思想

将高维的模式样本投影到最佳鉴别矢量空间，以达到抽取分类信息和压缩特征空间维数的效果。

投影后保证模式样本在新的子空间有最大的类间距离和最小的类内距离，即模式在该空间中有最佳的可分离性。

因此，它是一种有效的特征抽取方法。使用这种方法能够使投影后模式样本的类间散布矩阵最大，并且同时类内散布矩阵最小。

LDA与PCA

LDA与PCA都是常用的降维技术

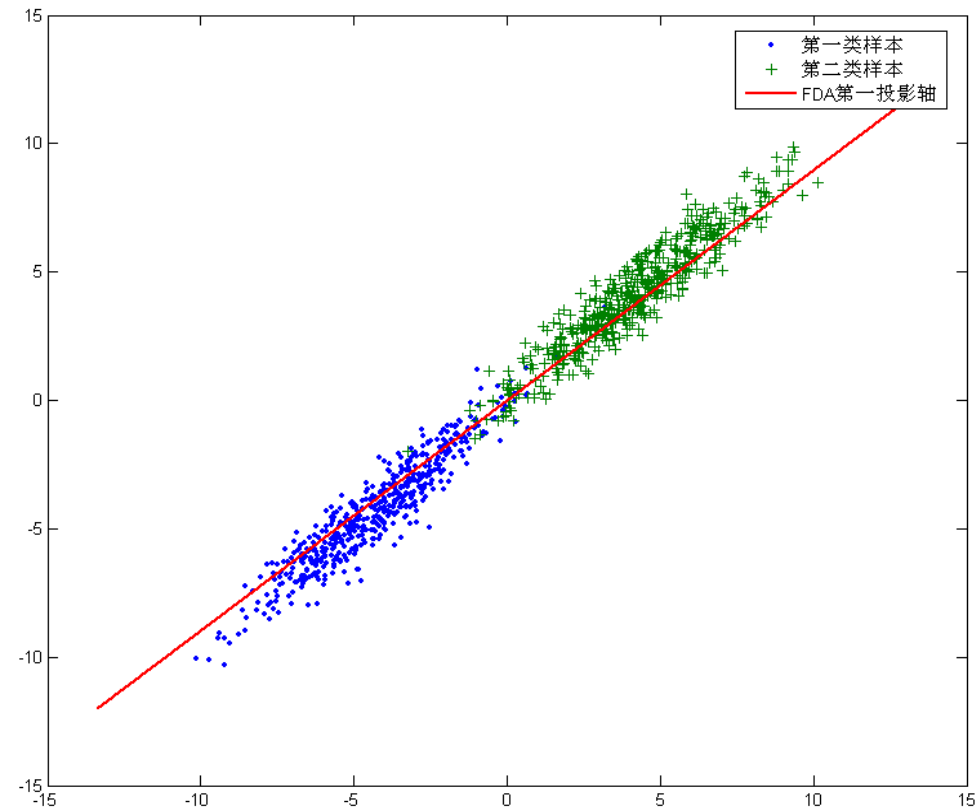
PCA是一种无监督的降维方法;

