



PCA 与 SVD 笔记

李向阳 d1142845997@gmail.com

目录	2
----	---

目录

1 引入	3
2 SVD	3
3 PCA	3
3.1 统计中的 PCA	3
3.2 利用 SVD 做 PCA	4
4 Kernel PCA	6
4.1 PCA 重述	6
4.2 核方法	7
5 总结	7
5.1 参考资料	7

1 引入

PCA 与 SVD 是一种降维技术, 在很多领域都有重要应用, 因此五花八门的介绍也很多. 为了建立自己的体系, 这里稍微介绍一下, 也方便回顾.

SVD(Singular Value Decomposition), 也就是奇异值分解, 严格的数学介绍是在大四的矩阵分析课上, 而 PCA(Principal Component Analysis), 即主成分分析, 是在大三的数据分析课上学的, 是偏统计的角度, 对细节也都做了证明. 因此, 这篇笔记会略去一些证明, 着重对关系进行梳理.

2 SVD

SVD 本身就具有丰富的内容, 详细可见维基 https://en.wikipedia.org/wiki/Singular_value_decomposition.

3 PCA

3.1 统计中的 PCA

统计中经常用 n 表示样本数, 用 p 表示特征数 (维数), 这里我们仍延续机器学习系列笔记的记号, 用 m 表示样本数, 用 n 表示特征数.

设总体为 n 维随机向量 $\mathbf{X} = (X_1, X_2, \dots, X_n)^T$, 主成分分析是构造原变量的一系列线性组合 Y_1, Y_2, \dots , 使其方差达到最大 (具体可回顾数据分析方法课本). 实际中, 是从总体中抽样得到了 m 个样本, 其中 $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{in})^T, \mathbf{x}_i \in \mathbb{R}^n, i = 1, 2, \dots, m$, 把样本按行并起来得到了 $m \times n$ 的观测数据矩阵

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_m^T \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \cdots & x_{mn} \end{pmatrix}$$

为了后面讨论的方便, 我们要求 \mathbf{X} 是去中心化后的矩阵 (即每一列已经减去了列均值), 即 $\sum_{i=1}^m x_{i,j} = 0, j = 1, 2, \dots, n$, 事实上可写为 $\sum_{i=1}^m \mathbf{x}_i = \mathbf{0}$, 这样样本的协方差阵为 (为了推导形式的方便, 分母取为 m 也可, 本身二者都可以)

$$\mathbf{S} = \frac{1}{m-1} \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T = \frac{1}{m-1} \mathbf{X}^T \mathbf{X}$$

数据分析课上讲过, 设 \mathbf{S} 为样本的协方差矩阵 (半正定矩阵), 其特征值按大小顺序排列为 $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n \geq 0$, 相应的正交单位化特征向量为 $\mathbf{e}_1, \cdots, \mathbf{e}_n$, 则总体的第 j 个主成分可表示为

$$y_j = \mathbf{e}_j^T \mathbf{x} = e_{j1}x_1 + e_{j2}x_2 + \cdots + e_{jn}x_n, j = 1, 2, \cdots, n$$

其中 $\mathbf{e}_j = (e_{j1}, e_{j2}, \cdots, e_{jn})^T$. 把第 i 个样本值 $\mathbf{x}_i = (x_{i1}, x_{i2}, \cdots, x_{in})^T$ 代入, 便得到该样本在第 j 个主成分上的得分 y_{ij} , 如下表所示

表 1: 原始数据及主成分得分

样本序号	原始变量				主成分			
	X_1	X_2	\cdots	X_n	Y_1	Y_2	\cdots	Y_n
1	x_{11}	x_{12}	\cdots	x_{1n}	y_{11}	y_{12}	\cdots	y_{1n}
2	x_{21}	x_{22}	\cdots	x_{2n}	y_{21}	y_{22}	\cdots	y_{2n}
\vdots	\vdots	\vdots		\vdots	\vdots	\vdots		\vdots
m	x_{m1}	x_{m2}	\cdots	x_{mn}	y_{m1}	y_{m2}	\cdots	y_{mn}

一般我们会选取前 $k(k < n)$ 个主成分, 用前 k 个主成分的得分替代原始数据做分析, 这样便达到了数据降维的目的. 其中第 j 个主成分的贡献率为 $\lambda_j / \sum_{i=1}^n \lambda_i$, 前 k 个主成分的累计贡献率为 $\sum_{i=1}^k \lambda_i / \sum_{i=1}^n \lambda_i$.

