

特征降维

翟婷婷

扬州大学
信息工程（人工智能）学院
zh tt@yzu.edu.cn

2023春

课程目标

- 理解特征降维的意义。
- 掌握主成分分析(PCA)和线性判别分析(FDA)两种特征降维方法的原理，并能实际运用。
- 理解PCA与FDA降维的差别。

引言

- 一般来说，如果用给定的一组特征获得的分类性能不够好，那么考虑添加新的特征是很自然的。
- 在实际应用中，不是特征越多分类精度越高。当特征维数增加到一个临界点后，继续增加反而会导致分类器性能变差，这种现象称为“**维数灾难**”问题。
- 缓解“维数灾难”经常使用的**办法**：
降低特征向量的维数，即**特征降维**

特征降维

➤ 降维的可行性：特征向量往往包含冗余信息：

① 有些特征可能与分类问题无关

② 特征之间存在着很强的相关性

➤ 降维方法：

重新设计特征：设计生成新的更少的有效特征。

特征选择：选择现有特征集的一个子集。

特征变换：对现有特征组合，形成新的更少的特征。

➤ 本节集中于通过特征变换进行降维的方法。

特征降维

- 两个经典的通过线性变换进行降维的方法：
 - ✓ 主成分分析 (Principal Component Analysis, PCA):
最小化重构误差 (误差平方和)
 - ✓ 线性判别分析 (Linear Discriminant Analysis, LDA):
最大化类别可分性
- 两个方法目的都是寻找一个好的线性变换/投影：
 - ✓ PCA: 寻找最能够表示(**represents**)原始数据的投影方法。——不考虑类别
 - ✓ LDA: 寻找最能够分开(**separates**)各类数据的投影方法。——考虑类别

主成分分析(PCA)

- PCA: 寻找数据的线性投影, 将高维特征向量投影到低维空间中, 是无监督的降维方法。
- 首先考虑问题1: 给定 d 维的样本集 $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, 要求采用一个向量 \mathbf{x} 来最好地表示该样本集, 那么, \mathbf{x} 是什么呢?

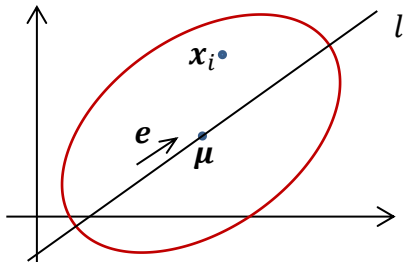
需要一个准则函数来评价表示的好坏程度, 经常使用误差平方和准则函数:

$$J(\mathbf{x}) = \sum_{i=1}^n \|\mathbf{x} - \mathbf{x}_i\|^2$$

使得准则函数 $J(\mathbf{x})$ 取得最小值的那个 \mathbf{x} , 就是我们要找的那个向量。令 $J(\mathbf{x})$ 关于 \mathbf{x} 的梯度为零, 很容易得到, 最优的向量 $\mathbf{x}^* = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$, 即为样本均值向量 $\boldsymbol{\mu}$ 。

主成分分析(PCA)

- **考虑问题2:** 给定 d 维的样本集 $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, 记 $\boldsymbol{\mu}$ 为样本集的均值向量, 记 l 为通过 $\boldsymbol{\mu}$ 的一条直线, \mathbf{e} 为直线的**方向向量**, 且 $\|\mathbf{e}\| = 1$ 。将 \mathcal{D} 中的样本投影到直线 l 上, 如何投影能使得投影点最好地表示原始样本集?