主要是从特征的协方差角度,
将数据投影到方差最大的几个相互正交的方向上。

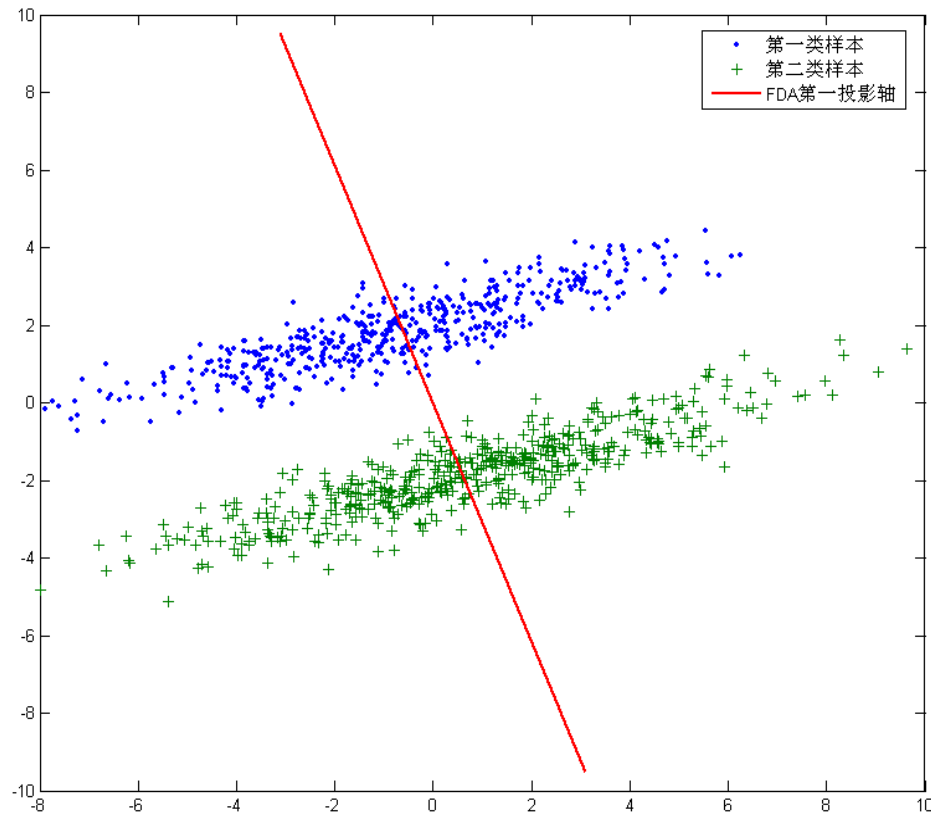
LDA是一种有监督的降维方式;

希望投影后不同类别之间数据点的距离更大, 同一类别的数据点更紧凑。

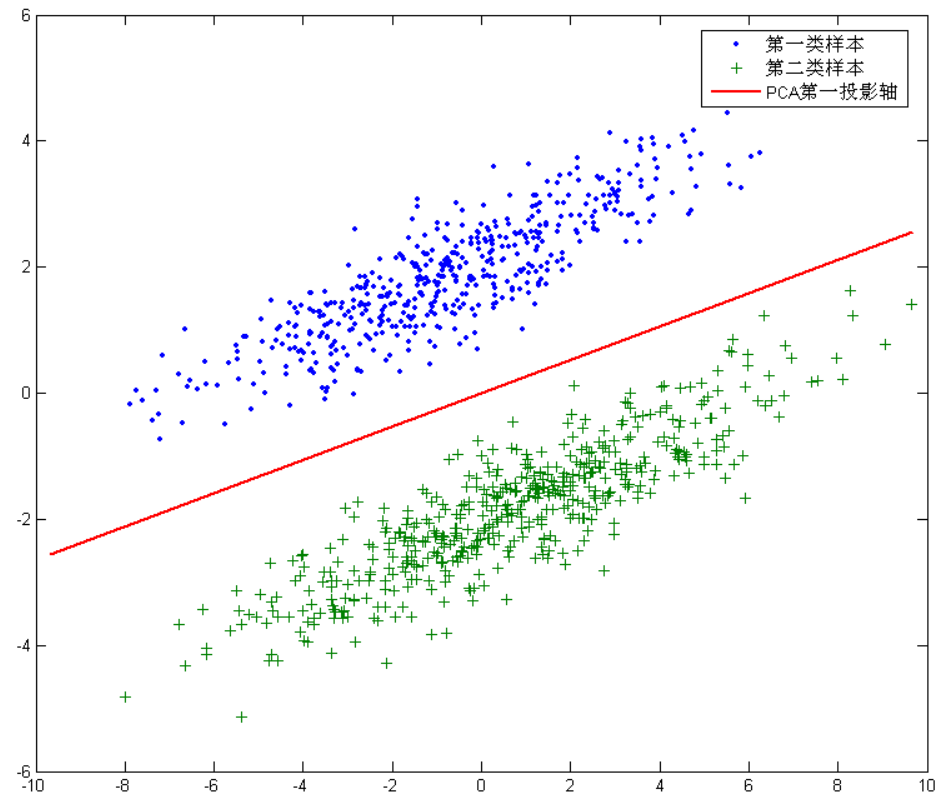
Linear discriminant analysis(LDA,FDA)



