

PART6 非监督式特征提取与低维可视化

2023-06

掌握:

1. 理解什么是特征提取? 什么是非监督式特征提取? 什么是线性/非线性特征提取?
2. 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?

3. PCA的全称?
4. 掌握利用PCA进行特征提取或特征降维的基本实现过程。

如何根据累积方差解释比确定特征提取的数目?

如何用PCA实现样本数据的低维可视化?

5. 能使用t-SNE进行数据的低维可视化。

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。

解:

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

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

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

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

对应本征向量 $a_i, i=1, 2, 3$

step3. 确定 3×2 的变换矩阵 $A = [a_1 \ a_2]$

step4. 对于任意观测 x , 提取该样本的前两个主成分: $\xi = A^T (x - \hat{\mu})$

注意:

观测 x 第1主成分: $\xi_1 = a_1^T (x - \hat{\mu})$

观测 x 第2主成分: $\xi_2 = a_2^T (x - \hat{\mu})$

其中: 第1主成分

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, 设计一个基于主成分分析法的特征降维算法, 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。

解:

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

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

step2. 确定 $\hat{\Sigma}$ 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 , 提取该样本r维新的特征向量: $\xi_r = A_r^T (x - \hat{\mu})$

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法的特征降维算法, 并满足累积方差解释比不低于0.9, 请确定新的特征空间特征维数r。请详细给出有关步骤和必要表达式。