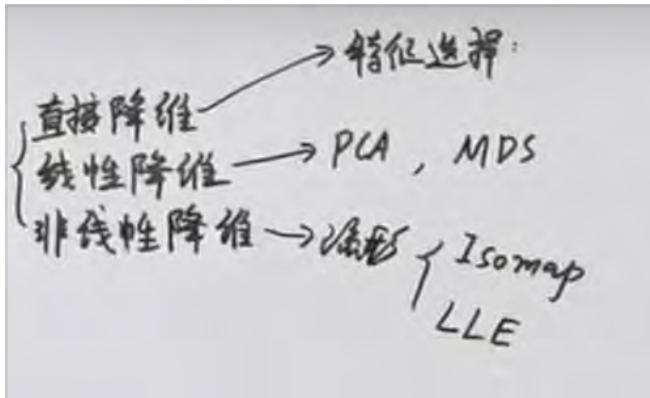


Week 13 笔记

PCA 主成分分析 (降维 dimension reduction)

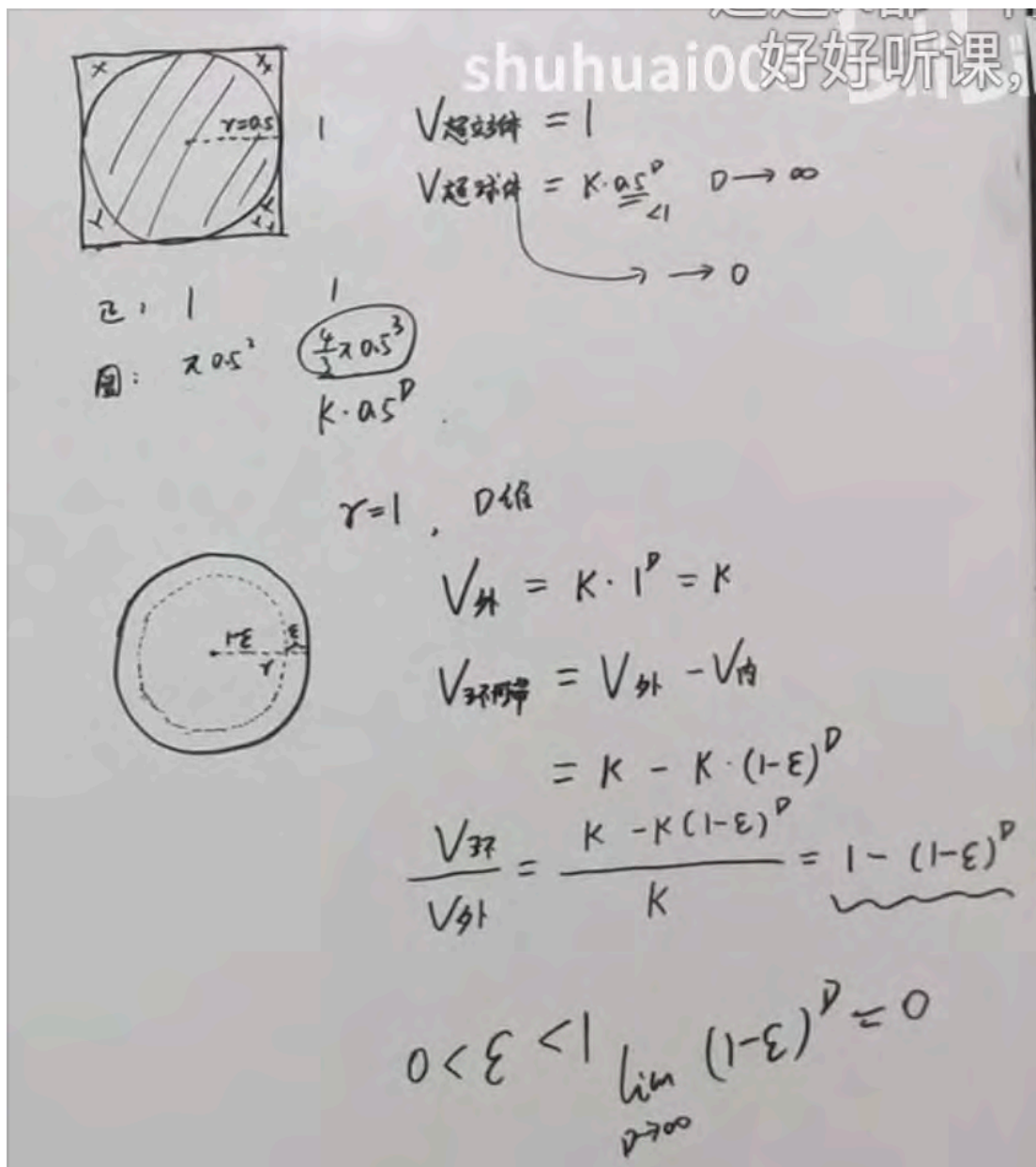
为什么要PCA? 防止过拟合 (提高模型泛化能力), 方法有:

1. 正则化 (将数据的 relative size 控制在同一数量级)
2. 增加样本量 (增加训练样本)
3. 降低维度 (防止维度灾难 dimension curse)
 1. 降维方法: 直接降维 (直接选择期望的特征 (特征降维))
 2. 线性降维: PCA MDS
 3. 非线性降维: 流形: Isomap LLE



维度灾难:

1. 几何角度



Lim $0 \rightarrow$ 正无穷 $(1 - (1 - \epsilon)^D) = 0$

在高维空间球体相对体积变小 会导致数据稀疏性 且不均匀

讲得太好了 救我狗命

(系列五) 降维2-样本均值&样本方差矩阵

预备知识

Data: $X = (x_1, x_2, \dots, x_N)^T_{N \times p} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{Np} \end{pmatrix}$
 $x_i \in \mathbb{R}^p, i=1,2,\dots,N$

Sample Mean: $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$

Sample Covariance: $S = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$

$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i = \frac{1}{N} \underbrace{(x_1, x_2, \dots, x_N)}_{X^T} \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}_{N \times 1} = \frac{1}{N} X^T \mathbf{1}_N$

$S = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$

$= \frac{1}{N} \underbrace{(x_1 - \bar{x} \quad x_2 - \bar{x} \quad \dots \quad x_N - \bar{x})}_{\substack{(x_1, x_2, \dots, x_N) - (\bar{x}, \bar{x}, \dots, \bar{x}) \\ X^T - \bar{x}(1 \ 1 \ \dots \ 1) \\ X^T - \bar{x} \mathbf{1}_N^T \\ X^T - \frac{1}{N} X^T \mathbf{1}_N \mathbf{1}_N^T}} \underbrace{\begin{pmatrix} (x_1 - \bar{x})^T \\ (x_2 - \bar{x})^T \\ \vdots \\ (x_N - \bar{x})^T \end{pmatrix}}_{\substack{\vdots \\ (X_N - \bar{x})^T}} \downarrow$

$= \frac{1}{N} X^T \underbrace{\left(\mathbf{I}_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T \right)}_{H_N \rightarrow \text{centering matrix}} \cdot \underbrace{\left(\mathbf{I}_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T \right)^T}_{\downarrow} X$

$= \frac{1}{N} X^T H \cdot H^T X$

$= \frac{1}{N} X^T H X$

给力

非常的 clean

Dimensionality Reduction

降维

$\mathbf{1}_N = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}_{N \times 1}$

$H = \left(\mathbf{I}_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T \right)$

$H^T = \left(\mathbf{I}_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T \right)^T = H$

$H^2 = H \cdot H = \left(\mathbf{I}_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T \right) \left(\mathbf{I}_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T \right)$

$= \mathbf{I}_N - \frac{2}{N} \mathbf{1}_N \mathbf{1}_N^T + \frac{1}{N^2} \mathbf{1}_N \mathbf{1}_N^T \mathbf{1}_N \mathbf{1}_N^T$

$= \mathbf{I}_N - \frac{1}{N} \begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{pmatrix}$