Projection axis of LDA

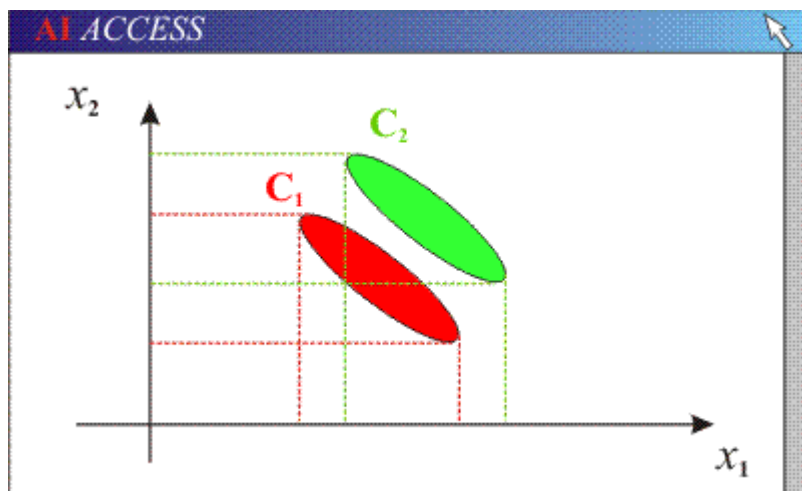


However, PCA performs badly in this case



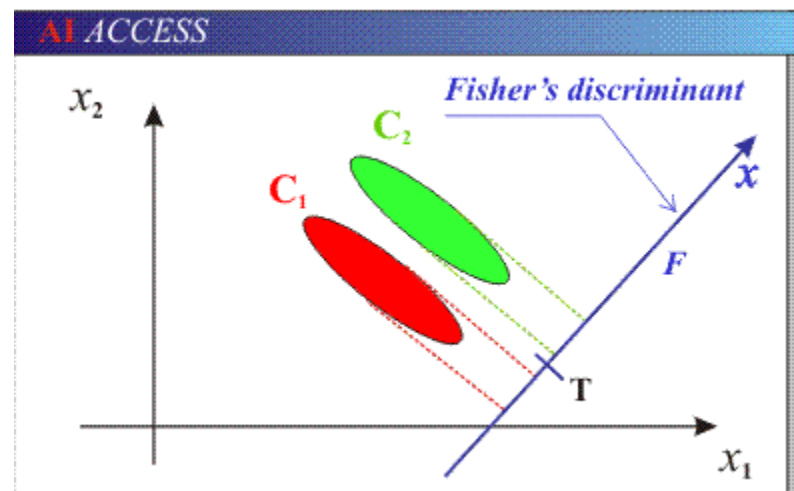


LDA的目标:



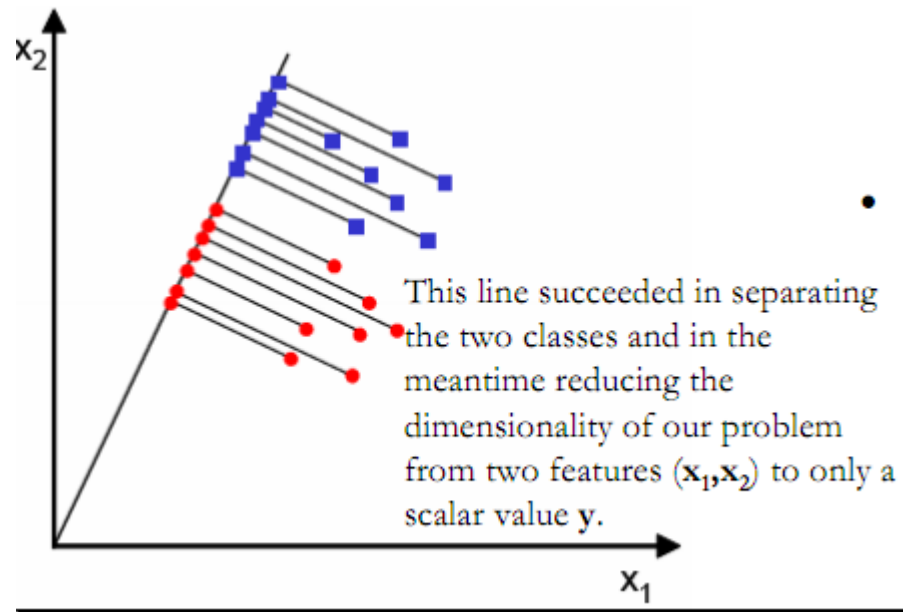
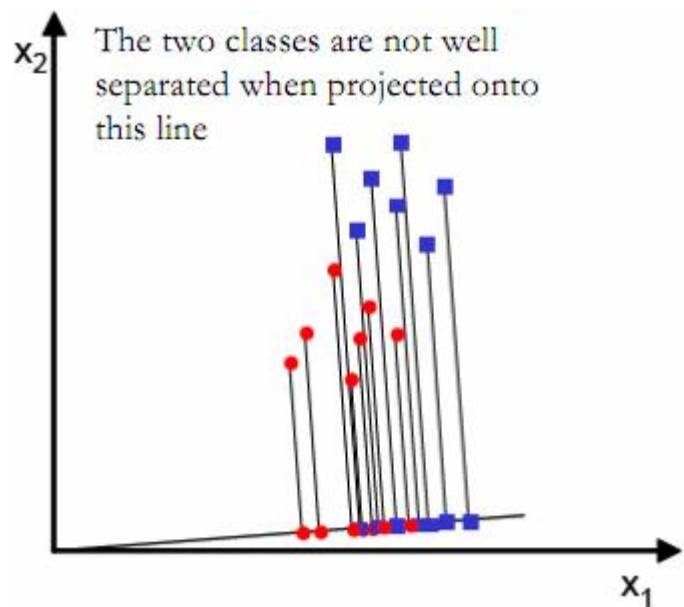
有绿色、红色两个类别。上图是两个类别的原始数据。