用矩阵形式表达, 即首先对于样本协方差阵 \mathbf{S} 对角化, 有

$$\mathbf{S} = \mathbf{Q} \text{diag}(\lambda_1, \cdots, \lambda_n) \mathbf{Q}^T$$

其中 $\mathbf{Q} = (\mathbf{e}_1, \mathbf{e}_2, \cdots, \mathbf{e}_n)$. 然后可得新的数据矩阵 $\mathbf{Y} = (y_{ij})$ 为

$$\mathbf{Y} = \mathbf{X}\mathbf{Q}$$

若取前 k 个主成分, 即令 $\mathbf{Q}_k = (\mathbf{e}_1, \cdots, \mathbf{e}_k)$, 降维矩阵 \mathbf{Y}_k 为

$$\mathbf{Y}_k = \mathbf{X}\mathbf{Q}_k$$

3.2 利用 SVD 做 PCA

关于 SVD 与 PCA 的关系, 可参考 <http://stats.stackexchange.com/questions/134282/relationship-between-svd-and-pca-how-to-use-svd-to-perform-pca>, 介绍的比较清楚.

这里说一下利用 SVD 做 PCA.

原始数据矩阵 \mathbf{X} 是一个 $m \times n$ 矩阵, 对其做奇异值分解

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$$

其中 $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_m)$ 为 $m \times m$ 阶正交阵, $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n)$ 为 $n \times n$ 阶正交阵, \mathbf{D} 为一个对角阵 (实际上为 $m \times n$), 其对角元为矩阵 \mathbf{X} 的奇异值 σ_i , 也就是 $\mathbf{X}^T \mathbf{X}$ 特征值的平方根. 于是可得

$$\mathbf{X}^T \mathbf{X} = \mathbf{V} \mathbf{D}^2 \mathbf{V}^T$$

这里 $\mathbf{D}^2 = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$ 表示 n 阶对角阵. 注意到样本协方差阵为 $\mathbf{S} = \frac{1}{m-1} \mathbf{X}^T \mathbf{X}$, 因此

$$\mathbf{S} = \mathbf{V} \frac{\mathbf{D}^2}{m-1} \mathbf{V}^T$$

这不正好是协方差阵 \mathbf{S} 的对角化吗? 其中 \mathbf{S} 的特征值为 $\lambda_i = \sigma_i^2 / (m-1)$. 那么根据上面主成分得分的求法, 新的得分矩阵为

$$\mathbf{Y} = \mathbf{X} \mathbf{V} = \mathbf{U} \mathbf{D} \mathbf{V}^T \mathbf{V} = \mathbf{U} \mathbf{D}$$

同样的, 若取前 $k (k < n)$ 个主成分, 即令 \mathbf{V}_k 为矩阵 \mathbf{V} 的前 k 列, \mathbf{U}_k 为矩阵 \mathbf{U} 的前 k 列, \mathbf{D}_{mk} 为矩阵 \mathbf{D} 的前 k 列, \mathbf{D}_k 为矩阵 \mathbf{D} 的左上 $k \times k$ 对角阵, 可得降维得分矩阵为

$$\mathbf{Y}_k = \mathbf{X} \mathbf{V}_k = \mathbf{U} \mathbf{D}_{mk} = \mathbf{U}_k \mathbf{D}_k$$

注意用分块矩阵可以去理解证明上式, $\mathbf{U}_k = (\mathbf{u}_1, \dots, \mathbf{u}_k)$, $\mathbf{V}_k = (\mathbf{v}_1, \dots, \mathbf{v}_k)$.

$$\begin{aligned} \mathbf{Y}_k &= \mathbf{X} \mathbf{V}_k = \mathbf{U} \mathbf{D} \mathbf{V}^T \mathbf{V}_k \\ &= \mathbf{U} \mathbf{D} \begin{pmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \vdots \\ \mathbf{v}_n^T \end{pmatrix} (\mathbf{v}_1, \dots, \mathbf{v}_k) = \mathbf{U} \mathbf{D} \begin{pmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \\ &= \mathbf{U} \mathbf{D}_{mk} = \mathbf{U}_k \mathbf{D}_k \end{aligned}$$

注 3.1. 若令 $\mathbf{X}_k = \mathbf{U}_k \mathbf{D}_k \mathbf{V}_k^T$, 即

$$\mathbf{X}_k = \sum_{j=1}^k \sigma_j \mathbf{u}_j \mathbf{v}_j^T$$

则 \mathbf{X}_k 仍是一个 $m \times n$ 矩阵, 只不过它的秩为 $k < n$, 在矩阵分析的课上我们证明 \mathbf{X}_k 是 \mathbf{X} 的一个最佳低秩逼近, 证明过程也可见 <http://stats.stackexchange.com/questions/130721/what-norm-of-the-reconstruction-error-is-minimized>

\mathbf{X}_k 也可看做是利用前 k 个主成分对矩阵 \mathbf{X} 的重建.

4 Kernel PCA

4.1 PCA 重述

PCA 实际上是将样本点在几个互相正交的方向上进行投影来达到降维的目的. 投影方向就是协方差阵 \mathbf{S} 的特征向量的方向.