- 记样本 \mathbf{x}_i 投影到直线 l 上的点记为 $\hat{\mathbf{x}}_i$, 则 $\hat{\mathbf{x}}_i$ 可以表示为 $\boldsymbol{\mu} + a_i \mathbf{e}$, 其中 a_i 是一个实数, 表示离开 $\boldsymbol{\mu}$ 的距离。

主成分分析(PCA)

- 误差平方和可表示为:

$$\begin{aligned} J(a_1, a_2 \cdots a_n) &= \sum_{i=1}^n \|\hat{\mathbf{x}}_i - \mathbf{x}_i\|^2 \\ &= \sum_{i=1}^n \|\boldsymbol{\mu} + a_i \mathbf{e} - \mathbf{x}_i\|^2 = \sum_{i=1}^n \|a_i \mathbf{e} - (\mathbf{x}_i - \boldsymbol{\mu})\|^2 \\ &= \sum_{i=1}^n a_i^2 - 2 \sum_{i=1}^n a_i \mathbf{e}^\top (\mathbf{x}_i - \boldsymbol{\mu}) + \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 \end{aligned}$$

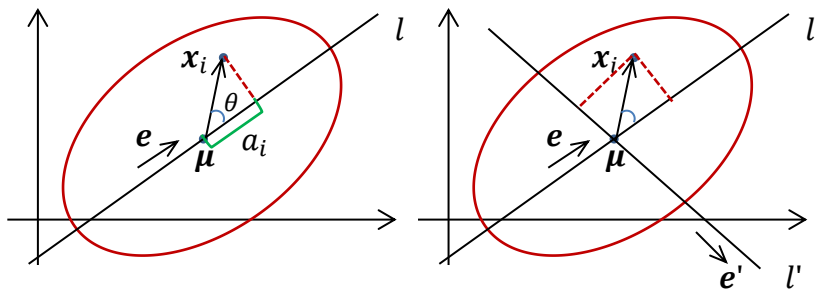
- 求出最优的系数 a_i :

$$\frac{\partial J(a_1, a_2 \cdots a_n)}{\partial a_i} = 2a_i - 2\mathbf{e}^\top (\mathbf{x}_i - \boldsymbol{\mu}) = 0$$

$$\Rightarrow a_i = \mathbf{e}^\top (\mathbf{x}_i - \boldsymbol{\mu}) = \|\mathbf{x}_i - \boldsymbol{\mu}\| \cos \theta \quad (\|\mathbf{e}\| = 1)$$

- 结果表明: 对每个样本垂直投影到直线 l 上, 会产生最小的误差平方和。

主成分分析(PCA)



➤ 考虑问题3:

如何找到一个最优的投影方向 e ，使得 \mathcal{D} 中的样本投影到该方向上的投影点能最好地表示原始样本集？

主成分分析(PCA)

➤ 把 $a_i = \mathbf{e}^\top (\mathbf{x}_i - \boldsymbol{\mu})$ 带入误差平方和准则函数:

$$\begin{aligned} J(\mathbf{e}) &= \sum_{i=1}^n a_i^2 - 2 \sum_{i=1}^n a_i^2 + \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 \\ &= - \sum_{i=1}^n \mathbf{e}^\top (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^\top \mathbf{e} + \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 \\ &= - \mathbf{e}^\top \mathbf{S} \mathbf{e} + \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 \end{aligned}$$

其中,

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

\mathbf{S} 称为散布矩阵(Scatter Matrix), 它是样本协方差矩阵的 $n-1$ 倍。 \mathbf{S} 是对称半正定矩阵。

主成分分析(PCA)

- 要找到满足如下条件的 \mathbf{e} :

$$\min_{\mathbf{e}} J(\mathbf{e}) \text{ 满足 } \|\mathbf{e}\| = 1$$

$$\text{等价于: } \min_{\mathbf{e}} -\mathbf{e}^T \mathbf{S} \mathbf{e} \text{ 满足 } \mathbf{e}^T \mathbf{e} = 1$$

- 应用拉格朗日乘子法:

$$f(\mathbf{e}, \lambda) = -\mathbf{e}^T \mathbf{S} \mathbf{e} + \lambda(\mathbf{e}^T \mathbf{e} - 1)$$

令对 \mathbf{e} 求偏导为零得到:

$$\frac{\partial f(\mathbf{e})}{\partial \mathbf{e}} = -2\mathbf{S}\mathbf{e} + 2\lambda\mathbf{e} = 0$$

$$\Rightarrow \mathbf{S}\mathbf{e} = \lambda\mathbf{e} \Rightarrow \mathbf{e}^T \mathbf{S} \mathbf{e} = \lambda \mathbf{e}^T \mathbf{e} = \lambda$$