如果将数据从二维降维到一维。直接投影到 x_1 轴或者 x_2 轴，不同类别之间会有重复，导致分类效果下降。



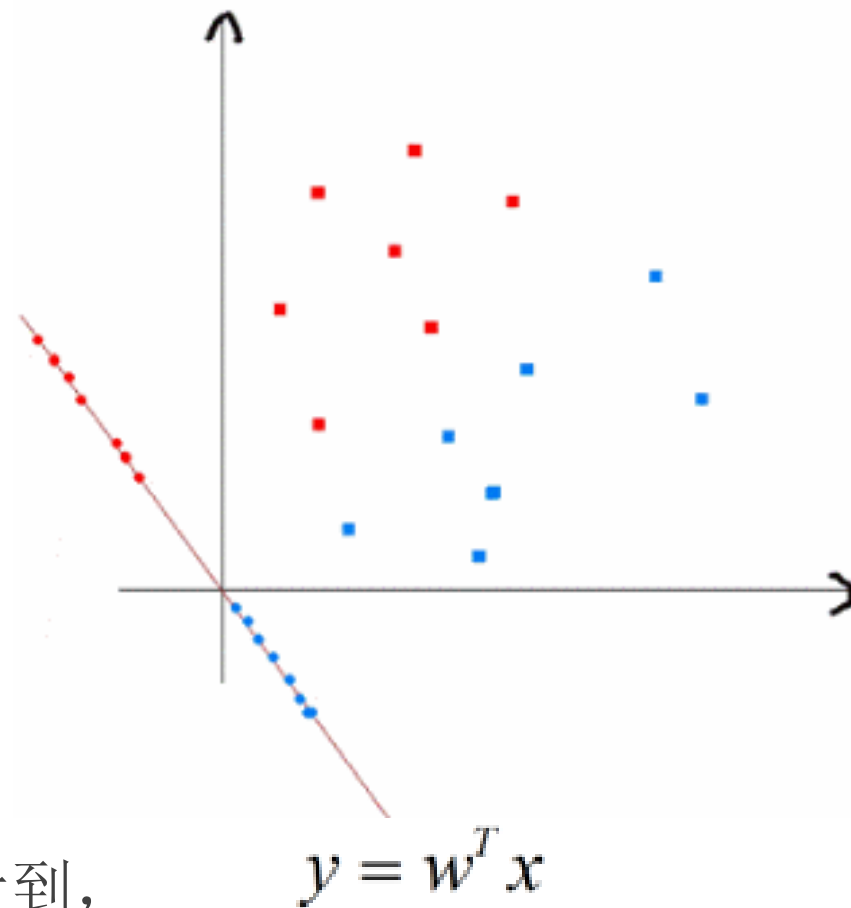
采用LDA方法将其映射到上图直线上

红色类别和绿色类别在映射之后之间的距离是最大的，而且每个类别内部点的离散程度是最小的（或者说聚集程度是最大的）。



二分类中

LDA的目标是：给出一个标注了类别的数据集，投影到了一条直线之后，能够使得点尽量的按类别区分开，当 $k=2$ 即二分类问题的时候，如图所示：



经过原点的直线就是投影的直线，从图上可以清楚的看到，红色的点和蓝色的点被原点明显的分开了。

用来区分二分类的直线（投影函数)为:

$$y = w^T x$$

类别i的中心点(均值)为: (D_i 表示属于类别i的点):

$$m_i = \frac{1}{n_i} \sum_{x \in D_i} x$$

类别i投影后的中心点为:

$$\widetilde{m}_i = w^T m_i$$

衡量类别i投影后，类别点之间的分散程度（方差）为:

$$\widetilde{s}_i^2 = \sum_{y \in Y_i} (y - \widetilde{m}_i)^2$$

不忘初心，方得始终

LDA分类的目标是使得不同类别之间的距离越大越好，同一类别之中的距离越小越好。

可以得到如下一个公式，表示LDA投影到 w 后的最大化目标优化函数：

$$J(w) = \frac{|\widetilde{m}_1 - \widetilde{m}_2|^2}{\widetilde{s}_1^2 + \widetilde{s}_2^2}$$

啥？啥？啥？



我能怎么办？
我也很绝望啊



$$J(w) = \frac{|\widetilde{m}_1 - \widetilde{m}_2|^2}{\widetilde{s}_1^2 + \widetilde{s}_2^2}$$

分母表示每一个类别内的方差之和，方差越大表示一个类别内的点越分散。（目标是类内投影点尽可能接近）

分子为两个类别各自的中心点的距离的平方。（目标是类间投影点尽可能远）

所以，最大化 $J(w)$ 就可以求出最优的 w

定义一个投影前的各类别分散程度的矩阵 S_i :
$$S_i = \sum_{x \in D_i} (x - m_i)(x - m_i)^T$$

则 $J(w)$ 分母化为:

$$\tilde{S}_i^2 = \sum_{x \in D_i} (w^T x - w^T m_i)^2 = \sum_{x \in D_i} w^T (x - m_i)(x - m_i)^T w = w^T S_i w$$

$$\tilde{S}_1^2 + \tilde{S}_2^2 = w^T (S_1 + S_2) w = w^T S_w w$$

S_w 被称为类内散度矩阵

同样的将 $J(w)$ 分子化为:

S_B 被称为类间散度矩阵

$$|\widetilde{m}_1 - \widetilde{m}_2|^2 = w^T (m_1 - m_2)(m_1 - m_2)^T w = w^T S_B w$$

目标优化函数可以化成下面的形式:

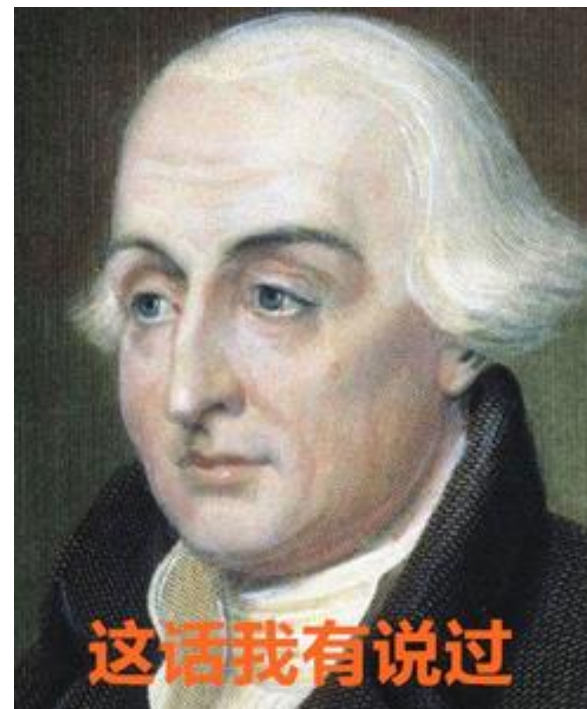
$$J(w) = \frac{w^T S_B w}{w^T S_w w}$$

如何得到 w 呢?



最棒的拉格朗日乘子法来了

但是还有一个问题，如果分子、分母是都可以取任意值的，那就会使得有无穷解，不失一般性，我们将分母限制为长度为1（这是用拉格朗日乘子法一个很重要的技巧，如果忘记了，请复习一下高数），并作为拉格朗日乘子法的限制条件。



可得到:

$$J(w) = \frac{w^T S_B w}{w^T S_w w}$$

$$c(w) = w^T S_B w - \lambda(w^T S_w w - 1)$$

$$\Rightarrow \frac{dc}{dw} = 2S_B w - 2\lambda S_w w = 0$$

$$\Rightarrow S_B w = \lambda S_w w$$

直接可得出

$$S_W^{-1} S_{BW} = \lambda w$$

w就是矩阵 $S_W^{-1} S_E$ 的特征向量。



过来，我告诉你
一个秘密



呕 天哪 我真是怕

别忘了

$$|\widetilde{m}_1 - \widetilde{m}_2|^2 = w^T (m_1 - m_2)(m_1 - m_2)^T w = w^T S_B w$$

得到

$$S_B W = (m_1 - m_2) (m_1 - m_2)^T W$$

$(m_1 - m_2)^T W$ 是一个标量哦

so!!!

$S_B W$ 与 $(m_1 - m_2)$ 同方向

可令

$$S_B W = \lambda (m_1 - m_2)$$

$$S_B W = \lambda (m_1 - m_2)$$



$$S_B w = \lambda S_w w$$

$$W = S_W^{-1} (m_1 - m_2)$$



这么简单呢

对于N(N>2)分类的问题，就可以直接写出以下的结论：

$$S_W = \sum_{i=1}^c S_i$$

$$S_B = \sum_{i=1}^c n_i (m_i - m)(m_i - m)^T$$

$$S_B w_i = \lambda S_W w_i$$

这同样是一个求特征值的问题，求出的第i大的特征向量，即为对应的 w_i 。

S_W^{-1} 一直都存在吗?



