

PART6 非监督式特征提取

2020-05-25

要点:

1. 特征工程；特征降维；特征选择、特征提取
2. 监督式、非监督式特征提取
3. 线性、非线性特征提取
4. PCA的全称？PCA有哪些应用？
5. 理解样本协方差矩阵的本征值与本征列向量的意义；各主成分的分布方差？
6. 掌握利用PCA进行特征提取、特征降维的基本实现过程。
 - (1) 针对特征空间任意观测样本，如何提取指定的主成分？第1主成分？第2主成分？...
 - (2) 如何根据**累积方差解释比**确定主成分数目（即：新的特征空间的特征维数）？

练习:

- 1 给定观测样本集 $D = \{x_i, i=1, \dots, N\}$ ，其中 $x_i \in R^3$. 请结合该样本集，设计一个基于主成分分析的特征降维方法，以便基于该算法，提取原始空间任意观测样本 $x \in R^3$ 的第1、第2主成分.

参考答案:

step1. 基于样本集 D ，估计**样本中心** μ 及**协方差矩阵** Σ .

$$\mu^* = \frac{1}{N} \sum_{i=1}^N x_i \quad \Sigma^* = \frac{1}{N} \sum_{i=1}^N (x_i - \mu^*)(x_i - \mu^*)^T$$

step2. 确定 Σ 的 $p=3$ 个**本征值**及**本征向量**.

得 p 个本征值 $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq 0$

对应本征向量 $a_i, i=1, 2, 3$

step3. 确定 3×2 的变换矩阵 $A_2 = [a_1 \ a_2]$

step4. 对于任意观测 x ，提取该样本的前两个主成分： $[\xi_1 \ \xi_2]^T = A_2^T (x - \mu^*)$

2. 给定数据集 $D = \{x_i, i = 1, \dots, m\}$, 其中 $x_i \in R^d$ 。请结合该样本集D, 设计一个基于PCA的特征降维算法, 以便基于该方法将任意观测 $x \in R^d$ 的降至r维。请详细给出有关步骤和必要表达式。

参考答案:

step1. 基于样本集D, 估计**样本中心** μ 及**协方差矩阵** Σ .

$$\mu^* = \frac{1}{N} \sum_{i=1}^N x_i \quad \Sigma^* = \frac{1}{N} \sum_{i=1}^N (x_i - \mu^*)(x_i - \mu^*)^T$$

step2. 确定 Σ^* 的前 **r ($r < d$)**个最大**本征值**及**本征向量**.

$$\text{得前 } r \text{ 个本征值} \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r$$

$$\text{对应本征向量} \quad a_i, i = 1, \dots, r$$

step3. 确定 $d \times r$ 的变换矩阵 $A_r = [a_1 \ a_2 \ \dots \ a_r]$

step4. 对于任意观测 x , 其新的 r 维空间的映射位置: $\xi_r = A_r^T (x - \mu^*)$

3. 给定数据集 $D = \{x_i, i = 1, \dots, m\}$, 其中 $x_i \in R^d$ 。请结合该样本集D, 设计一个基于主成分分析的特征降维方法, 并且使用新的特征描述样本时, 满足累积方差解释比不低于 $\alpha = 0.9$, 请确定新的特征空间特征维数r, 并将任意观测 $x \in R^d$ 降至r维。请详细给出有关步骤和必要表达式。

4. 给定鸢尾花数据集 $D = \{x_i, i = 1, \dots, m\}$, 其中 $x_i \in R^4$ 。请结合该样本集D, 基于主成分分析法实现上述数据集在新的二维空间可视化。请详细给出有关实现步骤和必要表达式。