λ 为 \mathbf{S} 的特征值, \mathbf{e} 为 \mathbf{S} 对应于特征值 λ 的特征向量

- 为了最小化误差平方和函数 $J(\mathbf{e})$, 选取散布矩阵 \mathbf{S} 最大的特征值对应的特征向量 \mathbf{e} 作为投影直线的方向。

主成分分析(PCA)

- 通过选取散布矩阵 S 最大的特征值对应的特征向量 \mathbf{e} 作为投影方向，可以将一个 d 维的训练样本 \mathbf{x}_i 投影到一维的空间中，它的坐标值为：

$$a_i = \mathbf{e}^T(\mathbf{x}_i - \boldsymbol{\mu})$$

- 考虑更一般的情况：投影到 d' 维的空间中， $d' < d$ 。设 $\mathbf{e}^1, \mathbf{e}^2 \dots \mathbf{e}^{d'}$ 分别为以 $\boldsymbol{\mu}$ 为原点的 d' 个方向向量， $\mathbf{e}^i \perp \mathbf{e}^j, \forall i \neq j, \|\mathbf{e}^i\| = 1$ ，则一个样本 \mathbf{x}_i 投影到 d' 维子空间中的点 $\hat{\mathbf{x}}_i$ 可以表示为：

$$\hat{\mathbf{x}}_i = \boldsymbol{\mu} + \sum_{j=1}^{d'} a_i^j \mathbf{e}^j$$

带入误差平方和函数，先求出投影后在每个方向上的坐标值 a_i^j ，然后再求出最优的投影方向 $\mathbf{e}^1, \mathbf{e}^2 \dots \mathbf{e}^{d'}$ 。

主成分分析(PCA)

➤ 误差平方和函数为:

$$\begin{aligned} J &= \sum_{i=1}^n \|\hat{\mathbf{x}}_i - \mathbf{x}_i\|^2 = \sum_{i=1}^n \left\| \sum_{j=1}^{d'} a_i^j \mathbf{e}^j - (\mathbf{x}_i - \boldsymbol{\mu}) \right\|^2 \\ &= \sum_{i=1}^n \left(\left\| \sum_{j=1}^{d'} a_i^j \mathbf{e}^j \right\|^2 - 2 \left(\sum_{j=1}^{d'} a_i^j \mathbf{e}^j \right)^\top (\mathbf{x}_i - \boldsymbol{\mu}) + \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 \right) \\ &= \sum_{i=1}^n \left(\sum_{j=1}^{d'} (a_i^j)^2 - 2 \sum_{j=1}^{d'} a_i^j \mathbf{e}^{j\top} (\mathbf{x}_i - \boldsymbol{\mu}) + \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 \right) \end{aligned}$$

其中:

$$\left\| \sum_{j=1}^{d'} a_i^j \mathbf{e}^j \right\|^2 = \left(\sum_{j=1}^{d'} a_i^j \mathbf{e}^j \right)^\top \left(\sum_{j=1}^{d'} a_i^j \mathbf{e}^j \right) = \sum_{j=1}^{d'} (a_i^j)^2$$

➤ 令准则函数对 a_i^j 求偏导数为零:

$$\frac{\partial J}{\partial a_i^j} = 2a_i^j - 2\mathbf{e}^{j\top} (\mathbf{x}_i - \boldsymbol{\mu}) = 0$$

$$a_i^j = \mathbf{e}^{j\top} (\mathbf{x}_i - \boldsymbol{\mu})$$

主成分分析(PCA)

➤ 代回误差平方和函数:

$$\begin{aligned} J &= \sum_{i=1}^n \left(\sum_{j=1}^{d'} (a_i^j)^2 - 2 \left(\sum_{j=1}^{d'} a_i^j \mathbf{e}^j \right)^\top (\mathbf{x}_i - \boldsymbol{\mu}) + \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 \right) \\ &= \sum_{i=1}^n \left(- \sum_{j=1}^{d'} (a_i^j)^2 + \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 \right) \\ &= - \sum_{i=1}^n \sum_{j=1}^{d'} \mathbf{e}^{j\top} (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^\top \mathbf{e}^j + \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 \\ &= - \sum_{j=1}^{d'} \mathbf{e}^{j\top} \mathbf{S} \mathbf{e}^j + \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 \quad (\mathbf{S} \text{ 为散布矩阵}) \end{aligned}$$