S_W 为奇异矩阵时可用奇异值分解得到的。

- 同样我们可以先使用PCA

- 1) PCA 用来降低数据的维度

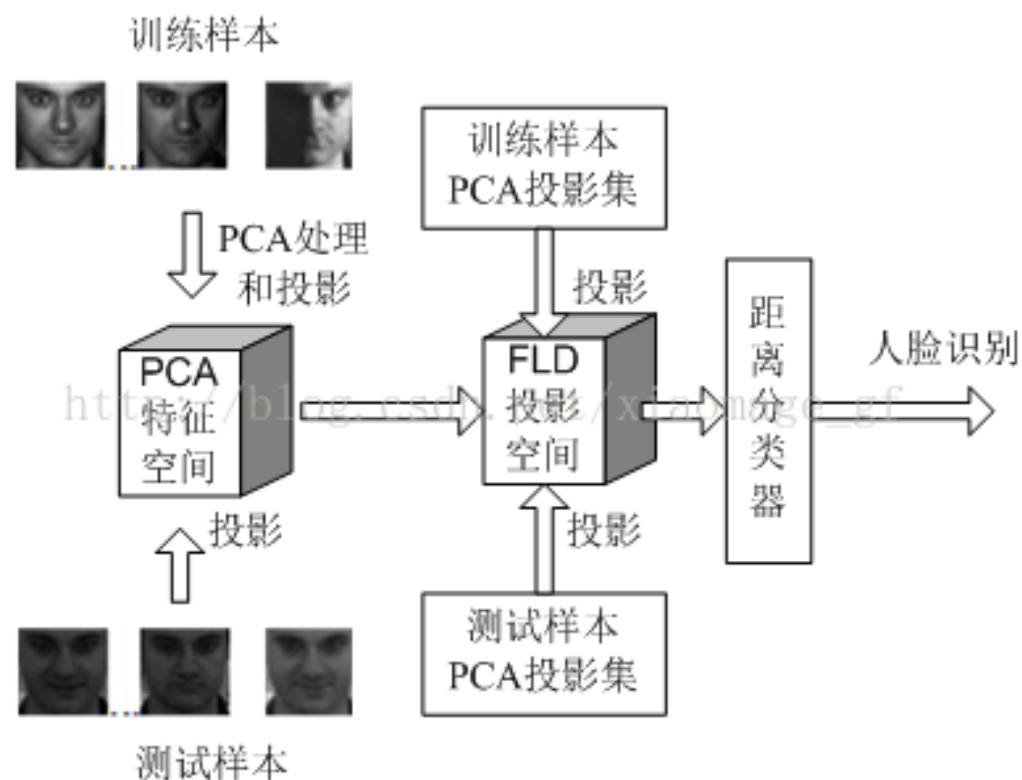
$$\begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_N \end{bmatrix} \dashrightarrow PCA \dashrightarrow \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_K \end{bmatrix}$$

- 2) LDA 寻找最具有判别性的方向

$$\begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_K \end{bmatrix} \dashrightarrow LDA \dashrightarrow \begin{bmatrix} z_1 \\ z_2 \\ \dots \\ z_{C-1} \end{bmatrix}$$

