

Degree of Freedom

Definition

Vertical between ϵ and y

reference: 统计学“自由度”详解 - 知乎 (zhihu.com)reference: (9条消息) 【线性回归】线性代数角度解释最小二乘法_Jesszen的博客-CSDN博客_最小二乘法 线性代数

Definition

第四种定义：自由度是一个随机向量的自由维度数，也就是一个向量能被完整描述所需的最少单位向量数。

Fisher 给“Student (t分布的发现者)”解释自由度的时候是这么来解释的：将 n 个样本随机变量构造成一个随机向量，那么这个向量可以看成是 n 维空间的一个点，每有一个约束条件，则向量的自由维度减1。比如 n 个样本在求样本方差的时候要先计算样本均值，所以最后一个变量就和前面 $n-1$ 个相关，这样随机向量只能有 $n-1$ 个元素可以在 $n-1$ 维空间自由取值。在Fisher指出老皮尔逊的卡方检验方法自由度计算错误的时候他是用“约束”这个词来解释的，这个“约束”有点像上面的定义二。

如果所研究的问题能抽象为模型，使用第四种定义计算自由度会容易很多。 n 个随机样本看成 n 个随机变量，这 n 个样本随机变量可以表示为

$Iy = y$
 y 是 n 个样本的随机向量, I 是 n 维单位矩阵

，因为 I 的列空间是 n 维， $Iy=y$ ，所以 y 一定在 I 的列空间中， y 的 n 个元素可以在 n 维空间自由取值，推出 y 的自由度是 n 。

下面计算线性回归拟合值（回归方程部分）的自由度。

$y = \hat{y} + \epsilon, \hat{y} = X\hat{\beta}$
 X 是设计矩阵, $\hat{\beta}$ 是估计出来的回归系数向量

很显然，拟合值向量 \hat{y} 一定在设计矩阵 X 的列空间中， X 列空间维度是回归系数个数，假设有 p 个预测变量，加上截距则回归系数个数是 $p+1$ ，所以拟合值向量的自由维度就是 $p+1$ 。回归平方和

$SSR = \sum_{i=1}^n (y_i - \hat{y})^2$

是样本因变量的平均值，需要估计出来，失去一个自由度，所以SSR的自由度是拟合值的自由度减1，即 p 。因变量 y 的自由度是 n ，拟合值的自由度是 $p+1$ ，那么残差向量的自由度是 $n-(p+1)$

y

$=$

\hat{y}

$+$

ϵ

残差向量

ϵ

垂直于“设计矩阵列空间”(后面有说明)，也就是它在设计矩阵列空间的垂直补空间中， y 是 n 维，设计矩阵列空间是 $p+1$ 维， $p+1 \leq n$ ，则残差向量维度是 $n-(p+1)$ ，也就是残差向量的自由度为 $n-(p+1)$ ，接着可以推出残差平方和

\$\$

$SSE = \sum_{i=1}^n (\epsilon_i)^2$

\$\$

的自由度是 $n-(p+1)$ 。

Vertical between ϵ and y

1. 我们知道 $X\beta = y'$ ，拟合的 y' 是 m 维向量。
观测值 y 同样是 m 维向量。
2. 观测值 y 和拟合值 y' 这两个向量，因为必然存在的误差致使 $X\beta = y$ 【观测值】无解【 y 观测值不在列空间】。
那么我们拟合的 y' ，只能尽可能接近 y 【观测值】。
3. 我们假设 y 和 y' 不再一个平面，我们知道 y' ，是由 $COL(X)$ 线性组合表示的，假设图中的超平面是列空间 $col(X)$ ，那么 y' ，必定落在这个平面。【多维，想象超平面】

perpendicular.png

4. 那么 $X\beta = y'$ 和 y 观测值，怎么才能最接近？换个角度，就是距离最短？

在这个空间我们用欧氏距离度量，我们知道欧氏距离的涉及到平方和的根号，所以‘最小二乘法’，中的‘二乘’就是这个概念。

那么最小二乘法，最小又该怎么理解？联想到距离的概念

5. 向量 $\mathbf{e} = \mathbf{y} - \mathbf{y}' = \mathbf{y} - \mathbf{X}\beta$

$|\mathbf{e}|$ 自然就是距离【其实这个对应到RSS也就是残差平方和】，距离最小，必然就是正交投影。

\mathbf{e} 向量自然属于 R^m 维的子空间，称为也属于左零空间【 A 转置的零空间】，左零空间垂直于列空间 $\text{COL}(\mathbf{X})$ 。