➤ 也就是说, 我们要找的最优投影方向 $\mathbf{e}^1, \mathbf{e}^2 \dots \mathbf{e}^{d'}$ 为:

$$\min_{\mathbf{e}^1, \mathbf{e}^2 \dots \mathbf{e}^{d'}} - \sum_{j=1}^{d'} \mathbf{e}^{j\top} \mathbf{S} \mathbf{e}^j \quad \text{满足} \mathbf{e}^{j\top} \mathbf{e}^j = 1, \forall j = 1, 2, \dots, d'$$

➤ 同样使用拉格朗日乘子法求解:

主成分分析(PCA)

$$f(\mathbf{e}^1, \mathbf{e}^2 \dots \mathbf{e}^{d'}, \lambda^1 \dots \lambda^{d'}) = - \sum_{j=1}^{d'} \mathbf{e}^{j\top} \mathbf{S} \mathbf{e}^j + \sum_{j=1}^{d'} \lambda^j (\mathbf{e}^{j\top} \mathbf{e}^j - 1)$$

➤ 令 f 对 \mathbf{e}^j 求偏导数为零:

$$\frac{\partial f}{\partial \mathbf{e}^j} = -2\mathbf{S}\mathbf{e}^j + 2\lambda^j \mathbf{e}^j = 0$$

$$\Rightarrow \mathbf{S}\mathbf{e}^j = \lambda^j \mathbf{e}^j \Rightarrow \mathbf{e}^{j\top} \mathbf{S} \mathbf{e}^j = \lambda^j$$

代回误差平方和函数得到:

$$J = - \sum_{j=1}^{d'} \lambda^j + \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}\|^2$$

因此, 为了使得误差平方和函数取得最小值, 选择的 d' 个投影方向应该是散布矩阵 \mathbf{S} 的前 d' 个最大的特征值对应的 d' 个特征向量。由于散布矩阵 \mathbf{S} 是实对称矩阵(d 维), 一定可以找到 d' 个正交的特征向量。

主成分分析(PCA)

- 求出的 $\mathbf{e}^1, \mathbf{e}^2 \dots \mathbf{e}^{d'}$ 称为**主成分**(Principal component), 它们构成一个**新的坐标系**。一个 **d 维**的向量 \mathbf{x}_i 在该坐标系中新的表示是一个 **d' 维**的向量 $(a_i^1, a_i^2, \dots a_i^{d'})^\top$ 其中 $a_i^j = \mathbf{e}^{j^\top}(\mathbf{x}_i - \boldsymbol{\mu})$ 是向量 $\mathbf{x}_i - \boldsymbol{\mu}$ 在 \mathbf{e}^j 上投影的长度。
- 重新表示 \mathbf{x}_i 降维后的向量为:

$$\begin{aligned} \begin{bmatrix} a_i^1 \\ a_i^2 \\ \vdots \\ a_i^{d'} \end{bmatrix} &= \begin{bmatrix} \mathbf{e}^{1^\top}(\mathbf{x}_i - \boldsymbol{\mu}) \\ \mathbf{e}^{2^\top}(\mathbf{x}_i - \boldsymbol{\mu}) \\ \vdots \\ \mathbf{e}^{d'^\top}(\mathbf{x}_i - \boldsymbol{\mu}) \end{bmatrix} = \begin{bmatrix} \mathbf{e}^{1^\top} \\ \mathbf{e}^{2^\top} \\ \vdots \\ \mathbf{e}^{d'^\top} \end{bmatrix} (\mathbf{x}_i - \boldsymbol{\mu}) \\ &= [\mathbf{e}^1 \ \mathbf{e}^2 \ \dots \ \mathbf{e}^{d'}]^\top (\mathbf{x}_i - \boldsymbol{\mu}) = \mathbf{Q}^\top (\mathbf{x}_i - \boldsymbol{\mu}) \end{aligned}$$