LDA在人脸识别中的应用

具体的LDA人脸识别流程图如下：



LDA算 法的主 要优点

在降维过程中可以使用类别的先验知识经验，而像**PCA**这样的无监督学习则无法使用类别先验知识。

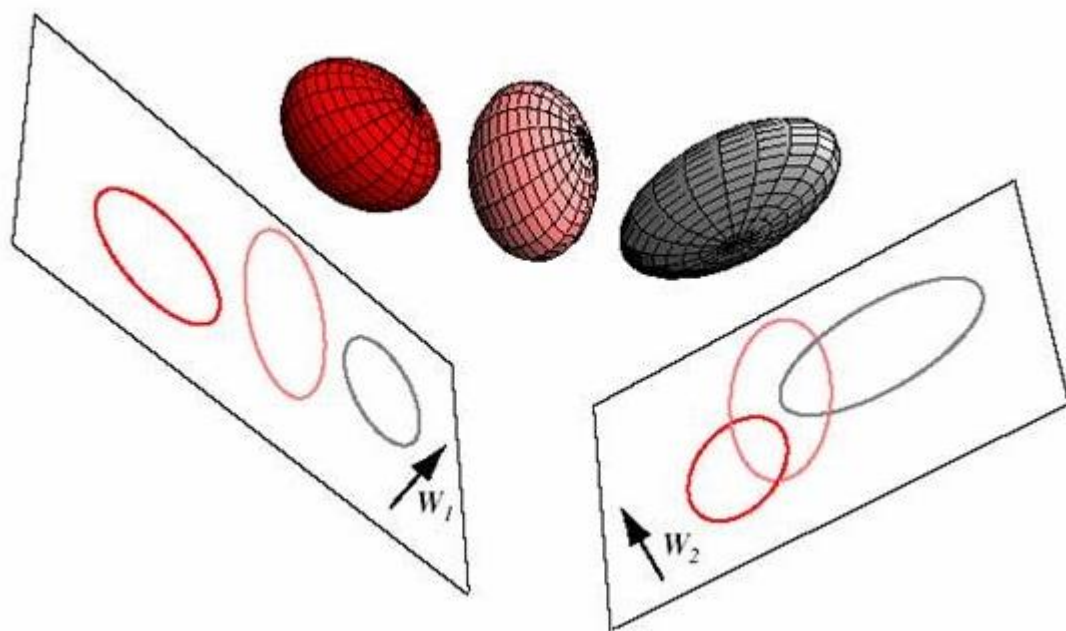
LDA算法的主要缺点

LDA不适合对非高斯分布样本进行降维，**PCA**也有这个问题。

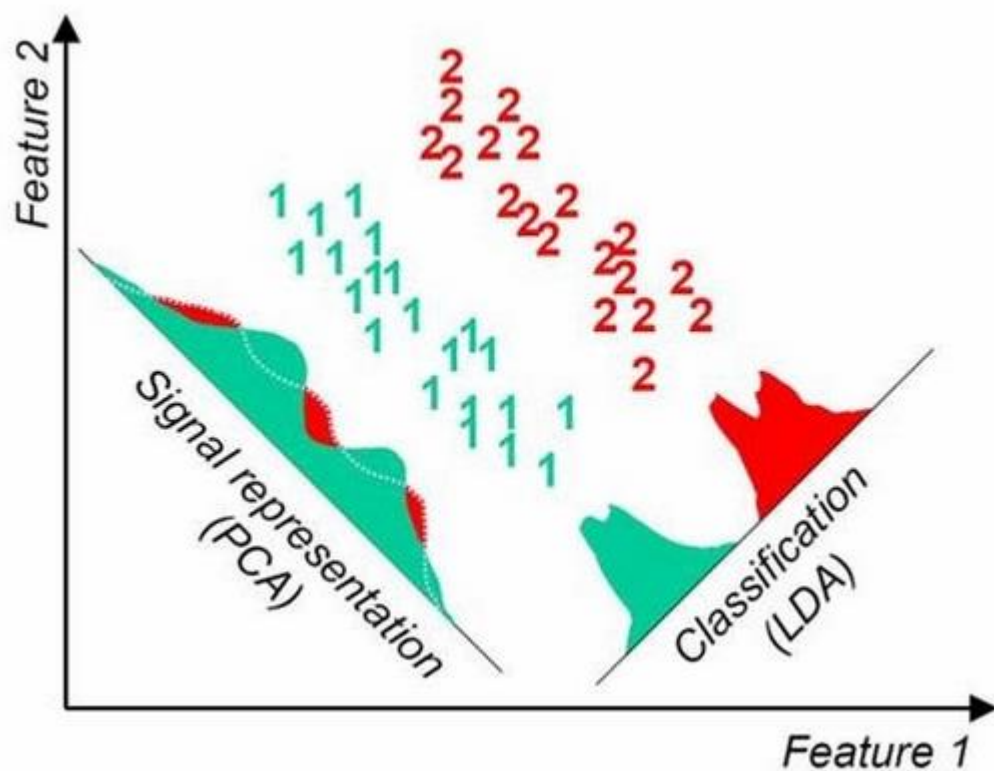
LDA在样本分类信息依赖方差而不是均值的时候，降维效果不好。

栗子

将3维空间上的球体样本点投影到二维上， W_1 相比 W_2 能够获得更好的分离效果。



PCA与LDA的降维对比:



PCA选择样本点投影具有最大方差的方向，
LDA选择分类性能最好的方向。