$= \mathbf{I}_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T$

$= H$

$H^N = H$

降维的矩阵表达

principal components analysis主成分分析

PCA

Data: $X = (x_1, x_2, \dots, x_N)^T = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{Np} \end{pmatrix}$
 $x_i \in \mathbb{R}^p, i=1,2,\dots,N$

Mean: $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i = \frac{1}{N} X^T \mathbf{1}_N$

Covariance: $S = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T = \frac{1}{N} X^T H X$

其中 $\mathbf{1}_N = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}_{N \times 1}, H = \mathbf{I}_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T, \bar{x} \in \mathbb{R}^p, S \in \mathbb{R}^{p \times p}$

Dimensionality Reduction

降维

一个中心：原始特征空间的重构
相关 \rightarrow 无关

两个基点：最大投影长度
最小重构距离

PCA降维方法 确保在降低维度后 在投影方向投影方差最大（即投影出来在一个维度上最稀疏）重构距离最小（在与投影维度正交的方向上最密集，实际上是对中心化

1. Max 投影方差:



对于一个 $n \times n$ 的框体内的取中位数来替代原本放置在中心位置，matlab代码实现：

% 输入:

%image: 原图

%m: 模板的大小 3*3 的模板, m=3

%输出:

%img: 中值滤波处理后的图像

%-----

```
n = m;
[ height, width ] = size(image);
x1 = double(image);
x2 = x1;
for i = 1: height-n+1
    for j = 1:width-n+1
        mb = x1( i:(i+n-1), j:(j+n-1) );%获取图像中n*n的矩阵
        mb = mb(:);%将mb变成向量化, 变成一个列向量
        mm = median(mb);%取中间值
        x2( i+(n-1)/2, j+(n-1)/2 ) = mm;

    end
end
```

```
img = uint8(x2);
```

```
end
```

2. 均值滤波:

均值滤波和中值滤波的区别在于取方框内的平均值来替代中位值, 所以图片会更加平滑 (锐度也更低)

3. 卡尔曼滤波:

卡尔曼滤波是用来解决现代控制问题的, 假设有状态空间:

$$\mathbf{x}_t = \mathbf{F}_t \mathbf{x}_{t-1} + \mathbf{B}_t \mathbf{u}_t + \mathbf{w}_t,$$

和观测空间:

$$\mathbf{z}_t = \mathbf{H}_t \mathbf{x}_t + \mathbf{v}_t,$$

我们根据上一个时刻的 \mathbf{x}_{t-1} 预测值和这一个时刻的观测值 \mathbf{z}_t , 假定这两个值都满足高斯分布, 且都是不完全准确的, 核心在于找得到卡尔曼增益 K , 来对这两个值进行求加权平均:

$$S = HP'H^T + R$$

$$K = P'H^TS^{-1}$$

如果在简单一维问题下，就简化为了：

Time Update (prediction)	Measurement Update (correction)
$\hat{x}_k^- = \hat{x}_{k-1}$ $P_k^- = P_{k-1}$	$K_k = \frac{P_k^-}{P_k^- + R}$ $\hat{x}_k = \hat{x}_k^- + K_k(z_k - \hat{x}_k^-)$ $P_k = (1 - K_k)P_k^-$

P_k 是噪声的分布

(参考文献 <https://ieeexplore.ieee.org/document/6279585>)

参考代码 <https://zhuanlan.zhihu.com/p/45238681>)

插值算法：主要用来在图片进行放大的时候怎么对于新的像素点来选择像素值，与采样互为逆运算

1. 最近邻插值：

```
def nearest(img, scale):
```

```
width, height, _ = img.shape
```

```
n_width = width*scale
```

```
n_height = height*scale
```

```
n_img = np.zeros((n_width, n_height, 3))
```

```
for k in range(3):
```

```
    for i in range(n_width):
```

```
        for j in range(n_height):
```

```
            #print(i, j, k)
```

```
            n_img[i,j, k] = img[round((i-1)/scale), round((j-1)/scale), k] #映射
```

```
    return Image.fromarray(np.uint8(n_img))
```