其中: $\mathbf{Q} = [\mathbf{e}^1, \mathbf{e}^2 \dots \mathbf{e}^{d'}]$ 称为**投影矩阵**。对任意一个 \mathbf{x}_i , 通过 $\mathbf{Q}^\top (\mathbf{x}_i - \boldsymbol{\mu})$ 可以将 \mathbf{x}_i 从 d 维降到 d' 维, 其中 $d' < d$ 。

主成分分析(PCA)

➤ 主成分分析步骤 (d 维降为 d' 维) :

1、计算散布矩阵 S (或样本的协方差矩阵)

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

2、计算 S 的特征值和特征向量;

3、按照特征值从大到小对相应的特征向量排序;

4、选择特征值最大的前 d' 个特征向量作为投影向量, 构成 $d * d'$ 维的**标准正交**的投影矩阵 Q ;

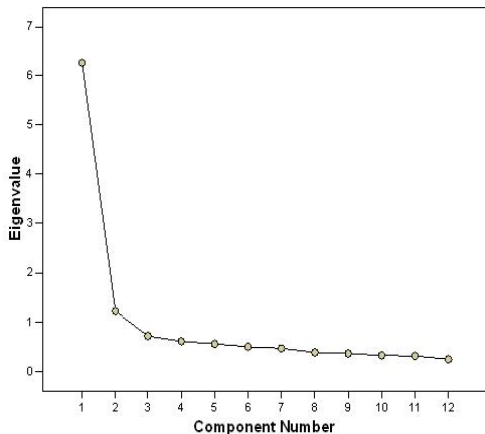
5、对任意的 d 维样本 \mathbf{x}_i , 用PCA降维后的 d' 维向量为:

$$\mathbf{y}_i = Q^\top (\mathbf{x}_i - \boldsymbol{\mu})$$

$$\text{重构公式: } \mathbf{x}_i \approx Q\mathbf{y}_i + \boldsymbol{\mu}$$

主成分分析(PCA)

- 通常，最大的几个特征值占据了所有特征值之和的绝大部分。将数据投影到少数几个最大特征值对应的特征向量方向上即可保留原数据中的绝大部分信息。



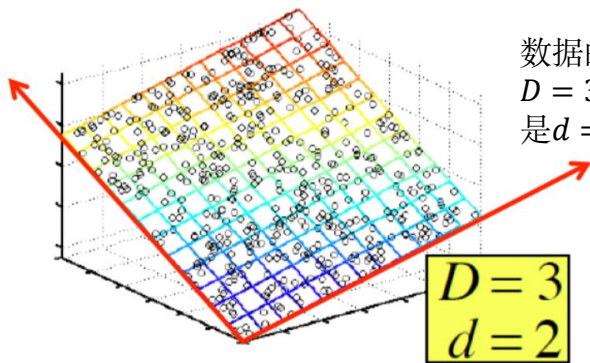
剩下的小部分信息，即较小的特征值对应的特征向量所表示的信息，通常可以认为是数据噪声而丢掉。

主成分分析(PCA)

- 原始数据的**本征维度**(intrinsic dimensionality)决定了能否在 d' 维($d' < d$)的子空间中充分表示给定的 d 维数据。
- 数据集的**本征维度**可以看作是表示数据集所需要的**最小变量数**。因为任意低维数据可简单地通过增加(如复制)不必要的维度或随机维将其转换至更高维空间中，所以数据集的本征维度小于等于数据集的原始维度。
- 给定一个 d 维的数据集，若数据集的**本征维度为 d' 维**， $d' < d$ ，可以将该数据集从 d 维削减至 d' 维，不会丢失重要信息。