设 $\mathbf{S}\mathbf{v} = \lambda\mathbf{v}$, 将特征向量 \mathbf{v} 按照特征值降序排列, 然后将样本点投影在特征向量上. 协方差阵 \mathbf{S} 的特征向量其实可表示为样本点的线性组合. 事实上

$$\lambda\mathbf{v} = \mathbf{S}\mathbf{v} = \frac{1}{m-1} \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T \mathbf{v} = \frac{1}{m-1} \sum_{i=1}^m (\mathbf{x}_i^T \mathbf{v}) \mathbf{x}_i = \frac{1}{m-1} \sum_{i=1}^m \langle \mathbf{x}_i, \mathbf{v} \rangle \mathbf{x}_i$$

其中 $\mathbf{x}_i^T \mathbf{v}$ 是一个数, 因此用了欧式内积表示, 于是可得

$$\mathbf{v} = \frac{1}{\lambda(m-1)} \sum_{i=1}^m \langle \mathbf{x}_i, \mathbf{v} \rangle \mathbf{x}_i = \sum_{i=1}^m \alpha_i \mathbf{x}_i = \mathbf{X}^T \boldsymbol{\alpha}$$

如何寻求系数 $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)^T$ 呢?

再次利用 $\lambda\mathbf{v} = \mathbf{S}\mathbf{v}$, 可得

$$(m-1)\lambda\mathbf{v} = \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T \mathbf{v}, \Rightarrow \mu\mathbf{v} = \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T \mathbf{v}$$

也即

$$\mu \sum_{i=1}^m \alpha_i \mathbf{x}_i = \sum_{i=1}^m \left(\mathbf{x}_i \mathbf{x}_i^T \cdot \sum_{j=1}^m \alpha_j \mathbf{x}_j \right)$$

变形可得

$$\sum_{i=1}^m \mu \alpha_i \cdot \mathbf{x}_i = \sum_{i=1}^m \left(\sum_{j=1}^m \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \right) \cdot \mathbf{x}_i$$

于是, 只需下式成立即可

$$\mu \alpha_i = \sum_{j=1}^m \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle, i = 1, 2, \dots, m$$

因此 $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_m)^T$ 满足 $\mathbf{K}\boldsymbol{\alpha} = \mu\boldsymbol{\alpha}$, 即 $\boldsymbol{\alpha}$ 是矩阵 \mathbf{K} 的特征向量, 其中 $\mathbf{K} \in \mathbb{R}^{m \times m}$ 为样本的内积矩阵

$$\mathbf{K} = \begin{pmatrix} \langle \mathbf{x}_1, \mathbf{x}_1 \rangle & \langle \mathbf{x}_1, \mathbf{x}_2 \rangle & \cdots & \langle \mathbf{x}_1, \mathbf{x}_m \rangle \\ \langle \mathbf{x}_2, \mathbf{x}_1 \rangle & \langle \mathbf{x}_2, \mathbf{x}_2 \rangle & \cdots & \langle \mathbf{x}_2, \mathbf{x}_m \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle \mathbf{x}_m, \mathbf{x}_1 \rangle & \langle \mathbf{x}_m, \mathbf{x}_2 \rangle & \cdots & \langle \mathbf{x}_m, \mathbf{x}_m \rangle \end{pmatrix}$$

当然, 此时的矩阵 \mathbf{K} 还可表示为 $\mathbf{K} = \mathbf{X}\mathbf{X}^T$.

4.2 核方法

考虑先将样本数据映射到另外一个空间, $\phi: \mathcal{X} \rightarrow \mathcal{H}, \mathbf{x} \mapsto \phi(\mathbf{x})$, 然后在新空间中对样本点 $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_m)$ 做 PCA.

假设数据已经中心化, 即 $\sum_{i=1}^m \phi(\mathbf{x}_i) = \mathbf{0}$, 如何将映射后的数据中心化后面再提. 此时样本的协方差为

$$\mathbf{S} = \frac{1}{m} \sum_{i=1}^m \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^T$$

接下来按照 PCA 的步骤是一样的, 只需要把上面的 \mathbf{x}_i 换为 $\phi(\mathbf{x}_i)$ 即可. 关键是计算样本的内积矩阵 \mathbf{K} , 其 (i, j) 元为 $\mathbf{K}(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$.

如同 SVM 一样, 在映射后的空间中计算内积是复杂的, 我们通过构造核函数 $\kappa(\cdot)$ 来计算这个内积 (相关回顾可看 SVM 笔记), 即有

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

核函数的选取跟 SVM 一样, 也有很多种选择, 比如多项式核、高斯核等等.

5 总结

5.1 参考资料

- (1) 博客: <http://www.cnblogs.com/LeftNotEasy/archive/2011/01/19/svd-and-applications.html>, 对奇异值分解的强大应用做了通俗介绍.
- (2) JerryLead 的博客: <http://www.cnblogs.com/jerrylead/archive/2011/04/18/2020216.html>, 对主成分分析做了不同解释.
- (3) PRML: 关于 Kernel PCA 的推导介绍, 可见 12.3 节.

附录