像素值映射公式为：

$$src_x = dst_x / scale$$

$$src_y = dst_y / scale$$

直接取最近的放缩后最近的一个像素点的值

2. 双线性插值

二维空间（坐标 (x,y)，图片的像素值为我们关心的 f(x)）中，我们需要计算出 P 点的像素值

取 P 点邻近的四个点 Q11,Q12,Q21,Q22，并假设在邻近范围内，点的像素值是呈线性变化的。

这时我们先在 x 方向上进行线性插值，计算出 R1,R2 的像素值，

再在 y 方向上进行单次线性插值，求出 P 的像素值。

这里的这些点都是位于原始图像上的，我们只需要找到一个映射公式，将放大图像上的点与 P 点对应即可。

当然这里的映射公式就是最近邻插值中的映射公式。与最近邻插值不同的是，双线性插值法并没有将映射点的像素值作为放大图像的像素值，而是将映射点周围的四个点的加权作为放大图像的像素值。

$$f(R1) = \frac{x2 - x}{x2 - x1} f(Q11) + \frac{x - x1}{x2 - x1} f(Q21)$$

$$f(R2) = \frac{x2 - x}{x2 - x1} f(Q12) + \frac{x - x1}{x2 - x1} f(Q22)$$

$$f(P) = \frac{y2 - y}{y2 - y1} f(R1) + \frac{y - y1}{y2 - y1} f(R2)$$

整合一下，就是：

$$f(P) = \frac{(x2 - x)(y2 - y)}{(x2 - x1)(y2 - y1)} f(Q11) + \frac{(x - x1)(y2 - y)}{(x2 - x1)(y2 - y1)} f(Q21) + \frac{(x2 - x)(y - y1)}{(x2 - x1)(y2 - y1)} f(Q12) + \frac{(x - x1)(y - y1)}{(x2 - x1)(y2 - y1)} f(Q22)$$

```
def double_linear(img, scale):
    width, height, _ = img.shape
    n_width = int(width*scale)
    n_height = int(height*scale)
    n_img = np.zeros((n_width, n_height, 3))
    for k in range(3):
        for i in range(n_width):
```

```

for j in range(n_height):
    src_x = i/scale
    src_y = j/scale
    src_x_0 = int(np.floor(src_x))
    src_y_0 = int(np.floor(src_y))
    src_x_1 = min(src_x_0 + 1, width - 1)
    src_y_1 = min(src_y_0 + 1, height - 1)
    #print(src_x, src_y, src_x_0, src_y_0, src_x_1, src_y_1)
    value0 = (src_x_1 - src_x)*img[src_x_0, src_y_0, k] + (src_x -
src_x_0)*img[src_x_1,src_y_0,k]
    value1 = (src_x_1 - src_x)*img[src_x_0, src_y_1, k] + (src_x -
src_x_0)*img[src_x_1,src_y_1,k]
    n_img[i, j, k] = int((src_y_1-src_y) * value0 + (src_y-src_y_0)*value1)
return Image.fromarray(np.uint8(n_img))

```

参考文献（最近邻插值、双线性插值与双三次插值 - ishihara的文章 - 知乎
<https://zhuanlan.zhihu.com/p/428523385>)

3. 拉格朗日插值

假定有 3 个点 $x_j, j=1,2,3$ ，定义 3 个二次函数 $f_i, i=1,2,3$ ，使得对应的二次函数在 x_i 处等于 1，在其他两个点处等于 0：

$$f_i(x_j), i = 1, 2, 3, j = 1, 2, 3$$

其中这三个二次函数需要满足：

$$f_i(x_j) = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases}$$

显然，我们构造的函数应该满足下列 2 式，第二个式子是更一般的场景：

$$f_1(x) = \frac{(x - x_2)(x - x_3)}{(x_1 - x_2)(x_1 - x_3)}$$

$$f_i(x) = \prod_{\substack{1 \leq j \leq 3 \\ j \neq i}} \frac{(x - x_j)}{(x_i - x_j)}$$

最终我们可以得到正确的 fitting 曲线：

$$f(x) = \sum_{i=1}^3 y_i f_i(x)$$