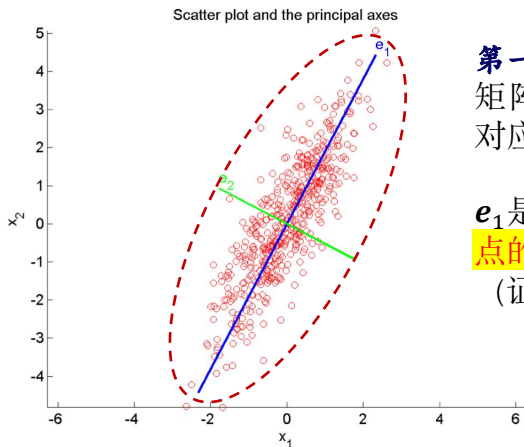
主成分分析(PCA)

- 本征维度的几何解释是整个数据集中的样本位于 d' 维超曲面拓扑上。



主成分分析(PCA)

- PCA的几何解释：样本在 d 维空间中形成一个椭圆形的云团，散布矩阵 S 的特征向量就是这个云团的主轴。



第一主成分 e_1 是散布矩阵 S 的最大特征值对应的特征向量。

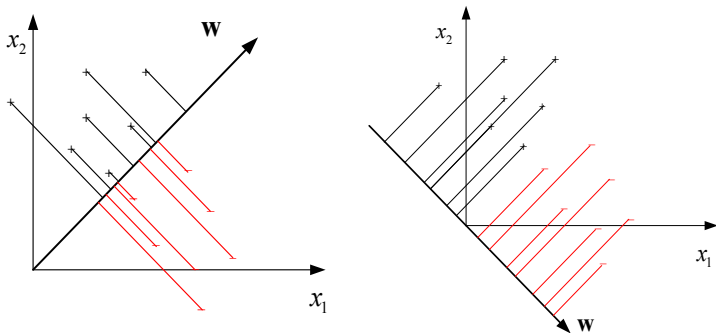
e_1 是所有样本的**投影点的方差最大的方向**。
(证明-课后作业)

Fisher线性判别分析(FDA)

- PCA寻找的是用来有效表示数据的主轴方向(主成分), 是**无监督**的特征降维方法, 没有考虑样本的类别信息。
- Fisher线性判别分析(FDA)寻找的是**能够有效分类的方向**, 是**有监督**的降维方法。

Fisher线性判别分析(FDA)

- 对于两类别的样本集合，FDA希望找到一个投影方向，使得不同类别的样本在该方向上的投影尽量分开。



对于两类别的样本集合，向不同的方向投影，产生的数据点的可分程度是不同的。

Fisher线性判别分析(FDA)

- 给定一个样本集合 $\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2\}$: \mathcal{D}_1 中的样本属于 ω_1 , 样本数为 n_1 ; \mathcal{D}_2 中的样本属于 ω_2 , 样本数为 n_2 。
- FDA的目标: 找到一个直线方向 \mathbf{w} , **将两类样本最好地分开**, 其中 $\|\mathbf{w}\| = 1$ 。样本 \mathbf{x} 投影后的点的坐标为 $y = \mathbf{w}^\top \mathbf{x}$ 。

FDA定义一个**准则函数**衡量两类样本分开的程度:

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

—— 不同类的投影点尽量分开
—— 同类的投影点尽量靠近

其中, \tilde{m}_i 是 \mathcal{D}_i 中样本**投影后的点的均值**;

\tilde{s}_i^2 是 \mathcal{D}_i 中样本**投影后的点的分散程度**, 即方差。

Fisher线性判别分析(FDA)

➤ \tilde{m}_i 是 \mathcal{D}_i 中样本投影后的点的均值:

$$\tilde{m}_1 = \frac{1}{n_1} \sum_{x_i \in \mathcal{D}_1} \mathbf{w}^\top x_i = \mathbf{w}^\top \mathbf{m}_1, \text{ 其中 } \mathbf{m}_1 = \frac{1}{n_1} \sum_{x_i \in \mathcal{D}_1} x_i$$

$$\tilde{m}_2 = \frac{1}{n_2} \sum_{x_i \in \mathcal{D}_2} \mathbf{w}^\top x_i = \mathbf{w}^\top \mathbf{m}_2, \text{ 其中 } \mathbf{m}_2 = \frac{1}{n_2} \sum_{x_i \in \mathcal{D}_2} x_i$$

$$(\tilde{m}_1 - \tilde{m}_2)^2 = \mathbf{w}^\top (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^\top \mathbf{w} = \mathbf{w}^\top S_B \mathbf{w}$$

其中,

$$S_B = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^\top$$

称 S_B 为类间散布矩阵 (**between-class scatter matrix**) , 它是对称半正定的。

Fisher线性判别分析(FDA)

➤ \tilde{s}_i^2 是 \mathcal{D}_i 中样本投影后的点的方差:

$$\begin{aligned}\tilde{s}_1^2 &= \frac{1}{n_1} \sum_{x_i \in \mathcal{D}_1} (\mathbf{w}^\top \mathbf{x}_i - \tilde{m}_1)^2 = \frac{1}{n_1} \sum_{x_i \in \mathcal{D}_1} (\mathbf{w}^\top \mathbf{x}_i - \mathbf{w}^\top \mathbf{m}_1)^2 \\ &= \frac{1}{n_1} \sum_{x_i \in \mathcal{D}_1} \mathbf{w}^\top (\mathbf{x}_i - \mathbf{m}_1) (\mathbf{x}_i - \mathbf{m}_1)^\top \mathbf{w} = \mathbf{w}^\top S_1 \mathbf{w}\end{aligned}$$

其中, $S_1 = \frac{1}{n_1} \sum_{x_i \in \mathcal{D}_1} (\mathbf{x}_i - \mathbf{m}_1) (\mathbf{x}_i - \mathbf{m}_1)^\top$

称为 ω_1 类的类内散布矩阵。

➤ 同理可得: $\tilde{s}_2^2 = \mathbf{w}^\top S_2 \mathbf{w}$

其中, $S_2 = \frac{1}{n_2} \sum_{x_i \in \mathcal{D}_2} (\mathbf{x}_i - \mathbf{m}_2) (\mathbf{x}_i - \mathbf{m}_2)^\top$, 称为 ω_2 类的类内散布矩阵。

Fisher线性判别分析(FDA)

- 定义总的类内散布矩阵(within-class scatter matrix):

$$S_w = S_1 + S_2$$

则可以得到:

$$\tilde{s}_1^2 + \tilde{s}_2^2 = \mathbf{w}^T S_1 \mathbf{w} + \mathbf{w}^T S_2 \mathbf{w} = \mathbf{w}^T S_w \mathbf{w}$$

- 代入准则函数得到:

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

广义瑞利熵
generalized
Rayleigh quotient

FDA的目标是找到一个最大化 $J(\mathbf{w})$ 的 \mathbf{w} 。

根据瑞利商函数的性质, $J(\mathbf{w})$ 的最大值为 $S_w^{-1} S_B$ 的最大特征值, 最小值为 $S_w^{-1} S_B$ 的最小特征值。

Fisher线性判别分析(FDA)

- 问题可以重写为

$$\min_{\mathbf{w}} -\mathbf{w}^T S_B \mathbf{w} \quad s.t. \mathbf{w}^T S_w \mathbf{w} = 1$$

- 拉格朗日乘子法: $f(\mathbf{w}, \lambda) = -\mathbf{w}^T S_B \mathbf{w} + \lambda(\mathbf{w}^T S_w \mathbf{w} - 1)$

- 对 \mathbf{w} 求偏导, 令梯度为0:

$$\frac{\partial f(\mathbf{w}, \lambda)}{\partial \mathbf{w}} = -2S_B \mathbf{w} + 2\lambda S_w \mathbf{w} = 0$$

- 解得: $S_B \mathbf{w} = \lambda S_w \mathbf{w}$ (S_B 的秩最多为1, S_w 在样本数大于维数时经常是非奇异矩阵)

$$\Rightarrow S_w^{-1} S_B \mathbf{w} = \lambda \mathbf{w}$$

$$\Rightarrow \mathbf{w} = \frac{1}{\lambda} S_w^{-1} S_B \mathbf{w} = \frac{1}{\lambda} S_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) \overset{\text{标量}}{\boxed{(\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}}}$$

$$\Rightarrow \mathbf{w} = \alpha S_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$

Fisher线性判别分析(FDA)

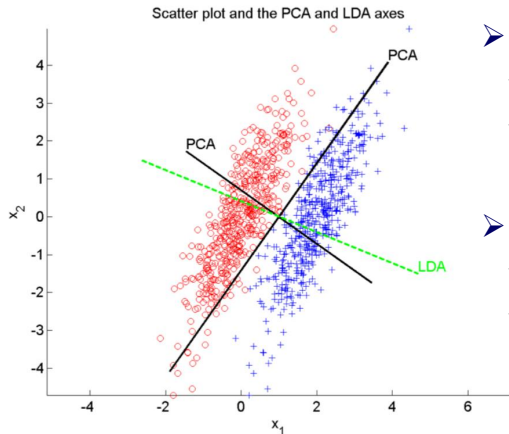
- 由于 \mathbf{w} 的模对问题本身无关紧要，因此使得准则函数最大化的 \mathbf{w} 取值为：

$$\mathbf{w} = S_w^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$$

- 一个 d 维的样本 \mathbf{x}_i 投影到方向 \mathbf{w} 上，成为一个1维的点： $y_i = \mathbf{w}^\top \mathbf{x}_i$ 。

Fisher线性判别分析(FDA)

PCA vs. FDA:



➤ 主成分分析(PCA):

- ✓ 无监督
- ✓ 寻找用来**有效表示数据**的投影

➤ 线性判别分析(FDA):

- ✓ 有监督
- ✓ 寻找用来**有效分类**的投影

多重FDA (课下了解)

- 如何将FDA拓展到多类别的数据集上?

多重FDA:

两类问题需要一个Fisher投影方向

c 类问题需要 $c-1$ 个Fisher投影方向

- 给定一个样本集 $\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2 \cdots \mathcal{D}_c\}$: \mathcal{D}_i 中的样本属于 ω_i 类, 样本数为 n_i , \mathcal{D} 中的总样本数为 $n = \sum_i n_i$.
- 多重FDA的目的是找到 $c-1$ 个投影方向 $\mathbf{w}_1, \mathbf{w}_2 \cdots \mathbf{w}_{c-1}$, 使得不同类别的样本投影后的点尽可能分开。
记 $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2 \cdots \mathbf{w}_{c-1}]$, \mathbf{W} 是 $d * (c-1)$ 维投影矩阵。

多重FDA (课下了解)

- 对 ω_i 类的类内散布矩阵进行拓展:

$$S_i = \sum_{\mathbf{x} \in \mathcal{D}_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^\top, \text{ 其中 } \mathbf{m}_i = \frac{1}{n_i} \sum_{\mathbf{x} \in \mathcal{D}_i} \mathbf{x}$$

总的类内散布矩阵为:

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

- 定义类间散布矩阵为:

$$S_B = \sum_{i=1}^c n_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^\top$$

其中, $\mathbf{m} = \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{D}} \mathbf{x}$ 为总的均值向量。

多重FDA (课下了解)

- 准则函数变为:

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

其中 $|\cdot|$ 表示行列式。

- 可以得到: 使得 $J(W)$ 最大化的 W 的列向量是矩阵 $S_w^{-1}S_B$ 的最大的特征值对应的特征向量。
- 一旦使用PCA或LDA完成从 d 维原始特征空间到低维子空间的转换, 就可以使用降维后的数据集来训练分类器。

小结

- 降低特征向量的维数，即特征降维，是缓解“维数灾难”经常的办法。
- 主成分分析 (PCA)和线性判别分析(LDA)是两种通过特征变换进行降维的方法。
- 两个方法目的都是寻找一个好的线性变换/投影：
 - ✓ PCA：寻找最能够表示(represents) 原始数据的投影方向。——不考虑类别
 - ✓ LDA：寻找最能够分开(separates)各类数据的投影方向。——考虑类别
- 使用PCA和LDA进行特征降维的步骤。

重点